

Connectionist learning of belief networks

Radford M. Neal

*Department of Computer Science, University of Toronto, 10 King's College Road,
Toronto, Ontario, Canada M5S 1A4*

Received January 1991
Revised November 1991

Abstract

Neal, R.M., Connectionist learning of belief networks, *Artificial Intelligence* 56 (1992) 71–113.

Connectionist learning procedures are presented for “sigmoid” and “noisy-OR” varieties of probabilistic belief networks. These networks have previously been seen primarily as a means of representing knowledge derived from experts. Here it is shown that the “Gibbs sampling” simulation procedure for such networks can support maximum-likelihood learning from empirical data through local gradient ascent. This learning procedure resembles that used for “Boltzmann machines”, and like it, allows the use of “hidden” variables to model correlations between visible variables. Due to the directed nature of the connections in a belief network, however, the “negative phase” of Boltzmann machine learning is unnecessary. Experimental results show that, as a result, learning in a sigmoid belief network can be faster than in a Boltzmann machine. These networks have other advantages over Boltzmann machines in pattern classification and decision making applications, are naturally applicable to unsupervised learning problems, and provide a link between work on connectionist learning and work on the representation of expert knowledge.

1. Introduction

The work reported here can be seen from two perspectives. From one point of view, it describes a connectionist network with capabilities comparable to those of the Boltzmann machine, but with better learning performance. From the other, it shows how belief networks can be learned from empirical data, as an alternative, or a supplement, to their specification by experts.

Correspondence to: R.M. Neal, Department of Computer Science, University of Toronto, 10 King's College Road, Toronto, Ontario, Canada M5S 1A4. E-mail: radford@cs.toronto.edu.

The original objective of the work was to find a network architecture that shared with Boltzmann machines [1,8] the capacity to learn arbitrary probability distributions over binary vectors, but that did not require the “negative phase” of Boltzmann machine learning. It was hypothesized that eliminating the negative phase would improve learning performance.

This goal was achieved by replacing the Boltzmann machine’s symmetric connections with directed, acyclic connections. In analogy with Boltzmann machines, the sigmoid function was used to compute the conditional probability of a unit being on from the weighted input from other units. The stochastic operation of such a network is somewhat more complex than for a Boltzmann machine, but is still possible using local communication. Maximum-likelihood, gradient-ascent learning can be done using locally available information.

These networks turn out to fall within the general class of “belief networks” studied by Pearl [15] and others as a means of representing probabilistic knowledge in expert systems. However, the specific network architectures considered by Pearl use a “noisy-OR” model for the probability of a unit being on, based on the states of units feeding into it. It is natural to ask whether a learning procedure can be developed for this model, as well as for that using the sigmoid function. A local learning rule was indeed found for a generalization of the noisy-OR model, though this time the gradient-ascent procedure must be constrained to avoid an invalid region of the weight space.

The representational power of the two types of belief network was investigated and compared to that of the Boltzmann machine. It turns out that each of these three networks can represent probability distributions over the full set of units that the other two networks cannot. With the help of “hidden” units, all these networks can represent arbitrary distributions over a set of “visible” units. Judicious placement of hidden and visible units can be used to constrain the representational capabilities of a belief network, in order to direct learning in a desired direction.

The presence of hidden units is an extreme case of “missing data”—data which is not always observed in the training cases. Learning procedures for belief networks that have been described previously have problems with missing data, but it is handled naturally by the method presented here (as well as by Boltzmann machines).

The learning capabilities of these networks were evaluated using a simple mixture distribution and an associated classification task. The sigmoid belief network was found to be capable of learning at a significantly higher rate than the Boltzmann machine. Additional experiments established that the sigmoid belief network’s advantage in learning speed is indeed due to the elimination of the negative phase. The noisy-OR belief network performed less well at the mixture modeling task, and in other cases showed a strong

tendency to get stuck at a local maximum. It did perform well at learning a distribution naturally expressed in the noisy-OR form, however.

This paper begins with reviews of Boltzmann machines and belief networks. I then define the sigmoid and noisy-OR varieties of belief network, derive gradient-ascent learning rules for them, and investigate their representational power. The experiment comparing the learning performance of these networks with Boltzmann machines is then described. Finally, I show how belief networks relate to other connectionist approaches to statistical modeling and to work on the representation of probabilistic knowledge in expert systems, and I discuss how these networks open up new possibilities for decision making, alternative learning procedures, and neural modeling.

2. A review of Boltzmann machines

The Boltzmann machine [1,8] is most naturally viewed as a device for modeling a probability distribution, from which conditional distributions for use in pattern completion and classification may be derived. In the limit as probabilities approach zero and one, deterministic input-output mappings can be represented as well. These capabilities would make the Boltzmann machine attractive in many applications, were it not that its learning procedure is generally seen as being painfully slow. Boltzmann machines have also been considered as a model of computation in the brain.

2.1. Definition of Boltzmann machines

A Boltzmann machine consists of some fixed number of two-valued units linked by symmetrical connections. In some formulations, the two possible values of a unit are 0 and 1; in other formulations the two values are -1 and $+1$. These alternate formulations are representationally equivalent, but the $-1/+1$ formulation is often found to have better learning performance.¹ Since, on the other hand, networks with 0/1-valued units are easier to understand, I will retain both formulations here.

The states of the units will be denoted by the vector \tilde{s} , with the state of unit i being s_i . This state vector will often be regarded as a realization of a corresponding random variable \tilde{S} . The weight on the connection between unit i and unit j will be denoted by w_{ij} . Since connections are symmetrical, $w_{ij} = w_{ji}$. Units do not connect to themselves. "Bias" weights, w_{i0} , from a fictitious unit 0 whose value is always 1 are also assumed to be present.

¹The benefit of a symmetric formulation for the related case of backpropagation networks is shown in [19].

In an analogy with thermodynamics, the “energy” of a network with state \tilde{s} is defined as follows:

$$E(\tilde{s}) = -\beta \sum_{j < i} s_i s_j w_{ij}, \quad (1)$$

where β is the constant 1 if units take on values of 0 and 1 or the constant $\frac{1}{2}$ if units take on values of -1 and $+1$. Intuitively, a state with low energy is more internally compatible than one with high energy.

The energy is used to define a “Boltzmann” probability distribution over states, in which low-energy states are more probable than high-energy states. Specifically,

$$P(\tilde{S} = \tilde{s}) = \exp(-E(\tilde{s})) / Z, \quad (2)$$

where Z is a normalization factor that makes the probabilities of all states sum to one:

$$Z = \sum_{\tilde{s}} \exp(-E(\tilde{s})). \quad (3)$$

Typically, some of the units in the network are “hidden”, and we are interested only in the marginal distribution of the other “visible” units. We then consider the state vector \tilde{s} to be split into the pair $\langle \tilde{h}, \tilde{v} \rangle$, and similarly the random variable \tilde{S} becomes $\langle \tilde{H}, \tilde{V} \rangle$. The distribution over the visible units is then

$$P(\tilde{V} = \tilde{v}) = \sum_{\tilde{h}} P(\tilde{S} = \langle \tilde{h}, \tilde{v} \rangle). \quad (4)$$

2.2. Gibbs sampling for Boltzmann machines

Since Z is the sum of an exponentially large number of terms, directly computing the probability of a given state vector is infeasible for networks of significant size. Even if this calculation could be performed efficiently, we would still need time exponential in the number of hidden units to calculate the marginal probability of a visible vector, or the probability distribution for a subset of visible units conditional on given values for the other visible units. These distributions can, however, be exhibited via a stochastic simulation procedure known as “Gibbs sampling”, a process which is fundamental to the operation of all the networks considered in this paper.²

The simulation starts with the network in an arbitrary state. Units are then repeatedly visited in turn, with a new value being selected on each visit

²The technique appears to have been first described in [11], in the form known as the “Metropolis algorithm”. General application of the method is discussed in [6].

according to the unit's probability distribution conditional on the values of all other units. For Boltzmann machines, this conditional distribution for unit i is as follows:

$$P(S_i = x \mid S_j = s_j : j \neq i) = \sigma\left(x^* \sum_{j \neq i} s_j w_{ij}\right). \quad (5)$$

The notation " $S_j = s_j : j \neq i$ " means the joint condition that $S_j = s_j$ for all j such that $j \neq i$. For $-1/+1$ -valued units, $x^* = x$, while for $0/1$ -valued units $x^* = 2x - 1$. The "sigmoid" function, $\sigma(t)$, is defined as $1/(1 + \exp(-t))$. Note that $\sigma(-t) = 1 - \sigma(t)$.

To produce a sample from the distribution over state vectors, the simulation is allowed to run for a length of time sufficient for it to settle to "equilibrium". A collection of state vectors taken at sufficiently widely separated times as the simulation continues to run will then form a sample from the distribution for \tilde{S} . Conditional distributions can be exhibited by clamping certain units to fixed values during the simulation and updating only the values of the remaining units. This allows the network to perform pattern completion and classification tasks.

Unfortunately, it is difficult to say how much time should be allowed for the simulation to reach equilibrium, or at what interval state vectors should subsequently be taken to form the sample. The technique of "simulated annealing" is often used to reach equilibrium faster. In this method, we make the probability distribution sampled from more uniform by raising the probability of each state to the power $1/T$ (and then renormalizing). T , the "temperature" parameter, is initially set high in order to make equilibrium easy to reach, and is then gradually reduced to 1, at which point we hope that the equilibrium distribution for the original probabilities will have been reached.

2.3. Learning in Boltzmann machines

The learning problem for Boltzmann machines is to adjust the weights so as to make the distribution over visible units match as closely as possible the distribution of some real-world attributes, as evidenced by a set of training cases.

Adopting the maximum-likelihood approach to such estimation, we attempt to maximize the log-likelihood given the training cases, defined as

$$L = \log \prod_{\tilde{v} \in \mathcal{T}} P(\tilde{V} = \tilde{v}) = \sum_{\tilde{v} \in \mathcal{T}} \log P(\tilde{V} = \tilde{v}), \quad (6)$$

where \mathcal{T} is the collection of training cases (which may contain repetitions).

The partial derivative of L with respect to a particular weight can be expressed as follows:

$$\frac{\partial L}{\partial w_{ij}} = \beta \sum_{\tilde{v} \in \mathcal{T}} \left(\sum_{\tilde{s}} P(\tilde{S} = \tilde{s} \mid \tilde{V} = \tilde{v}) s_i s_j - \sum_{\tilde{s}} P(\tilde{S} = \tilde{s}) s_i s_j \right). \quad (7)$$

The above formula provides the basis for a gradient-ascent learning procedure involving two parallel Gibbs sampling simulations for each training case. In the “positive phase” simulation, the visible units are clamped to the values they take in the training case, with the result that the simulation produces a sample consisting of some number of states from the conditional distribution of \tilde{S} given $\tilde{V} = \tilde{v}$. In the “negative phase” simulation, no units are clamped, producing an (equal size) sample from the unconditional distribution for \tilde{S} . For each state vector \tilde{s}^+ in the positive phase samples, the weight w_{ij} is incremented by a small amount proportional to $s_i^+ s_j^+$. For each state vector \tilde{s}^- in the negative phase samples, w_{ij} is decremented by an amount in the same proportion to $s_i^- s_j^-$. This procedure is repeated until convergence is reached.

2.4. Need for the negative phase

Intuitively, the need for a negative as well as a positive phase in Boltzmann machine learning arises from the presence of the normalizing factor, Z , in the expression for the probability of a state vector. Because of this, the direction of steepest descent in energy is not the same as that of steepest ascent in probability. The negative phase of the learning procedure is needed to account for this effect.

Looked at another way, the negative phase provides the mechanism by which the learning comes to a stop—once the correct distribution over visible units has been learned, this distribution is exhibited in the negative phase, just as it is forced in the positive phase. The positive phase increments and negative phase decrements then balance, on average, and the weights become stable.

The presence of the negative phase has several disadvantages:

- (1) It directly increases computation by a factor of more than two.
- (2) It may make the learning procedure more sensitive to statistical errors.
- (3) It may reduce any neurological plausibility the scheme possesses.

Note that since the negative phase simulations have more unclamped units, they take longer to run than the positive phase simulations. The presence of a negative phase may make it necessary to collect a larger sample of state

vectors from the simulations in order to reduce the variance in the estimate of the gradient of L , which will be the sum of the variances of the positive and the negative phase statistics. Taking the difference of the statistics from two phases may also exacerbate the ill-effects of not reaching equilibrium in the simulations.

On the other hand, the negative phase can be exploited to control how network resources are utilized. In particular, the network can be forced to learn a mapping between a group of visible input units and a group of visible output units while ignoring the distribution of the input units themselves. This is done by clamping the input units in the negative as well as the positive phase. It will turn out that in belief networks, where the negative phase has been eliminated, control over what the network learns can be exercised by other means.

3. A review of belief networks

Belief networks, also known as “Bayesian networks”, “causal networks”, “influence diagrams”, and “relevance diagrams”, are designed, like Boltzmann machines, to represent a probability distribution over a set of attributes. Study of these networks by Pearl [15] and others [13] has been motivated principally by the desire to represent knowledge obtained from human experts, however. Accordingly, hard-to-interpret parameters such as the weights in a Boltzmann machine have been avoided in favour of more intuitive representations of conditional probabilities.

3.1. Definition of belief networks

Sticking as closely as possible to the terminology of the previous sections, we can view the state of a belief network as a vector, \tilde{S} , with s_i being the state of unit i . In this paper, the units will always be two-valued. When belief networks are applied to expert system design, the units represent propositions concerning the problem situation that are meaningful to the expert.

The probability of a state vector is defined in terms of what I will call “forward condition probabilities”—the probability of a unit having a particular value conditional on the values of the units that precede it:

$$P(\tilde{S} = \tilde{s}) = \prod_i P(S_i = s_i \mid S_j = s_j : j < i). \quad (8)$$

The conditional probabilities above are assumed to have been given by an expert. Typically, only a subset of the units preceding unit i will be “connected” to it, and only these will be relevant in specifying its forward

conditional probabilities. Note that the ordering of units in the state vector is crucial, since it determines which conditional probabilities must be specified.

3.2. Gibbs sampling for belief networks

In contrast with Boltzmann machines, computing the probability of a particular state vector for a belief network is straightforward. One can also easily generate a sample from the distribution for \tilde{S} . However, making predictions by computing conditional probabilities or sampling from conditional distributions are in general difficult problems. Various methods for computing exact conditional probabilities in belief networks have been proposed [9,15,17], but all are either restricted to special forms of network or have exponential time complexity in the worst case.

It appears that the only plausible method of sampling from conditional distributions in belief networks with high connectivity is Gibbs sampling, introduced in this context by Pearl [14,15]. As with Boltzmann machines, a step in the simulation requires selecting a new value for unit i from its distribution conditional on the values of the other units. For a belief network, this distribution is given by the proportionality

$$\begin{aligned} &P(S_i = x \mid S_j = s_j : j \neq i) \\ &\propto P(S_i = x \mid S_j = s_j : j < i) \\ &\quad \cdot \prod_{j>i} P(S_j = s_j \mid S_i = x, S_k = s_k : k < j, k \neq i). \end{aligned} \quad (9)$$

For this procedure to be guaranteed to work (in the limit as the number of simulation passes grows), the forward conditional probabilities should be nonzero. The time to reach equilibrium in the simulation can be reduced by using simulated annealing, as described for Boltzmann machines.

A “short-cut” simulation method is possible when the units whose values are known happen to be the first ones in the state vector. In this case, rather than employ the full Gibbs sampling procedure, we can simply select new values for each unclamped unit in a single forward pass, using the forward conditional probabilities. The selection for each unit depends only on the values for preceding units, and the values in the previous state vector have no effect on the result. Accordingly, no settling to equilibrium is required, and the state vectors obtained in successive passes are all independent. This short-cut can be exploited when belief networks are used for pattern classification or for other tasks that have the form of an input–output mapping.

3.3. The noisy-OR model of conditional probabilities

So far, forward conditional probabilities have been assumed to be given explicitly. In fact, this will generally not be feasible, since explicitly specifying the conditional distribution for S_i given the values of the preceding units requires 2^{i-1} parameters. Even if some of the preceding units are not connected to unit i , more compact specifications will generally be necessary.

One method, termed the “noisy-OR” model [7,15], views the units as 0/1-valued OR gates with the preceding units as inputs. An input of 1 does not invariably force a unit to take on the value 1, however. Rather, there is a certain probability, q_{ij} , that even though unit j has the value 1, it will fail to force a unit i that it feeds into to go to 1. Under this model, the forward conditional probabilities can be expressed in terms of the q_{ij} as follows:

$$P(S_i = 1 \mid S_j = s_j : j < i) = 1 - \prod_{j < i, s_j = 1} q_{ij}. \quad (10)$$

Once again, a fictitious unit 0 whose value is always 1 has been assumed, along with associated parameters q_{i0} . Note that if q_{ij} is one, unit j is effectively not connected to unit i .

4. Two varieties of belief network

We are now in a position to describe the two types of belief network that are investigated in this paper. The first, “sigmoid”, variety was designed in analogy with Boltzmann machines.³ When the connection with belief networks was realized, the second variety was developed as a generalization of the “noisy-OR” model for specifying conditional probabilities. In contrast with the use of belief networks in expert systems, the units in these networks will not necessarily be seen as representing propositions that would be meaningful in human terms. As in other connectionist systems, the units in these networks may interact in ways that are useful in solving a problem without mimicking its usual symbolic structure.

4.1. Definition of sigmoid belief networks

Two formulations of sigmoid belief networks will be considered. In one, units take on the values 0 and 1, in the other, they take on the values -1 and $+1$. Directed forward connections connect the units. The weight on the connection from unit j to unit i will be denoted by w_{ij} . A bias unit, 0, set permanently to 1 is assumed to exist, with associated weights, w_{i0} . The

³Belief networks of this type have also been discussed in [18]. They can be seen as generalizations of the “logistic regression” model well-known in statistics.

forward conditional probabilities for sigmoid belief networks can then be defined as follows:

$$P(S_i = s_i \mid S_j = s_j : j < i) = \sigma\left(s_i^* \sum_{j < i} s_j w_{ij}\right). \quad (11)$$

Here again, for $-1/+1$ -valued units, $s_i^* = s_i$, while for $0/1$ -valued units, $s_i^* = 2s_i - 1$. Note the analogy with equation (5) for the Boltzmann machine. The above gives the distribution for a unit's value conditional only on the units *preceding* it, however, not on all other units.

One can easily verify that a network of $0/1$ -valued units with forward conditional probabilities defined as above can be converted to an equivalent network of $-1/+1$ -valued units with weights w'_{ij} by the transformation:

$$w'_{i0} = w_{i0} + \sum_{0 < j < i} w_{ij}/2, \quad (12)$$

$$w'_{ij} = w_{ij}/2 \quad \text{for } 0 < j < i. \quad (13)$$

This transformation is easily inverted. The two formulations thus have equal representational power. A similar equivalence applies to Boltzmann machines.

The probability of a state vector, \tilde{s} , is defined in terms of the forward conditional probabilities:

$$P(\tilde{S} = \tilde{s}) = \prod_i P(S_i = s_i \mid S_j = s_j : j < i) \quad (14)$$

$$= \prod_i \sigma\left(s_i^* \sum_{j < i} s_j w_{ij}\right). \quad (15)$$

As with Boltzmann machines, we are often interested in the marginal distribution over a subset of "visible" units, given by

$$P(\tilde{V} = \tilde{v}) = \sum_{\tilde{h}} P(\tilde{S} = \langle \tilde{h}, \tilde{v} \rangle), \quad (16)$$

where \tilde{S} has been split into $\langle \tilde{H}, \tilde{V} \rangle$. We are also interested in conditional distributions involving subsets of visible units, as these allow one to perform tasks such as pattern completion and classification.

To exhibit these marginal and conditional distributions via Gibbs sampling, we must repeatedly select a new value for each unit from its distribution conditional on the rest of the network. This distribution is given by the proportionality

$$P(S_i = x \mid S_j = s_j : j \neq i) \propto \sigma\left(x^* \sum_{j < i} s_j w_{ij}\right) \prod_{j > i} \sigma\left(s_j^* \left(x w_{ji} + \sum_{k < j, k \neq i} s_k w_{jk}\right)\right). \quad (17)$$

To select a new value from the above distribution, unit i must have available both its own total input: $\sum_{j<i} s_j w_{ij}$, and the input to each unit, j , that it feeds into, exclusive of its own contribution: $\sum_{k<j, k \neq i} s_k w_{jk}$. The procedure is thus somewhat more complex than that for a Boltzmann machine (see equation (5)), but the information required can still be made available through local network communication, provided data can pass both ways along the directed connections.

4.2. Noisy-OR belief networks

In the “noisy-OR” form of belief network described above, the probabilities, q_{ij} , that an input of 1 from unit j into unit i will be ineffective in forcing unit i to 1 can be replaced with weights defined by $w_{ij} = -\log q_{ij}$. The forward conditional probabilities of equation (10) can then be written as follows:

$$P(S_i = 1 \mid S_j = s_j : j < i) = 1 - \exp\left(-\sum_{j<i} s_j w_{ij}\right). \quad (18)$$

Here, units take on values of 0 and 1, and a unit 0 set permanently to 1 exists.

In the above formulation, all weights are nonnegative. However, the conditional probability specification will be valid even when some weights are negative, provided that the weighted input to a unit cannot be negative, no matter what states the preceding units have. For a network with 0/1-valued units, this is equivalent to the constraint that, for all i ,

$$w_{i0} + \sum_{j<i, w_{ij}<0} w_{ij} \geq 0. \quad (19)$$

With this generalization, units can behave not only as OR gates, but also as OR gates with some or all inputs negated. For example, if unit 3 has input weights of $w_{30} = +20$, $w_{31} = -10$, and $w_{32} = -10$, it will behave as a slightly noisy OR gate with negated inputs from units 1 and 2 (i.e. as a NAND gate).

Noisy-OR belief networks can also be formulated with $-1/+1$ -valued units. Forward conditional probabilities are defined as above, with the constraint that to be valid, the weights must satisfy the following, for all i :

$$w_{i0} - \sum_{0<j<i} |w_{ij}| \geq 0. \quad (20)$$

The same equivalence between 0/1 and $-1/+1$ formulations that applied to sigmoid networks applies to noisy-OR networks as well.

Conditional probability distributions for noisy-OR belief networks can be exhibited using Gibbs sampling in a process entirely analogous to that described above for sigmoid networks.

5. Learning belief networks from empirical data

The particular formulations of belief networks that have been described are meant to be learnable from empirical data, rather than being constructed from expert knowledge. Except for the lack of a negative phase, learning is similar to that in a Boltzmann machine.

These learning procedures are all based on the widely used method of maximum-likelihood estimation. One should realize that this method is prone to “overfitting” the data when the amount of training data is small in relation to the number of free network parameters, with the result that the network generalizes poorly to future cases. (The use of Bayesian and cross-validation methods to avoid this is briefly discussed later.) Also, all these procedures use gradient-ascent to try to find a set of weights maximizing the likelihood. This may lead to the learning getting stuck at a point that is a local but not global maximum of the likelihood.

5.1. Learning in sigmoid networks

In the learning scenario assumed here, we have a collection, \mathcal{T} , of training cases drawn from the distribution of interest. Each training case consists of the values for certain attributes, assumed here to be two-valued. Exact repetitions are possible, indeed expected, in proportion to how common a particular combination of attributes is.

In order to model the distribution from which the training sample was drawn, we first decide on some size for a state vector, \tilde{S} , for our network, and then select some subset, \tilde{V} , of units in the state vector to represent the attributes in the training cases. The remaining, “hidden”, units constitute the set \tilde{H} . Note that since the ordering of units in the state vector is significant for belief networks, different selections for the subset of visible units may give different results. This is discussed further in Section 6.4.

Next, we must find values for the network weights that maximize the likelihood given the training cases, though to avoid overfitting or to reduce computation we might decide to fix certain weights at zero based on *a priori* knowledge. Other weights will be set to zero (or to small random values if we wish to break symmetry faster) and then adjusted by gradient-ascent so as to maximize the log-likelihood:

$$L = \log \prod_{\tilde{v} \in \mathcal{T}} P(\tilde{V} = \tilde{v}) = \sum_{\tilde{v} \in \mathcal{T}} \log P(\tilde{V} = \tilde{v}). \quad (21)$$

For a sigmoid belief network, the partial derivatives of the log-likelihood with respect to the weights may be found as follows:

$$\frac{\partial L}{\partial w_{ij}} = \sum_{\tilde{v} \in \mathcal{T}} \frac{1}{P(\tilde{V} = \tilde{v})} \frac{\partial P(\tilde{V} = \tilde{v})}{\partial w_{ij}} \quad (22)$$

$$= \sum_{\tilde{v} \in \mathcal{T}} \frac{1}{P(\tilde{V} = \tilde{v})} \sum_{\tilde{h}} \frac{\partial P(\tilde{S} = \langle \tilde{h}, \tilde{v} \rangle)}{\partial w_{ij}} \quad (23)$$

$$= \sum_{\tilde{v} \in \mathcal{T}} \sum_{\tilde{h}} P(\tilde{S} = \langle \tilde{h}, \tilde{v} \rangle \mid \tilde{V} = \tilde{v}) \cdot \frac{1}{P(\tilde{S} = \langle \tilde{h}, \tilde{v} \rangle)} \frac{\partial P(\tilde{S} = \langle \tilde{h}, \tilde{v} \rangle)}{\partial w_{ij}} \quad (24)$$

$$= \sum_{\tilde{v} \in \mathcal{T}} \sum_{\tilde{s}} P(\tilde{S} = \tilde{s} \mid \tilde{V} = \tilde{v}) \frac{1}{P(\tilde{S} = \tilde{s})} \frac{\partial P(\tilde{S} = \tilde{s})}{\partial w_{ij}} \quad (25)$$

$$= \sum_{\tilde{v} \in \mathcal{T}} \sum_{\tilde{s}} P(\tilde{S} = \tilde{s} \mid \tilde{V} = \tilde{v}) \cdot \frac{1}{\sigma(s_i^* \sum_{k < i} s_k w_{ik})} \frac{\partial \sigma(s_i^* \sum_{k < i} s_k w_{ik})}{\partial w_{ij}} \quad (26)$$

$$= \sum_{\tilde{v} \in \mathcal{T}} \sum_{\tilde{s}} P(\tilde{S} = \tilde{s} \mid \tilde{V} = \tilde{v}) s_i^* s_j \sigma(-s_i^* \sum_{k < i} s_k w_{ik}). \quad (27)$$

The last step uses the fact that $\sigma'(t) = \sigma(t)\sigma(-t)$.

These partial derivatives can be evaluated by running a separate Gibbs sampling simulation of the network for each training case, clamping the visible units to the values they take in that training case and observing the state vectors that arise as a result. If the simulation is run “long enough”, these observations will form a sample from the conditional distribution for \tilde{S} given the values in the current training case. Incrementing each weight, w_{ij} , by a small amount proportional to the average value of $s_i^* s_j \sigma(-s_i^* \sum_{k < i} s_k w_{ik})$ over the combined samples for all training cases will then move the weights along the gradient toward a local maximum of the likelihood. Various detailed implementations of this procedure are possible, as is discussed in Section 7.2.

Intuitively, only a single phase is needed for learning in a sigmoid belief network because normalization of the probability distribution over state vectors is accomplished locally at each unit via the sigmoid function, rather than globally via the hard-to-compute normalization constant, Z . The role of the Boltzmann machine’s negative phase in stopping learning once the distribution has been correctly modeled is taken over by the factor $\sigma(-s_i^* \sum_{k < i} s_k w_{ik})$ used to weight the learning increments. In the limit-

ing case where unit i is learning a deterministic function of the preceding units, for example, this factor is the probability that unit i would be set to the wrong value in an unclamped network. As the correct function is approached, this factor becomes zero, and the learning stops.

5.2. Learning in noisy-OR belief networks

Learning in noisy-OR belief networks is analogous to that in sigmoid belief networks, with the added complication that the gradient-ascent procedure must be constrained to the region of the weight space that produces valid probabilities for state vectors.

The partial derivatives of the log-likelihood with respect to the weights in a noisy-OR belief network can be expressed as follows, starting from equation (25):

$$\begin{aligned} \frac{\partial L}{\partial w_{ij}} &= \sum_{\tilde{v} \in \mathcal{T}} \sum_{\tilde{s}} P(\tilde{S} = \tilde{s} \mid \tilde{V} = \tilde{v}) \frac{1}{P(\tilde{S} = \tilde{s})} \frac{\partial P(\tilde{S} = \tilde{s})}{\partial w_{ij}} \quad (28) \\ &= \sum_{\tilde{v} \in \mathcal{T}} \sum_{\tilde{s}} P(\tilde{S} = \tilde{s} \mid \tilde{V} = \tilde{v}) \cdot \left\{ \begin{array}{l} -s_j + s_j / \left(1 - \exp\left(-\sum_{k < i} s_k w_{ik}\right)\right), \\ \quad \text{if } s_i = 1, \\ -s_j, \quad \text{if } s_i \neq 1. \end{array} \right\}. \quad (29) \end{aligned}$$

This formula is valid both for networks with 0/1-valued units and for those with $-1/+1$ -valued units.

These derivatives are computed via Gibbs sampling and used to perform gradient-ascent learning as described above for sigmoid belief networks. For noisy-OR belief networks, however, we must also ensure that the weights always define a valid probability distribution. In fact, in order for the simulations to reach equilibrium in a reasonable amount of time, it is desirable to further constrain the weights so that the conditional probability of unit i being 1 given the values of the preceding units is at least some minimum. For a noisy-OR network with 0/1-valued units, this will be so provided that

$$w_{i0} + \sum_{j < i, w_{ij} < 0} w_{ij} \geq \eta, \quad (30)$$

where η is some small positive constant. This can be ensured by applying the procedure in Fig. 1 to the weights for unit i after each movement along the gradient. One can show⁴ that this procedure moves the weights to the set of

⁴See Claim A.1 in Appendix A.

```

Loop:
   $C \leftarrow \{0\} \cup \{j : 0 < j < i \ \& \ w_{ij} < 0\}$ 
   $t \leftarrow \sum_{j \in C} w_{ij}$ 
  If  $t \geq \eta$  then exit loop
   $d \leftarrow (\eta - t)/|C|$ 
  For each  $j \in C - \{0\}$ : if  $|w_{ij}| < d$  then  $d \leftarrow |w_{ij}|$ 
  For each  $j \in C$ :  $w_{ij} \leftarrow w_{ij} + d$ 
End loop

```

Fig. 1. Procedure to move the weights into unit i to the valid region.

valid values that is closest in Euclidean distance to the previous set. Since the valid region of the weight space is convex, it follows that if one starts with a valid initial set of weights, moves in the direction of the gradient, and then applies the above procedure, the resulting total movement will have a positive projection in the direction of the gradient whenever this is possible.

An analogous constraint procedure exists for noisy-OR belief networks with $-1/+1$ -valued units.

6. Representational power of belief networks

In this section, I will investigate how powerful the various forms of belief network are at representing probability distributions. I will first consider sigmoid belief networks, and then discuss the noisy-OR variety. I also describe how the placement of hidden units in a belief network can be used to constrain its representational power in order to control learning.

Recall that, as mentioned earlier, there is no difference in the representational power of networks with $0/1$ -valued units and those with $-1/+1$ -valued units.

6.1. Representing distributions over the full set of units

Consider first the relative capacity of sigmoid belief networks and Boltzmann machines to represent probability distributions over the full state vector, \tilde{S} .

One can show that any distribution over one or two units can be approximated arbitrarily closely by either a Boltzmann machine or a sigmoid belief network, while for networks of three units, the restricted set of possi-

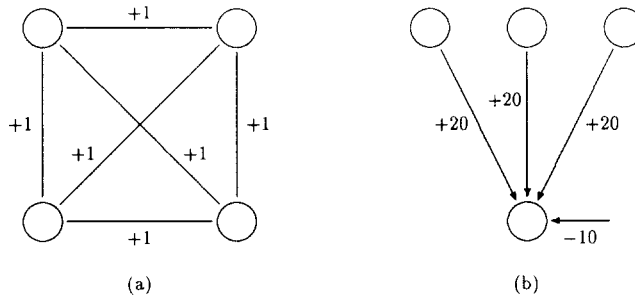


Fig. 2. An untranslatable Boltzmann machine (a) and sigmoid belief network (b).

ble probability distributions turns out to be the same for the two types of network.⁵

With four or more units, however, both Boltzmann machines and sigmoid belief networks can represent probability distributions that the other cannot. This is illustrated in Fig. 2, which shows a Boltzmann machine that cannot be translated into a sigmoid belief network, and a sigmoid belief network that cannot be translated into a Boltzmann machine.⁶ In both cases, 0/1-valued units are used, and absent connection weights are assumed to be zero. An unattached connection is from the bias unit.

Intuitively, the belief network cannot express the symmetric compatibility relations in the Boltzmann machine of Fig. 2(a), while the Boltzmann machine is not capable of equalizing the probabilities for all patterns over the top three units in the belief network of Fig. 2(b).

6.2. Representing mixture distributions

Suppose we are interested in the probability distribution over a vector of “visible” units, \tilde{V} . As seen above, in general, not all such distributions will be representable in a net consisting of these visible units alone. This problem can be overcome by including additional “hidden” units in a network.

In particular, sigmoid belief networks and Boltzmann machines can use hidden units to represent visible distributions that are expressed as “mixtures” of several other distributions. Such a mixture distribution can be written as follows:

$$P(\tilde{V} = \tilde{v}) = \sum_m P(\tilde{V} = \tilde{v} \mid M = m) P(M = m). \quad (31)$$

The hidden variable M identifies a “component” of the mixture. Each component produces its own distribution for \tilde{V} ; these are then combined in

⁵See Claims A.2 and A.3 in Appendix A.

⁶See Claims A.4 and A.5 in Appendix A.

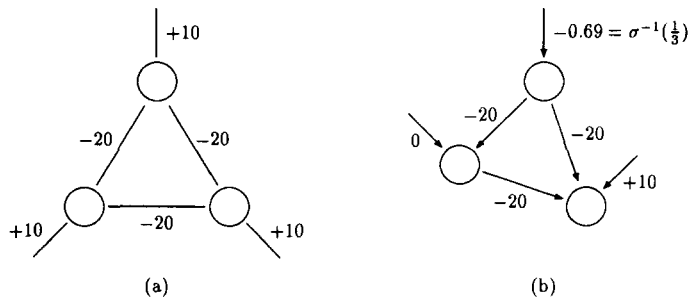


Fig. 3. 1-in-3 clusters in a Boltzmann machine (a) and a sigmoid belief network (b).

the proportions $P(M)$. In this paper, the component distributions will be such that the visible variables are independent—i.e.

$$P(\tilde{V} = \tilde{v} \mid M = m) = \prod_i P(V_i = v_i \mid M = m).$$

Such mixture distributions are commonly encountered and much studied [21].

To represent a mixture distribution in a network, we need first to represent the mixture variable, M . For a mixture of n components, one way to do this is via a cluster of n units, exactly one of which is on at any time. Figure 3 shows how a three-unit cluster of this sort can be constructed for both a Boltzmann machine and a sigmoid belief network (with 0/1-valued units). In both cases, the three state vectors with exactly one unit on have nearly equal probabilities, and all other state vectors have very small probability. The constructions generalize to clusters of any size, and to clusters in which the possible state vectors have unequal probabilities.

Using a 1-in- n cluster, we can implement a mixture in which the component distributions assign independent probabilities to the various visible units. For belief networks, all that is required is to connect each cluster unit to the visible units using weights that produce the required conditional probabilities. For Boltzmann machines, after making these connections to visible units, one must also adjust the bias weights to the cluster units in order to re-create the correct mixture proportions.

Note that any distribution over k visible units can be represented as a mixture of 2^k component distributions, each of which generates but a single vector. It follows that sigmoid belief networks and Boltzmann machines can approximate any distribution over k visible units arbitrarily closely, provided one is prepared to employ 2^k hidden units.

6.3. Power of noisy-OR belief networks

Unlike sigmoid belief networks, the ability of a noisy-OR network to represent a distribution is sensitive to negation of the unit values. For example, there is no way to make a noisy-OR unit behave as an AND gate, but one can make one behave as a NAND gate. This sensitivity is of significance only for visible units, since the output of a hidden unit can always be implicitly negated by negating the weights on all its outgoing connections.

In view of the above, there are certainly distributions over \tilde{S} that both a Boltzmann machine and a sigmoid belief network can implement but which a noisy-OR belief network cannot. Conversely, one can show⁷ that there are distributions over a three-unit noisy-OR network that cannot be duplicated by either a Boltzmann machine or a sigmoid belief network with only three units.

A 1-in- n cluster similar to that of Fig. 3(b) but with the unit values negated can be constructed from noisy-OR units. Such a cluster can be used to represent a mixture distribution over visible units using a noisy-OR belief network, just as with sigmoid belief networks and Boltzmann machines.

6.4. Manipulating the representational power of belief networks

When training a Boltzmann machine, one can control what the network learns by clamping certain units in the negative as well as the positive phase. This is commonly done when only the mapping from a set of “input” attributes to a set of “output” attributes is of interest—i.e. when we wish to do “supervised” rather than “unsupervised” learning. Clamping the input units in both phases forces the hidden units to model the conditional distribution of the output given the input, rather than the distribution of the input itself.

The same technique could be used with belief networks, but this would naturally require re-introduction of a negative learning phase, the elimination of which was the original motivation for this work. Fortunately, one can achieve similar control via judicious placement of input, output, and hidden units within a belief network, in such a way as to limit its representational power to that which one wishes it to learn.

Four network architectures that illustrate the control possible are shown in Fig. 4, using a medical diagnosis problem as an example. In all cases, a set of visible “symptom” units is used to represent various attributes of a patient, and a set of visible “disease” units is used to encode a diagnosis. There are assumed to be no connections among the units within each visible set.

⁷See Claim A.6 in Appendix A.

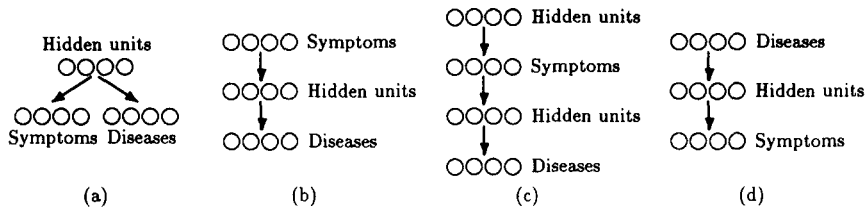


Fig. 4. Four network architectures for a medical diagnosis problem.

Training data is assumed to be available giving the true sets of symptoms and diseases for a sample of patients.

The network in Fig. 4(a) is designed for unsupervised learning—for modeling the data without any particular task in mind. The hidden units feed into both sets of visible units. As a result of training, these units may come to model correlations among symptoms, among diseases, or between symptoms and diseases. If the network succeeds in modeling the total distribution perfectly, it will be capable of performing any sort of pattern completion task. For example, one could clamp a set of symptoms and then observe the most likely diseases as the Gibbs sampling procedure is run, or, conversely, one could clamp a set of diseases and observe the most likely symptoms. However, if the number of hidden units is insufficient to model the total distribution, the network will end up modeling whichever correlations are strongest, and these might not be the ones that are most important for diagnosis.

The network in Fig. 4(b) is designed for supervised learning aimed at the diagnostic task. The hidden units are placed between the symptom units and the disease units. This forces the hidden units to learn to model the conditional distribution of the diseases given the symptoms. One could then clamp a set of symptoms and observe the most likely diseases. In fact, this can be done using the short-cut simulation procedure described in Section 3.2, since the clamped symptom units precede all the unclamped units (the full simulation procedure is still required during learning). The converse operation of clamping a set of diseases and observing likely symptoms no longer works well, however, since there are no hidden units in a position to model correlations among symptoms.

The network in Fig. 4(c) adds a set of hidden units prior to the symptom units in order to capture such correlations. This network has capabilities comparable to those of network (a), with the difference that the number of hidden units devoted to modeling each type of correlation is under the control of the network designer. Network (c) might be appropriate for a diagnosis application in which knowledge of correlations among symptoms is sometimes needed in order to fill in missing symptom values.

None of the above networks express the usual causal view that diseases

cause symptoms (not the other way around), and that the presence of one disease only weakly affects whether other diseases are also present. This illustrates the fact that the arrows in a belief network are a device for expressing probabilities, and need not correspond to real influences. In some circumstances we may wish to learn a network that does correspond to our causal view—we may feel that such a network would be easier to interpret, for example. Figure 4(d) shows how such a network can be set up for the disease/symptom example. Hidden units could be added prior to the disease units if we wish to model the weak correlations between diseases. This architecture would also be appropriate if the training data gives only the patients' symptoms, and we wish to *discover* a set of diseases that explains these symptoms in an unsupervised fashion.

7. Empirical comparison with Boltzmann machines

In this section, I will describe an experiment in which the learning procedures for belief networks and for Boltzmann machines were compared on the task of learning a simple mixture distribution and classifying items derived from it. Further details on this experiment may be found in [12].

7.1. Objectives of the experiment

This experiment is intended to answer the following questions:

- (1) Are the learning procedures for belief networks capable in practice of learning an approximation to a nontrivial distribution, based on a set of training cases?
- (2) If so, how does the speed of learning in sigmoid belief networks compare to the speed of learning in a Boltzmann machine?
- (3) Can differences in learning speed between sigmoid belief networks and the Boltzmann machine be attributed to the lack of a negative phase in the learning procedure for the belief networks?
- (4) Are there differences in the learning performances of networks with 0/1-valued units and those with $-1/+1$ -valued units?
- (5) How does learning in noisy-OR belief networks compare to learning in sigmoid belief networks?
- (6) How well do the solutions learned by the various networks on the basis of training data generalize to the true distribution?

Regarding points (2) and (3), the expectation is that the negative phase adds additional noise to the estimation of the gradient, and that this noise is detrimental to the learning process in Boltzmann machines. The magnitude of this effect is hard to judge, however. The added noise could even be beneficial, if it allows the network to escape local maxima during learning.

7.2. The learning procedure used

Numerous variations of the Boltzmann machine learning procedure have been tried [5], each of which requires fixing a number of parameters, such as the learning rate, and the temperatures in an annealing schedule. This presents a problem in comparing learning in Boltzmann machines to learning in belief networks—a valid comparison would require searching for the optimal parameter settings for each type of network, which would be a rather large undertaking.

The approach I have adopted is to train both types of network using a simple method that has only one adjustable parameter—the learning rate, ϵ . A complete picture of the performance of each type of network for various values of ϵ can be obtained with a reasonable number of runs, and the relative performance of the different networks with their best ϵ can then be compared.

The procedure used can be characterized as follows:

- (1) Learning was done in “batch” mode—i.e. each change to the weights was made on the basis of the entire set of training cases.
- (2) Each training case was clamped into a separate copy of the network, where a separate Gibbs sampling simulation was run.⁸ For Boltzmann machines, there was also an unclamped negative phase copy of the network associated with each training case.
- (3) No annealing was done.
- (4) The state of each copy of the network was retained after each change to the weights, on the assumption that if the weight changes are “small”, these existing simulation states will be close to equilibrium, and be good starting points for the next pass.
- (5) Changes to the weights were made after each simulation pass, based on the sample consisting of the current state vectors from the simulations for all training cases (plus the state vectors from the negative phase simulations, for Boltzmann machines).
- (6) Weight changes were scaled by a learning rate parameter, ϵ .
- (7) Weights were set to zero initially. Symmetry was broken by the stochastic nature of the simulation.

In detail, the weights in the Boltzmann machines were changed by

$$\Delta w_{ij} = \frac{\beta\epsilon}{N} \left(\sum_{\tilde{s}^+ \in \mathcal{T}^+} s_i^+ s_j^+ - \sum_{\tilde{s}^- \in \mathcal{T}^-} s_i^- s_j^- \right). \quad (32)$$

⁸This aspect of the learning procedure appears to be advantageous from an engineering point of view, but is quite implausible in a biological context.

Here, \mathcal{T}^+ is the set of current state vectors from the positive phase simulations (one per training case), and \mathcal{T}^- is the set of state vectors from the corresponding negative phase simulations. N is the number of training cases.

Similarly, the weights in the sigmoid belief networks were changed by

$$\Delta w_{ij} = \frac{\varepsilon}{N} \sum_{\tilde{s} \in \mathcal{T}} s_i^* s_j \sigma \left(-s_i^* \sum_{k < i} s_k w_{ik} \right), \quad (33)$$

and those in the noisy-OR belief networks by

$$\Delta w_{ij} = \frac{\varepsilon}{N} \sum_{\tilde{s} \in \mathcal{T}} \begin{cases} -s_j + s_j / \left(1 - \exp \left(- \sum_{k < i} s_k w_{ik} \right) \right), \\ \quad \text{if } s_i = 1, \\ -s_j, \quad \text{if } s_i \neq 1. \end{cases} \quad (34)$$

Weight changes in noisy-OR networks were limited to a magnitude of no more than 1 to avoid the possibility of huge weight changes resulting from the division above. After all changes were made, the constraint procedure of Fig. 1 with $\eta = 2^{-7}$ was applied.

The lack of annealing in this procedure is unconventional, as is the changing of weights based on a single state vector from each training case. The rationale behind these choices is that as ε approaches zero, the simulations will necessarily approach equilibrium, as they will run for many passes with the weights essentially unchanged. Furthermore, the cumulative effect of many changes with a small ε that are based on a single state vector from each training case will be equivalent to a single change with a larger ε that is based on a larger sample. As ε approaches zero, the learning procedure used will thus “do the right thing”.

Whether this procedure is better or worse than previous methods is not important for this experiment, however, provided only that any differences in learning performance between the various networks seen using this procedure will show up in some guise in any other implementation.

7.3. The task learned

The networks were evaluated on the task of learning the mixture distribution shown in Table 1. There are four equally probable mixture components, each of which produces a distribution over nine visible attributes in which each attribute is independent of the others (given knowledge of the mixture component).

All the networks tested had a similar structure. Six interconnected hidden units were provided to allow the network to model the mixture variable, using a cluster such as in Fig. 3. (Four hidden units would have sufficed; six

Table 1
The mixture distribution to be learned.

m	$P(M = m)$	$P(V_i = v_i M = m), i = 1, \dots, 9$								
1	0.25	0.8	0.8	0.8	0.8	0.2	0.2	0.2	0.2	1.0
2	0.25	0.2	0.2	0.2	0.2	0.8	0.8	0.8	0.8	1.0
3	0.25	0.8	0.8	0.2	0.2	0.8	0.8	0.2	0.2	0.0
4	0.25	0.2	0.2	0.8	0.8	0.2	0.2	0.8	0.8	0.0

were provided to help avoid problems with local maxima.) These hidden units were connected to a set of nine visible units. For the belief networks, these connections were directed from the hidden to the visible units. The visible units were not connected to each other. All units had a bias connection.

Since the task is to model the entire distribution, the negative phase in Boltzmann machines was left completely unclamped.

The entropy of the target distribution is 7.67 bits. For this experiment, the number of training cases used, N , was 250. The particular set of training cases generated at random and used in these runs had an average value of $-\log_2 P^*(\tilde{V} = \tilde{v})$ of 7.87 bits, where $P^*(\cdot)$ is the true probability distribution. This is close to the entropy, as expected. This value is the target for $-L/N$ (the log-likelihood per training case) in network training, but due to overfitting, the training procedures might well reach values even lower than this.

The networks were also evaluated on the task of guessing the last attribute given the values of the other eight attributes. With knowledge of the real distribution, the optimal error rate on this classification task is 18.6%. Note that performance on this task is not the formal learning objective, and need not, in fact, be monotonically related to the actual objective of maximizing the likelihood.

7.4. Evaluation method

Typically, Gibbs sampling is used when applying networks such as these to a problem instance, as well as when training them. For example, the classification task would be performed by clamping the values of the eight known attributes and observing which value for the unknown ninth attribute shows up most often as the network is simulated.

This method was *not* used for most of the evaluations in this experiment, however. Instead, the exact probabilities of all 2^{15} states of the trained network were computed, and from these, the log-likelihood given the training data, its analogue for the real distribution, the performance on the classifi-

cation task for the training data, and the performance for items drawn from the real distribution were all calculated.

Of course, this method is infeasible for networks that are even slightly larger than the ones used here. It is convenient for this experiment since it eliminates statistical noise from the evaluations. In the tests of generalization performance, the classification task was performed with Gibbs sampling as well as with exact probabilities, and results were similar, as reported below.

7.5. Comparing sigmoid belief networks and Boltzmann machines

The 250 training cases drawn from the mixture distribution were used to train both sigmoid belief networks and Boltzmann machines, using values for ε of $\frac{1}{4}$, $\frac{1}{2}$, 1, 2, and so on until network behaviour became unstable. Networks with 0/1-valued units and those with $-1/+1$ -valued units were both tried.

Illustrative results are shown in Fig. 5, for $-1/+1$ -valued units and ε of $\frac{1}{4}$ and 1. Three runs are shown, in which different random seeds were used in the simulations. During each run, the log-likelihood, L , was computed exactly after 25, 50, 100, 200, 400, and 800 simulation passes. (Recall that each pass consists of a (potential) change to each unit value in each simulation, and that weights are changed after each pass.) The value of $-L/N$ in bits (i.e. using base-2 logarithms) is plotted. It is nine bits initially, since with zero weights all of the nine-element visible vectors are equally probable.

With $\varepsilon = \frac{1}{4}$, the Boltzmann machine and the sigmoid belief network behaved similarly. As ε was increased, however, the Boltzmann machine became unstable. This is seen in the figure for $\varepsilon = 1$, where the Boltzmann machine reached the 8.25-bit performance level, but thereafter failed to improve consistently. In contrast, the sigmoid belief network with $\varepsilon = 1$ simply learned at four times the rate that it did with $\varepsilon = \frac{1}{4}$. For larger ε , the Boltzmann machine became even more unstable, while the sigmoid belief network tolerated learning rates up to $\varepsilon = 4$ before becoming unstable at $\varepsilon = 8$ and above.

The instability of the Boltzmann machine with $\varepsilon = 1$ was examined at a finer time scale by evaluating the network after every learning pass for one of the runs, as performance went from 8.29 bits for $-L/N$ at pass 25 to 8.58 bits at pass 50. Changes in $-L/N$ of as much as 0.41 bits were seen after single learning passes, and $-L/N$ ranged in value from 8.25 bits to 8.83 bits during this interval. Examination at this time scale of learning in a sigmoid belief network simply shows steady improvement.

Results using 0/1-valued units were similar, except that learning was slower for a given value of ε in both types of network. This was largely compensated for with sigmoid belief networks by the fact that a larger ε

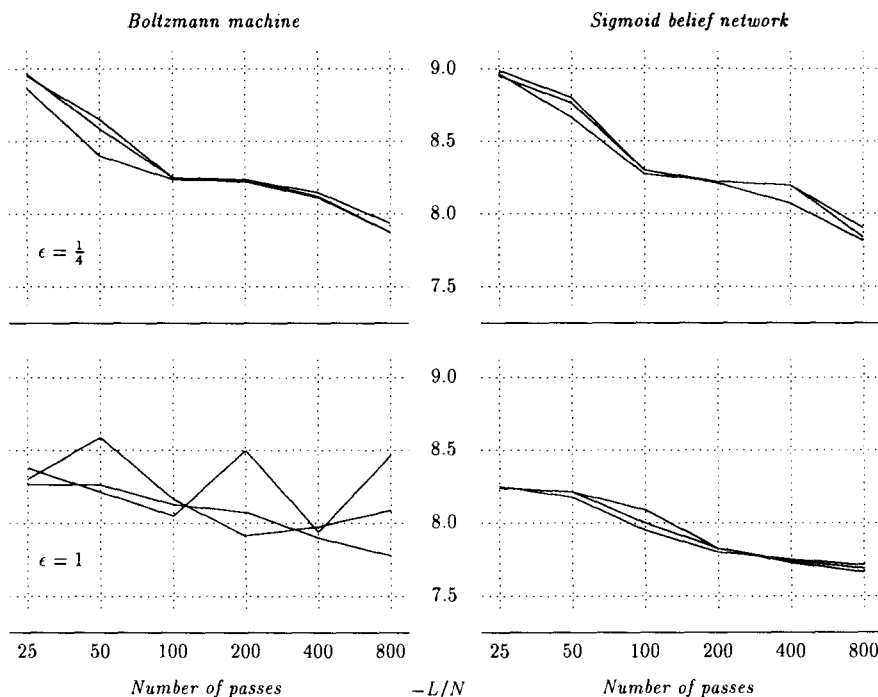


Fig. 5. Learning performance of Boltzmann machines and sigmoid belief networks with $-1/+1$ -valued units, for ϵ of $1/4$ and 1 . Three runs with different random number seeds are shown.

could be used before instability set in. With Boltzmann machines, however, there appeared to be some net advantage for the $-1/+1$ formulation.⁹

The relative performance of these networks is shown in Table 2, under the assumption that learning must be stopped after 200 passes. This would produce a fair comparison if the computation time per pass was the same for all networks. In fact, Boltzmann machine passes require somewhat more time, as would be expected from the need to simulate negative phase cases, so the comparison is somewhat biased in favour of Boltzmann machines. (The times shown, measured on a machine rated at approximately 20 MIPS, should not be taken too seriously, since they are affected by many implementation factors that may not be of general significance.)

The entries in Table 2 were produced by selecting the value of ϵ that gave the best value of $-L/N$ after 200 passes, averaged over the three runs that were done. These values are shown along with the corresponding error rates

⁹Clear differences between the $0/1$ and $-1/+1$ formulations are seen in other problems. For example, with both Boltzmann machines and sigmoid belief networks, learning to assign high probability to only those 4-bit visible vectors with odd parity, using four hidden units to express correlations, is much easier with $-1/+1$ units than with $0/1$ units.

Table 2
Best performances after 200 learning passes (three runs each).

Type of network	Time/pass	Best ε	Values of $-L/N$			Error rates		
Boltzmann machine (0/1 units)	0.39 s	1/2	8.30	8.22	8.23	35%	30%	33%
Boltzmann machine (-1/+1 units)	0.58 s	1/2	7.94	8.25	8.10	19%	37%	25%
Sigmoid belief network (0/1 units)	0.30 s	4	7.76	7.82	7.77	19%	19%	17%
Sigmoid belief network (-1/+1 units)	0.35 s	2	7.72	7.74	7.74	17%	17%	16%

when guessing the last attribute of the training cases from the first eight attributes. The superiority of the sigmoid belief networks is evident. The high error rates for the Boltzmann machine (especially using 0/1-valued units) is due to the fact that all these networks initially learn correlations among the first eight attributes, and only later discover how the ninth attribute relates to these. Four of the six Boltzmann machine runs had not progressed far into the second stage after 200 passes.

Further experiments established that the superiority seen for sigmoid belief networks over Boltzmann machines was not related to the fact that weights were updated based on a single state vector. The instability of Boltzmann machine learning for large values of ε was also found to be only slightly reduced by the use of annealing. Details of these experiments are reported in [12].

7.6. Interpretation of the results

These experiments show that a sigmoid belief network can learn the target mixture distribution faster than a Boltzmann machine. This difference is due to the sigmoid belief network's tolerance of a high learning rate that causes instability in the Boltzmann machine. Since this instability is apparent at the time scale of a single learning pass, and since it is not due to the lack of annealing, it appears that it results simply from the sampling noise in the calculation of the gradient from the results of positive and negative phase Gibbs sampling simulations.

One advantage of belief networks in this respect may be seen clearly when there are no hidden units. In this case, the positive phase, clamped simulations are completely deterministic, while the negative phase, unclamped simulations remain stochastic. Learning in a belief network, for which only

the positive phase simulation is necessary, will then take place with no noise disturbing the measurement of the gradient. Learning in a Boltzmann machine, which requires a negative phase, will still be subject to noise. When hidden units are present, the estimate of the gradient in the belief network will have some noise, but still not as much as in the Boltzmann machine.

This is not the full explanation of the difference, however, as was seen in experiments where the sigmoid belief network was trained with a redundant unclamped phase, using the “short-cut” simulation method (see Section 3.2) to ensure that state vectors came from the true equilibrium distribution. Regardless of whether this redundant phase was negative (as in Boltzmann machines) or positive (equivalent to a second set of 250 unclamped training cases) its inclusion did *not* induce instability, but merely introduced a bit more variability in the progress of learning (see Fig. 6, below).

The difference appears to result from the way weights are changed in the two networks. In a belief network, each change to w_{ij} is weighted by the forward conditional probability of S_i having a value different from that it presently has. As learning progresses, these weighting factors tend to decrease, leading to stability. In Boltzmann machines, the magnitude of each change remains constant; it is only the balance between positive phase increments and negative phase decrements that, in theory, brings learning to a stable halt, but this balance is sensitive to noise.

7.7. Effect of failure to reach equilibrium

Although it was not a major factor in the experiment described here, failure to reach equilibrium during Boltzmann machine learning is noted as a problem in [8], where it is observed that after a period of good progress learning can “go sour”, as weights are built up to values where they form large energy barriers that inhibit settling to the equilibrium distribution. The authors prescribe “weight decay” as a partial solution.

One would think that learning could “go sour” in belief networks as well as in Boltzmann machines, but such problems have not been observed. However, it *is* possible to make the sigmoid belief network go sour in the mixture distribution experiment by adding a redundant, unclamped negative phase, simulated in the normal manner (i.e. with no short-cut). This is seen in Fig. 6. Using $-1/+1$ -valued units with $\epsilon = 2$, learning with the redundant negative phase closely matches that without a negative phase for about the first 100 passes, but then becomes unstable. Interestingly, adding a redundant, unclamped *positive* phase does *not* cause the learning to go sour.

These results can be understood by picturing the effects of failing to sample from the true equilibrium distribution in the various phases. In a clamped positive phase, the effect will be to confine the state vectors seen to

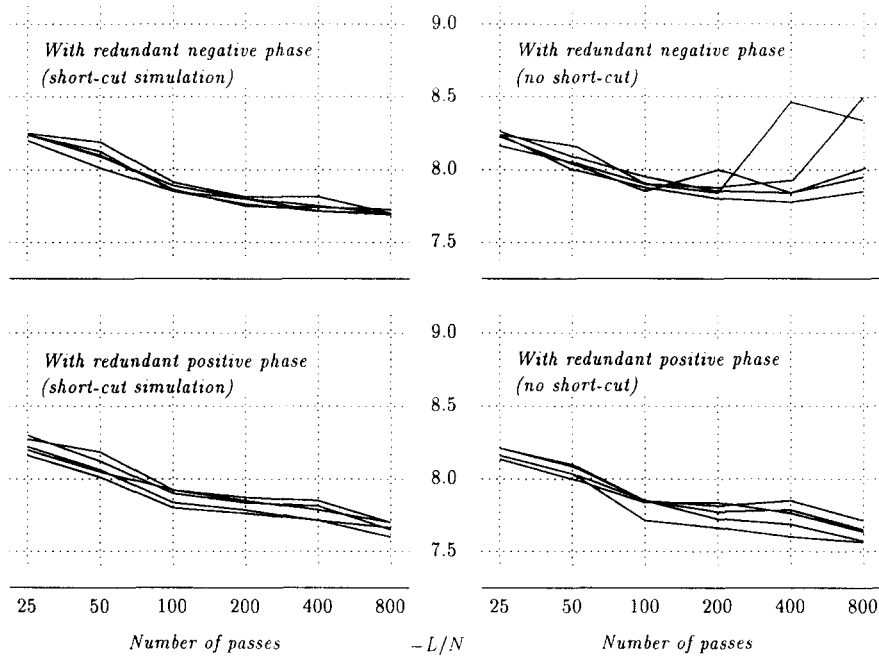


Fig. 6. Effect of redundant phases on learning performance in sigmoid belief networks ($-1/+1$ -valued units, $\epsilon = 2$). Five runs with different random number seeds are shown.

a subset of those high probability state vectors that are compatible with the clamped visible units. The learning increments that result from this sample will still tend to increase the probability of the clamped training data, albeit at a lesser rate than would be the case if *all* the compatible high probability state vectors had been seen.

In an unclamped negative phase, failure to sample from the equilibrium distribution will produce state vectors that do not represent all those of high probability. Once learning has made some progress, there will be a group of high probability state vectors compatible with each training case. In a non-representative sample, some of these groups may not be sampled from at all, while other groups will contribute more than their share of state vectors to the sample. The learning decrements that occur in the negative phase will then unfairly decrease the probability of the training cases compatible with the over-sampled groups, more than offsetting the increments in the positive phase and producing instability.

Similarly, an extra unclamped positive phase results in some training cases being over-sampled, and thus weighted more heavily in the learning. This may produce suboptimal progress, but not instability. In fact, runs done with an extra unclamped positive phase (simulated without use of the short-cut) were notable for a high variability—some runs did significantly better than

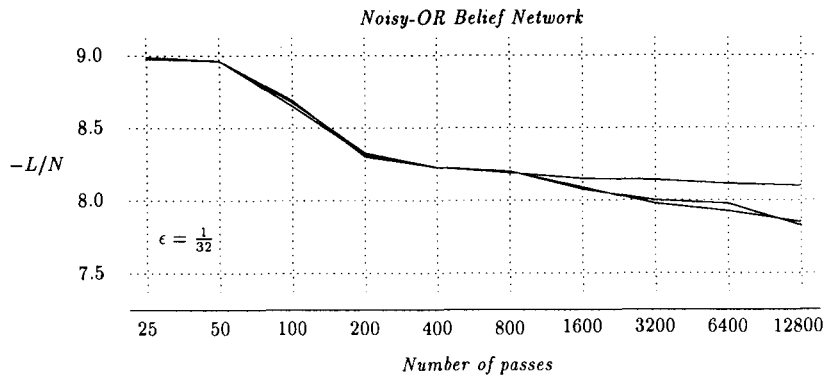


Fig. 7. Learning performance of a noisy-OR belief network with $-1/+1$ -valued units for ϵ of $1/32$. Three runs with different random number seeds are shown.

those without the extra phase, while others did rather worse. Nevertheless, even the less successful runs showed nearly steady improvement as the simulations progressed, in contrast to the drastic worsening seen at times with an extra negative phase.

Thus, it appears that the consequences of failure to reach equilibrium are more serious in a negative phase than in a positive phase. This gives belief networks a qualitative advantage in circumstances where equilibrium is hard to reach—learning may be adversely affected, but the instability that can occur with Boltzmann machines does not arise.

7.8. Performance of noisy-OR belief networks

Noisy-OR belief networks were also applied to the task of learning the mixture distribution, with rather disappointing results. Performance was both poorer and more erratic than for sigmoid belief networks.

Figure 7 shows the progress of three runs using $-1/+1$ -valued units, with $\epsilon = \frac{1}{32}$. The networks appear to have difficulty learning to reduce $-L/N$ to less than 8 bits. Increasing ϵ sometimes improved learning speed, but not reliably so. Performance of noisy-OR networks with $0/1$ -valued units was essentially similar, except that a higher value of ϵ was desirable.

In other experiments, noisy-OR networks sometimes showed a strong tendency to get stuck at a local maximum (or at a point where the gradient was so small that learning essentially stopped). For example, attempts to train a noisy-OR network with $0/1$ -valued units to compute XOR using two hidden units between inputs and output (the minimum required) succeeded in only 1 out of 20 tries.¹⁰ Somewhat better results were obtained using

¹⁰Some details: The two hidden units were connected to the inputs, but not to each other. The output unit was connected to the inputs and to the hidden units. Training was done for 5000 passes with $\epsilon = 1/8$.

	0.05	0.05	0.05	0.05	0.05	0.05	0.05	0.05	0.05	0.05	
(a)	0.29	0.00	1.20	1.20	0.00	0.00	0.00	0.00	0.00	1.20	0.00
	0.29	0.00	0.00	1.20	0.00	0.00	0.00	1.20	0.00	1.20	1.20
	0.29	1.20	0.00	1.20	1.20	0.00	0.00	0.00	0.00	0.00	1.20
	0.29	0.00	0.00	1.20	0.00	0.00	0.00	1.20	1.20	0.00	0.00
		0.04	0.01	0.05	0.04	0.01	0.04	0.03	0.02	0.05	0.07
(b)	0.32	0.00	1.80	1.22	0.12	0.03	0.02	0.14	0.10	1.32	0.05
	0.19	0.00	0.00	0.88	0.02	0.08	0.00	1.51	0.27	1.44	0.96
	0.24	1.80	0.03	1.65	1.02	0.01	0.00	0.00	0.00	0.10	2.06
	0.57	0.11	0.05	1.16	0.00	0.00	0.00	0.75	0.52	0.08	0.00
	0.11	0.00	0.00	0.00	0.39	0.30	0.23	0.00	0.02	0.00	0.27

Fig. 8. Learning a two-level noisy-OR network: (a) the true network; (b) network learned from 200 training cases. In each case, the left column contains the bias weights of the four (five) hidden “disease” units, the top row contains the bias weights of the ten visible “symptom” units, and the body of the table contains the weights from disease units to symptom units. The hidden units in (b) have been manually re-ordered to correspond to those in (a).

−1/+1-valued units—success in 10 out of 20 attempts. Sigmoid networks almost never get stuck when solving this problem, even with only one hidden unit (the minimum needed with sigmoid networks).

An additional experiment was done to test whether the learning procedure for the noisy-OR network performs better when the distribution to be learned can be simply represented in the noisy-OR form. A two-level, 0/1-valued, noisy-OR belief network of the type used for medical diagnosis in [7] was constructed. In this network, four hidden units represent “diseases” which occur independently in each patient with 25% probability. Ten visible units represent “symptoms”. Each disease has three or four potential symptoms, each of which is produced with 70% probability. Symptoms also have a 5% probability of occurring spontaneously. The weights embodying this network are shown in Fig. 8(a).

This manually constructed network was used to randomly generate 200 training cases, which were then used to train a 0/1-valued noisy-OR network of the same form (but with five hidden units rather than four, to allow escape from local maxima). In order to ease interpretation of the network learned, weights were constrained to be nonnegative.¹¹ Note that only the symptoms of each “patient” were used for training—the network was not told the true

¹¹Similar results were obtained without this constraint, except that the presence of a disease was sometimes represented by a hidden unit being 0, rather than 1. Results using −1/+1 units were also similar.

Table 3

Performance on training data (first line) and on items from the real distribution (second line) for various networks trained to near-convergence. Results from three runs with different random number seeds are shown.

Type of network	Passes	ϵ	Values of $-L/N$			Error rates (exact ~ simulated)					
Boltzmann machine (-1/+1 units)	1600	1/4	7.80	7.81	7.77	19%	17%	18% ~ 20%	18%	18%	
			7.90	7.85	7.84	19%	19%	19% ~ 20%	17%	20%	
Sigmoid belief network (-1/+1 units)	800	2	7.64	7.66	7.70	15%	17%	18% ~ 18%	18%	20%	
			7.91	7.89	7.87	19%	20%	19% ~ 19%	19%	20%	
Noisy-OR belief network (-1/+1 units)	12800	1/32	7.81	7.85	8.08	19%	19%	35% ~ 18%	19%	33%	
			7.82	7.86	8.12	19%	19%	35% ~ 18%	20%	36%	
Mixture model (EM algorithm)	100	-	7.73	7.74	7.72				18%	18%	20%
			7.82	7.85	7.86				19%	22%	19%

set of diseases underlying a training case, nor even the number of diseases present in the population.

The network that resulted from training for 2000 passes with $\epsilon = \frac{1}{16}$ is shown in Fig. 8(b). A reasonably close correspondence with the weights in the true network is apparent, with the extra hidden unit in the trained network being largely unused. The learning procedure may thus be said to have discovered the essential structure of this distribution.

7.9. Generalization performance

All the results concerning the mixture distribution shown so far give the performance of the networks on the training cases. Generally, the true objective is good performance on items drawn from the real distribution of which the training cases are a sample.

Table 3 shows the performance of all the network types, using -1/+1-valued units, on both the training data and on items from the real distribution. (Results for 0/1-valued units were similar or worse; they may be found in [12].) Each type of network was trained with a reasonable value of ϵ until performance on the training data approached convergence. (The choice of ϵ and the point of near-convergence were both subjectively determined.) The value of $-L/N$ and the classification error rates for the training data are shown, along with the corresponding figures for the real distribution. (The analogue of $-L/N$ for the real distribution is $-\sum_{\tilde{v}} P^*(\tilde{V} = \tilde{v}) \log P(\tilde{V} = \tilde{v})$, where $P^*(\cdot)$ is the real distribution, and $P(\cdot)$ that given by the network.)

Classification error rates shown in the table were calculated in two ways. The first calculation uses the exact, real distribution, and assumes that

classification is based on the exact probabilities defined by the networks. The second calculation is based on a sample of 1000 test items drawn from the real distribution which were classified by clamping the first eight attributes and using Gibbs sampling (with annealing) to observe the resulting values for the ninth attribute. Results of the two methods were similar, showing that classification performance is not dependent on very small differences in probabilities that would be swamped by noise when Gibbs sampling is used.

For comparison, results from a maximum-likelihood fit of a mixture model with six components to the training data using the EM algorithm [21] are given as well, evaluated on a test sample of 5000 items.

The sigmoid belief network, the Boltzmann machine, and the mixture model all show signs of overfitting the data, since their values for $-L/N$ on the training data are less than the value of 7.87 bits that the true model would give. Accordingly, it is not surprising that their performance on the real distribution is not quite as good as on the training data. The mixture model appears to have overfitted to a lesser extent than the networks. This is expected, since it is a more restricted model that nevertheless can exactly represent this particular distribution, but the penalty in overfitting paid for the generality of the network models does not seem large. The EM algorithm does take considerably less time than any of the network training procedures, however.

The noisy-OR belief network did not show such definite signs of overfitting the training data. One run did poorly on both the training data and on the real distribution. The other two performed well on the real distribution—slightly better, in fact, than the other two networks. This is probably an ironic consequence of the noisy-OR network's generally inferior learning performance, which would make convergence to an overfitted solution more difficult, though it could possibly be due to the particular representational capabilities of the noisy-OR network matching this problem well.

Generalization performance for all these networks might well be improved by using a cross-validation criterion [20] to stop learning before convergence, or to select an optimal number of hidden units. Use of weight decay [8] might also help.

7.10. Summary of empirical results

To summarize, the experimental results show that the sigmoid and noisy-OR belief networks are capable of learning to model a nontrivial distribution, that the sigmoid belief network can learn at a higher rate than the Boltzmann machine, and that this advantage over the Boltzmann machine is due to the elimination of the negative learning phase. The $-1/+1$ formulations of all networks appeared to outperform the $0/1$ formulations, though this

point was not investigated in detail. The generalization performance of the networks was found to be broadly similar, though again only preliminary investigations were undertaken.

The noisy-OR belief network learned the mixture distribution considerably less well than did the sigmoid belief network, and also did poorly at several other tasks. However, the noisy-OR network did perform well when learning a distribution that was naturally represented in the noisy-OR form.

How well the learning procedures for belief networks perform on larger, real-world problems can only be determined by experience. However, I expect the superiority of the sigmoid belief network over the Boltzmann machine to be at least as great for large networks as for the moderate-size networks examined here. The Boltzmann machine's problems when equilibrium is hard to reach are likely to be more apparent for larger networks, since the state space that needs to be explored is larger. The problem of sampling noise might appear to be lessened for larger networks, which will have correspondingly larger training sets, but contrary to this, the magnitude of the gradients that must be estimated tend (initially) to be smaller for complex networks.

As an aside, it is interesting that the weights learned for the mixture distribution generally bore only a vague resemblance to those that would result from manually solving the problem using the clusters of Fig. 3 to represent mixture components.

8. Discussion

I conclude by discussing how the learning procedures for belief networks described in this paper relate to other connectionist approaches to statistical modeling and to work on the representation of expert knowledge. I also outline some areas in which this work appears to open up new possibilities.

8.1. Relation to deterministic classifier networks

Problems such as speech or handwriting recognition are fundamentally statistical in nature. Although some *a priori* knowledge of the task may be available, much of the information required to solve the problem must come from training data. The preferred output for such classification problems is a probability distribution over possible classes, conditional on the attributes presented as input.

A deterministic feedforward network, trained by a method such as back-propagation [16], can represent a distribution over two classes by simply producing the probability of one of the classes as its output. Such a network that uses the sigmoid function to compute the output of a unit from its

weighted input appears very similar to a sigmoid belief network with the same structure. Indeed, the two networks are essentially equivalent if there are no hidden units. However, in the general case, this is not so, since the hidden units in a deterministic network take on fixed real values, while those in a stochastic network represent a distribution over binary vectors.

Distributions over more than two classes can be represented in a deterministic network using a cluster of output units, one for each possible classification. The output of unit c , representing $P(\text{Class} = c \mid \text{Input})$, is set to $\exp(X_c) / \sum_i \exp(X_i)$, where X_i is the total input of unit i [2]. Distributions over a vector of output attributes can be represented using several such clusters, under the assumption that the probabilities for the various attributes are independent.

Stochastic networks, such as belief networks and Boltzmann machines, have the more general capacity to exhibit distributions over a large output vector in which there are arbitrary dependencies among attributes. For example, in a medical diagnosis context, a stochastic network can represent a diagnosis that the patient has either disease A, or disease B, but likely not both. The improved learning speed of belief networks over Boltzmann machines may make use of such networks feasible in practice. In a belief network where the input units precede all the hidden and output units, as in Fig. 4(b), the short-cut simulation method can be used to produce possible classifications for a given input without the need to settle to equilibrium, at a speed comparable to that of a deterministic network. Settling to equilibrium is still necessary during learning, when the output units are clamped.

A further advantage of stochastic over deterministic networks is their superior ability to cope with missing data, especially when architectures such as those of Fig. 4 (a), (c), or (d) are used. For problems where the advantages of a stochastic network are not relevant, however, deterministic networks are likely to remain the best choice, since the exact calculation of gradients they support will generally allow faster learning.

8.2. Application to unsupervised learning

The uses of belief networks and other stochastic networks are not confined to classification problems. They are also naturally applicable to unsupervised learning, in which the objective is simply to discover the underlying structure of the data, without addressing any explicit classification task (though the resulting network may well be useful for classification). One natural measure of success in unsupervised learning is how well the probability distribution of the data has been modeled, and this is the formal objective of the learning procedure for belief networks. The experimental evaluations of learning in belief networks in Section 7 were of an unsupervised nature, with the

tasks being to model the mixture distribution of Table 1 and the two-level disease/symptom distribution of Fig. 8(a).

Standard statistical methods such as the EM algorithm can be applied to unsupervised learning problems, as is done for the case of mixture models in the AutoClass system of [3]. As was seen in Section 7, the learning procedures for belief networks are capable of discovering mixture models when they are appropriate, but they also have the capacity to learn models with a componential structure, such as that of the two-level belief network of Fig. 8(a), in which the various diseases can occur independently, but jointly influence the symptoms observed. To represent the distribution of Fig. 8(a) by a mixture model would require 2^4 mixture components, one for each possible combination of diseases. This would be impractical for large numbers of diseases.

A similar situation arises in the context of Hidden Markov Models, which are much used in speech recognition [10]. The mean field variety of the Boltzmann machine was used in [22] to learn a model that economically represented Hidden Markov Model states possessing a componential structure. The learning procedures for belief networks described here should also be applicable to this problem.

8.3. Relation to expert systems

In applications such as medical diagnosis, experts have extensive knowledge relevant to the task. Empirical data, while valuable, may be of limited extent, or may have been acquired under circumstances different from those currently prevailing. The need in such applications to integrate knowledge derived from experts with that derived from empirical data has been recognized by workers in the area (see the discussion in [9], for example). The learning procedures described in this paper may contribute to solving this problem.

One possible approach would be for the expert to specify the structure of a belief network, while leaving the numeric values of the forward conditional probabilities to be estimated empirically. If training data is available in which all attributes are known, this will be straightforward. It is likely, however, that the belief network will contain units whose values were not always measured, or which are not directly observable (such as the true underlying disease a patient suffered from). In this case, the gradient-ascent learning procedures of this paper could be applied, perhaps starting with weight values derived from an expert's tentative assessment of the probabilities. The expert might also constrain probabilities to some interval in order to guard against training data that is not extensive enough, or that is not representative of all possible contexts in which the system might be used. Another possibility would be for the expert to

construct artificial, “textbook” training cases to supplement the empirical data.

More ambitiously, in parts of the network where causal connections are not clear to the expert a pool of hidden units could be included and their weights trained from empirical data. A problem with this approach is that the resulting networks may be hard to interpret. Using “weight decay” [8] to encourage some weights to go to zero might help.

The desire to keep the network’s operation intelligible to the experts might also lead one to use the noisy-OR model for conditional probabilities, regardless of whether the learning performance of sigmoid units might be better. The particular properties of the noisy-OR model might also be desirable for technical reasons; they are exploited in the heuristic diagnostic search algorithm of [7], for example.¹² Noisy-OR and sigmoid units can also be mixed in the same network, and for that matter, incorporating a Boltzmann machine as a subnetwork is not impossible.

8.4. Making decisions

Belief networks are compatible with the “influence diagrams” used to formulate decision problems. An algorithm of Shachter [17] exploits the structure of these diagrams to find decisions that maximize expected utility. Unfortunately, this algorithm can sometimes take exponential time. I will describe here a method of making simple decisions using Gibbs sampling that also exploits properties of belief networks.

Consider a network with three sets of visible units—a “context” set, \tilde{C} , an “action” set \tilde{A} , and a “result” set, \tilde{R} . Using empirical data, we can train this network to represent the conditional probabilities that \tilde{R} will result given that we perform action \tilde{A} in context \tilde{C} . Suppose now that we wish to bring about some “goal”, \tilde{g} , at a time when the context is \tilde{c} . Our best bet is to perform an action \tilde{a} that maximizes $P(\tilde{R} = \tilde{g} \mid \tilde{A} = \tilde{a}, \tilde{C} = \tilde{c})$.

We could find the action that maximizes the probability of our goal by running a separate Gibbs sampling simulation for every possible action. In each simulation, the action and context units would be clamped, and we would observe how often the goal shows up in the result units. We would then choose the action that leads to the goal showing up most often. However, this method is infeasible if there are many actions, represented by a large number of units. (Consider the number of possible medical treatments when twenty drugs can be given in combination, for example.)

¹²If this algorithm is to be used, the noisy-OR network must be constrained to allow only nonnegative weights.

However, we can transform the problem by rewriting the probability to be maximized using Bayes' rule:

$$\begin{aligned} P(\tilde{R} = \tilde{g} \mid \tilde{A} = \tilde{a}, \tilde{C} = \tilde{c}) \\ = \frac{P(\tilde{A} = \tilde{a} \mid \tilde{R} = \tilde{g}, \tilde{C} = \tilde{c}) P(\tilde{R} = \tilde{g} \mid \tilde{C} = \tilde{c})}{P(\tilde{A} = \tilde{a} \mid \tilde{C} = \tilde{c})}. \end{aligned} \quad (35)$$

Now, *provided* that $P(\tilde{A} = \tilde{a} \mid \tilde{C} = \tilde{c})$ is the same for all \tilde{a} , we can choose the best action by running a Gibbs sampling simulation in which we clamp the context units to \tilde{c} and the result units to \tilde{g} , and then observe which value of \tilde{a} turns up most often in the action units.

Ensuring that $P(\tilde{A} = \tilde{a} \mid \tilde{C} = \tilde{c})$ is the same for all \tilde{a} is easy in a belief network—we simply set up the network so that units in \tilde{A} have no incoming connections, ensuring equal probabilities for each \tilde{a} in an unclamped network, and we further arrange that there is no directed path from a unit in \tilde{A} to a unit in \tilde{C} , which ensures that clamping \tilde{C} will not change these probabilities.¹³ Producing these equal probabilities in a Boltzmann machine is not so easy. We could try to train the Boltzmann machine to satisfy the constraint, but there is no guarantee that we will succeed very well, and the attempt may interfere with learning the distribution of \tilde{R} given \tilde{A} and \tilde{C} .

Unfortunately, the transformed method does not completely solve this decision problem. It is possible that with a particular context and goal clamped, a different action will show up after every simulation pass, even during a long simulation. We will then have no basis for deciding which (if any) of these actions is best. Accordingly, the method is most applicable in situations where only a small, but unknown, subset of actions have a significant probability of producing the goal.

It is tempting to try to solve this problem using simulated annealing by “cooling” the simulation down to a temperature of zero in order to find the most probable state vector compatible with the clamped context and goal. However, the action with highest *marginal* probability (i.e. probability after summing over all possible hidden unit values) need not be the same as the action part of the most probable total state vector, so the annealing method is guaranteed to work only if there are no hidden units in the network.

¹³Note that $P(\tilde{A} = \tilde{a} \mid \tilde{C} = \tilde{c})$ exists only in a formal sense, and may thus be manipulated in any fashion that is convenient. With a usual degree-of-belief interpretation of probability, we do not assess how likely we are to perform action \tilde{a} , we simply decide whether or not to do it.

8.5. Potential for new learning procedures

Gradient-ascent learning has the advantage that it is simple, and that it can be performed in an “on-line” manner if desired. However, it can be rather slow, and can get stuck at a local maximum. For Boltzmann machines, there appears to be no reasonable alternative to gradient-ascent, but for belief networks the fact that the probability of a full state vector can be explicitly calculated allows one to contemplate other learning procedures.

For noisy-OR belief networks, one possibility is to apply a stochastic version of the EM algorithm [4]. This seems feasible provided that the efficacy of each input to a unit in forcing the unit to take on the value 1 is made explicit in a set of auxiliary units that are simulated along with the main units. Probabilities can then be iteratively estimated from co-occurrence counts.

It also seems feasible to implement Bayesian learning by applying Gibbs sampling to the learning process as a whole. This method may avoid both the problem of overfitting the training data and the possibility of getting stuck in a local maximum. It may work especially well for noisy-OR networks with the auxiliary units described above, since it turns out that one can then avoid having to simulate distributions over continuous parameters.

One advantage of all the learning methods based on Gibbs sampling is the ability to easily handle missing data, which is an inherent aspect of learning whenever the network contains hidden units. This is a problem for previously described learning methods for belief networks, such as those of [18].

8.6. Neural modeling

The connectionist learning procedures for belief networks also provide additional options for modeling of real neural processes. Although analogies between the negative phase of Boltzmann machine learning and dream sleep are speculated upon in [8], it may well turn out that the negative phase is biologically implausible. The work here shows that this would not necessarily be fatal to the idea that gradient ascent using Gibbs sampling plays some role in learning in the brain. The somewhat greater complexity of Gibbs sampling in belief networks may be a barrier to their incorporation in models of neural processing, however.

Appendix A. Proofs of claims

Claim A.1. *The procedure of Fig. 1 moves the weights into unit i to the point closest in Euclidean distance that satisfies the constraint*

$$w_{i0} + \sum_{j < i, w_{ij} < 0} w_{ij} \geq \eta. \quad (\text{A.1})$$

Proof. Let w_{ij} be the original set of weights, and let $w'_{ij} = w_{ij} + \delta_j$ be the set of weights satisfying the constraint for which $\Delta^2 = \sum_j \delta_j^2$ is minimal. We can prove a number of properties of the δ_j .

First, all the δ_j are nonnegative, since decreasing a weight will certainly not help satisfy the constraint. Also, for $j > 0$, $\delta_j = 0$ if $w_{ij} \geq 0$ and $\delta_j \leq |w_{ij}|$ if $w_{ij} < 0$, since once δ_j is large enough to make w'_{ij} zero, making it any larger does not help satisfy the constraint.

Next, if $w_{ij} \leq w_{ik} < 0$, then $\delta_j \geq \delta_k$. Otherwise replacing both δ_j and δ_k by $\frac{1}{2}(\delta_j + \delta_k)$ would reduce Δ^2 while keeping the constraint satisfied. Similarly, $\delta_0 \geq \delta_j$ for all j , since otherwise there would be an advantage in replacing them both by $\frac{1}{2}(\delta_0 + \delta_j)$.

We can therefore renumber the units before i in such a way that for some n :

$$\begin{aligned} \delta_0 &\geq \delta_1 \geq \dots \geq \delta_n > 0, \\ w_{i1} &\leq \dots \leq w_{in} < 0, \end{aligned} \tag{A.2}$$

and $w_{ij} \geq 0$ and $\delta_j = 0$ for $j > n$. One can now show, by arguments similar to those above, that there is an m such that for all j up to m , $\delta_j = \delta_0$ and $\delta_j < |w_{ij}|$, while for $m < j \leq n$, $\delta_j = |w_{ij}|$.

The entire set of optimal changes, δ_j , is therefore determined by the value of δ_0 . The other δ_j are either equal to δ_0 , or are less, if a lesser value suffices to make w'_{ij} nonnegative.

The procedure of Fig. 1 is now easily seen to be a search for the appropriate value of δ_0 . \square

Claim A.2. *The weights in both a Boltzmann machine and in a sigmoid belief network consisting of only one or two units can be set so as to produce any probability distribution over state vectors. (Except that distributions in which some state vectors have zero probability can only be approached as the weights go to infinity.)*

Proof. We need only consider networks with 0/1-valued units. Clearly, any distribution over a network of one unit can be produced by simply adjusting the single bias weight.

To produce a given distribution over a sigmoid belief network with two units, start by setting the bias weight for the first unit to produce the required marginal probability distribution for that unit. Then set the bias weight for the second unit to produce the required conditional probability distribution given that the first unit has value 0. Lastly, set the weight on the connection from the first to second unit to produce the correct conditional probability distribution given that the first unit has value 1, taking into account the value of bias weight determined earlier.

For a two-unit Boltzmann machine, we must find weights that give energies to the four possible states that produce the required distribution. The energy of state $\langle 0, 0 \rangle$ is zero irrespective of the weight values. We can arrange for states $\langle 1, 0 \rangle$ and $\langle 0, 1 \rangle$ to have the appropriate energies relative to that of $\langle 0, 0 \rangle$ by adjusting their bias weights. The energy of state $\langle 1, 1 \rangle$ can then be made whatever we wish by setting the weight on the connection between the two units, taking account of the bias weights. \square

Claim A.3. *The set of probability distributions that can be produced over a network of three units is the same for Boltzmann machines and sigmoid belief networks.*

Proof. To translate a three-unit Boltzmann machine to a sigmoid belief network, start by setting up the first two units of the sigmoid network so as to duplicate the marginal distribution over (any) two units of the Boltzmann machine. Claim A.2 guarantees that this is possible. Now add a third unit after these two, connected to the first two using the same weights as in the Boltzmann machine. This duplicates the required conditional probabilities for the third unit, without disturbing the distribution over the first two units.

To translate a three-unit sigmoid belief network to a Boltzmann machine, start by setting the weights to the third unit in the Boltzmann machine to be the same as those into the last unit in the sigmoid network. This duplicates the conditional probabilities for this unit given the values of the other two units. Now we need to set up the weights between the remaining two units so as to produce the same marginal distribution as in the sigmoid network, taking into account the biasing effects of the third unit. This can be done because, again, all things are possible with only two units. \square

Claim A.4. *The probability distribution produced by the 0/1-valued Boltzmann machine of Fig. 2(a) cannot be duplicated by a sigmoid belief network with the same number of units.*

Proof. We can assume that the sigmoid belief network also uses 0/1-valued units. Due to the symmetry of the Boltzmann machine, there is no choice in ordering the units when trying to find an equivalent sigmoid belief network. The last unit in the sigmoid belief network (unit 4) must have the same weights as in the Boltzmann machine in order to reproduce the conditional probabilities for that unit's value given the values in the rest of the network.

Now consider how we must set the weights into the second-to-last unit in the sigmoid belief network (unit 3). By symmetry, the two weights from the earlier units must be equal; call their value w . There is also a bias weight, b . Consider the odds in favour of unit 3 having the value 1 when unit 4

has the value 1 and there are zero, one, or two units with value 1 before unit 3. Equating these odds in the sigmoid belief network to the odds in the Boltzmann machine produces the constraints, respectively:

$$\exp(b) \frac{\sigma(1)}{\sigma(0)} = \exp(1), \quad (\text{A.3})$$

$$\exp(b + w) \frac{\sigma(2)}{\sigma(1)} = \exp(2), \quad (\text{A.4})$$

$$\exp(b + 2w) \frac{\sigma(3)}{\sigma(2)} = \exp(3). \quad (\text{A.5})$$

By taking logarithms, one obtains a system of linear equations in w and b which numerical calculation shows to be inconsistent. \square

Claim A.5. *The probability distribution produced by the 0/1-valued sigmoid belief network of Fig. 2(b) cannot be duplicated by a Boltzmann machine with the same number of units.*

Proof. We can restrict consideration to Boltzmann machines with 0/1-valued units. Consider the unit in a candidate Boltzmann machine corresponding to the bottom unit in the sigmoid network. The weights into this unit must be the same as those in the sigmoid network, in order to reproduce the conditional probabilities for this unit given the various combinations of other unit values. Note that the value of the bottom unit is effectively a deterministic function of the values of the upper three units—i.e. there are only eight state vectors with significant probability.

Now consider the constraints placed on the weights in the Boltzmann machine by the requirement that all combinations of values for the upper three units be equally probable, as they are in the sigmoid network. In the Boltzmann machine, this translates to the requirement that the energy of the network be the same for all eight possible state vectors. In particular, since the energy of the state vector $\langle 0, 0, 0, 0 \rangle$ is zero, the energy of the other seven state vectors must be zero as well. Applying this constraint to the three state vectors $\langle 1, 0, 0, 1 \rangle$, $\langle 0, 1, 0, 1 \rangle$, and $\langle 0, 0, 1, 1 \rangle$, we find that the bias weights into the three upper units must be -10 . Applying it to the the three state vectors $\langle 0, 1, 1, 1 \rangle$, $\langle 1, 0, 1, 1 \rangle$, and $\langle 1, 1, 0, 1 \rangle$, we get that the weights between the upper units must all be -10 as well. The energy of the state $\langle 1, 1, 1, 1 \rangle$ is now determined to be -10 , showing that a proper set of weights is impossible. \square

Claim A.6. *The probability distribution produced by a three-unit noisy-OR belief network with 0/1-valued units in which $w_{31} = w_{32} = 1$ and all other*

weights are zero cannot be duplicated by either a Boltzmann machine or a sigmoid belief network with only three units.

Proof. In view of Claim A.3, it suffices to show that the noisy-OR network cannot be duplicated by a Boltzmann machine with 0/1-valued units.

Consider the unit in a candidate Boltzmann machine corresponding to unit 3 in the noisy-OR network. In the noisy-OR network, this unit is always zero when the other units are both zero. To approximate this in the Boltzmann machine, the bias weight for this unit must be very large and negative. The probability of this unit being zero when one of the other two units is one is only e^{-1} , however. The weights from the other two units must therefore be very large and positive, in order to nearly cancel out the large negative bias in this case. Now, however, the probability of the unit being one when *both* other units are one is nearly 1 in the Boltzmann machine, but only $1 - e^{-2}$ in the noisy-OR network. Thus no set of weights for the Boltzmann machine can produce (or even closely approximate) the required distribution. \square

Acknowledgement

I thank Geoff Hinton and the other members of the Connectionist Research Group at the University of Toronto for many helpful discussions. This research was supported by the Natural Sciences and Engineering Research Council of Canada and the Ontario Information Technology Research Centre.

References

- [1] D.H. Ackley, G.E. Hinton and T.J. Sejnowski, A learning algorithm for Boltzmann machines, *Cogn. Sci.* **9** (1985) 147–169.
- [2] J.S. Bridle, Probabilistic interpretation of feedforward classification network outputs, with relationships to statistical pattern recognition, in: F. Fogelman-Soulie and J. Héroult, eds., *Neuro-computing: Algorithms, Architectures, and Applications* (Springer, Berlin, 1989).
- [3] P. Cheeseman, J. Kelly, M. Self, J. Stutz, W. Taylor and D. Freeman, AutoClass: a Bayesian classification system, in: *Proceedings Fifth International Conference on Machine Learning*, Ann Arbor, MI (1988).
- [4] A.P. Dempster, N.M. Laird and D.B. Rubin, Maximum likelihood from incomplete data via the EM algorithm (with discussion), *J. Roy. Stat. Soc. B* **39** (1977) 1–38.
- [5] M. Derthick, Variations on the Boltzmann machine learning algorithm, Tech. Rept. CMU-CS-84-120, Department of Computer Science, Carnegie-Mellon University, Pittsburg, PA (1984).
- [6] A.E. Gelfand and A.F.M. Smith, Sampling-based approaches to calculating marginal densities, *J. Am. Stat. Assoc.* **85** (1990) 398–409.

- [7] M. Henrion, Towards efficient probabilistic diagnosis in multiply connected belief networks, in: R.M. Oliver and J.Q. Smith, eds., *Influence Diagrams, Belief Nets and Decision Analysis* (Wiley, Chichester, England, 1988).
- [8] G.E. Hinton and T.J. Sejnowski, Learning and relearning in Boltzmann machines, in: D.E. Rumelhart, J.L. McClelland and the PDP Research Group, eds., *Parallel Distributed Processing: Explorations in the Microstructure of Cognition, Vol. 1: Foundations* (MIT Press, Cambridge, MA, 1986) 282–317.
- [9] S.L. Lauritzen and D.J. Spiegelhalter, Local computations with probabilities on graphical structures and their application to expert systems (with discussion), *J. Roy. Stat. Soc. B* **50** (2) (1988) 157–224.
- [10] S.E. Levinson, L.R. Rabiner and M.M. Sondhi, An introduction to the application of the theory of probabilistic functions of a Markov process to automatic speech recognition, *Bell Syst. Tech. J.* **62** (4) (1983).
- [11] N. Metropolis, A.W. Rosenbluth, M.N. Rosenbluth, A.H. Teller and E. Teller, Equation of state calculations by fast computing machines, *J. Chem. Phys.* **21** (6) (1953) 1087–1092.
- [12] R.M. Neal, Learning stochastic feedforward networks, Tech. Rept. CRG-TR-90-7, Connectionist Research Group, Department of Computer Science, University of Toronto, Ont. (1990).
- [13] R.M. Oliver and J.Q. Smith, *Influence Diagrams, Belief Nets and Decision Analysis* (Wiley, Chichester, England, 1988) (*Proceedings Conference on Influence Diagrams for Decision Analysis, Inference, and Prediction*, Berkeley, CA (1988)).
- [14] J. Pearl, Evidential reasoning using stochastic simulation of causal models, *Artif. Intell.* **32** (2) (1987) 245–257.
- [15] J. Pearl, *Probabilistic Reasoning in Intelligent System: Networks of Plausible Inference* (Morgan Kaufmann, San Mateo, CA, 1988).
- [16] D.E. Rumelhart, G.E. Hinton and R.J. Williams, Learning representations by back-propagating errors, *Nature* **323** (1986) 533–536.
- [17] R.D. Shachter, Probabilistic inference and influence diagrams, *Oper. Res.* **36** (4) (1988) 589–604.
- [18] D.J. Spiegelhalter and S.L. Lauritzen, Sequential updating of conditional probabilities on directed graphical structures, *Networks* **20** (1990) 579–605.
- [19] W.S. Stornetta and B.A. Huberman, An improved three-layer back propagation algorithm, in: *Proceedings First IEEE International Conference on Neural Networks*, San Diego, CA (1987).
- [20] M. Stone, Cross-validatory choice and assessment of statistical predictions (with discussion), *J. Roy. Stat. Soc. B* **36** (1974) 111–147.
- [21] D.M. Titterton, A.F.M. Smith and U.E. Makov, *Statistical Analysis of Finite Mixture Distributions* (Wiley, Chichester, England, 1985).
- [22] C.K.I. Williams and G.E. Hinton, Mean field networks that learn to discriminate temporally distorted string, in: *Connectionist Models: Proceedings of the 1990 Summer School* (Morgan Kaufmann, San Mateo, CA, 1990).