

VICTORIA UNIVERSITY OF WELLINGTON  
*Te Whare Wānanga o te Ūpoko o te Ika a Māui*



School of Engineering and Computer Science  
*Te Kura Mātai Pūkaha, Pūrorohiko*

PO Box 600  
Wellington  
New Zealand

Tel: +64 4 463 5341  
Fax: +64 4 463 5045  
Internet: [office@ecs.vuw.ac.nz](mailto:office@ecs.vuw.ac.nz)

**Richer Restricted Boltzmann  
Machines**

Max Godfrey

Supervisors: Marcus Freen

Submitted in partial fulfilment of the requirements for  
Bachelor of Engineering with Honours.

**Abstract**

TODO WRITE THE ABSTRACT



# Acknowledgments

Any acknowledgments should go in here, between the title page and the table of contents. The acknowledgments do not form a proper chapter, and so don't get a number or appear in the table of contents.



# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Problem . . . . .	1
1.1.1	Deep Belief Networks can achieve state of the art performance . . . . .	1
1.1.2	DBNs have no mechanism for separating sources . . . . .	1
1.1.3	Restricted Boltzmann Machines cannot separate sources either . . . . .	1
1.1.4	Sigmoid Belief Networks; Intractably rich in practice . . . . .	1
1.2	Solution . . . . .	2
1.2.1	Trading tractability for Source Separation . . . . .	2
1.3	Results . . . . .	2
<b>2</b>	<b>Backgroud</b>	<b>3</b>
2.1	Source Separation, nature can do it . . . . .	3
2.1.1	An example, the cocktail party problem . . . . .	3
2.2	Generative Models . . . . .	3
2.2.1	Terminology in Generative Models observable and hidden variables . . . . .	3
2.2.2	PGMs as a tool reasoning about generative models . . . . .	4
2.3	Sampling and inverting the model . . . . .	4
2.3.1	Gibbs sampling, a subset of Markov Chain Monte Carlo . . . . .	4
2.3.2	Reconstructions, visualising what the model has learnt . . . . .	5
2.4	An intractable model for causes . . . . .	5
2.4.1	Sigmoid Belief Networks . . . . .	5
2.4.2	Explaining Away creates a trade off between richness and tractability . . . . .	6
2.4.3	Boltzmann Machines . . . . .	7
2.5	The Current Approach: A Strong assumption . . . . .	7
2.5.1	Restricted Boltzmann Machines . . . . .	7
2.5.2	Deep Learning . . . . .	9
2.5.3	Inference . . . . .	9
2.5.4	Evaluating Restricted Boltzmann Machines . . . . .	9
2.6	A New Approach - The ORBM . . . . .	11
2.6.1	Architecture . . . . .	11
2.6.2	Inference In the ORBM . . . . .	11
2.6.3	Source Separation - Reconstructions in the ORBM . . . . .	11



# Figures

1.1	A graphical representation showing the proposed generative model for capturing two causes, the ORBM. . . . .	2
2.1	An example PGM, showing an observed variable ‘A’ and it’s hidden cause ‘B’.	4
2.2	The famous Burglar, Earthquake, Alarm network showing a minimal case of explaining away. . . . .	6
2.3	A Boltzmann Machine, the blue shaded nodes representing the observed variables, and the non-shaded nodes the latent variables. . . . .	7
2.4	A Boltzmann Machine, the blue shaded nodes representing the observed variables, and the non-shaded nodes the latent variables. . . . .	8
2.5	Good Hinton Diagram . . . . .	10
2.6	Bad Hinton Diagram . . . . .	10





# Chapter 1

## Introduction

### 1.1 Problem

#### 1.1.1 Deep Belief Networks can achieve state of the art performance

Deep Belief networks are powerful models that have proven to achieve state of the art performance in many domains. For instance a non-exhaustive list is image classification, dimensionality reduction, natural language recognition, Document classification, Semantic Analysis and .

DBNs capture non-linear interactions between low level features, in the context of image classification the lower layers can capture image filters.

#### 1.1.2 DBNs have no mechanism for separating sources

Despite a DBNs expressiveness, there is no way to extract these interactions. If an input has multiple sources then the complex combination is instead learnt, the network has no mechanism for extracting multiple causes. This is the motivation for this project, to be able to separate the sources of data in a new model.

[1]

#### 1.1.3 Restricted Boltzmann Machines cannot separate sources either

Restricted Boltzmann Machines are two layer, fully connected, unsupervised neural networks. DBNs are constructed by stacking RBMs. Being the building block of the powerful DBN, RBMs are a natural starting point for representing multiple sources. RBMs make the assumption that the features of the input data are dependant in the prior, as they are independant in the posterior. The latter makes them tractable to use in practice, but also means they model/encode a single representation. Again using the example of images, an input image will map to a single representation, again there is a lack of mechanism for modelling sources that are acting independantly.

#### 1.1.4 Sigmoid Belief Networks; Intractably rich in practice

The Sigmoid Belief Network, the parameterized version of a Bayesian/Belief network appears as a natural choice for modelling independant sources in that it makes a polar assumption to the RBM; – Warning Semicolon Use – Every feature has an independant cause. The sigmoid belief networks assumption could capture data that has multiple sources, but this is intractable in practice.

## 1.2 Solution

### 1.2.1 Trading tractability for Source Separation

Frean and Marsland propose a generative model that aims to trade a small amount of the RBMs performance for richness, finding a middle ground between the sigmoid belief network and the restricted boltzmann machine. Frean and Marsland also propose an algorithm to invert this model, separating the sources of an input.

The new generative model, referred to onwards as an ORBM, uses an RBM to model each source and a sigmoid belief network to capture their combination to form data. This project explores the ORBM use for separating two causes.

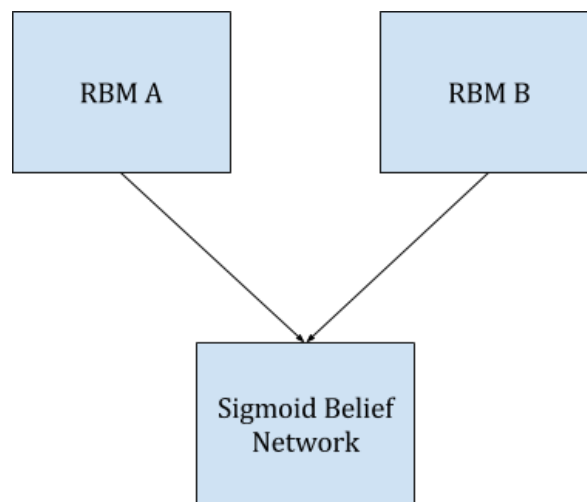


Figure 1.1: A graphical representation showing the proposed generative model for capturing two causes, the ORBM.

Given the proposed model and algorithm, this project answers the following questions:

- Can this model encode data comprised of more than one cause as it's constituted causes? That is, can the model and new algorithm for inverting it, perform source separation.
- Is the ORBMs two cause structure to rich to be tractible in practice?

## 1.3 Results

**TODO Draft Points Still to write TODO**

Spoil the results here. Good on smaller cases, less good on larger case

# Chapter 2

## Backgroud

As the ORBM builds on the previous work of Restricted Boltzmann Machines and Sigmoid Belief networks, the concepts and previous work in source separation, generative models as well as background on RBMs and Sigmoid Belief Networks need to be introduced.

### 2.1 Source Separation, nature can do it

#### 2.1.1 An example, the cocktail party problem

A famous example that illustrates the idea of source separation is the cocktail party problem. At a cocktail party many conversations are taking part at the same time, creating noise, a composition of all the conversations. Despite this, a partygoer is able to separate their conversation from cacophany, separating the sources.

The applications of source separation are far wider than talkative partygoers. In the field of signal processing

TODO Draft Points Still to write TODO

### 2.2 Generative Models

Generative models are a powerful way to model data. [TODO WORDING \(TODO-GRAB-THOSE-GENERATIVE-MODEL-USES-CITATIONS\)](#) Basically justify generative models.

The ORBM proposed in this project aims to represent data generated by two indepedanlty acting causes and does so by combining two existing generative models, the Restricted Boltzmann Machine, and the Sigmoid Belief Network.

#### 2.2.1 Terminology in Generative Models observable and hidden variables

Generative models are comprised of variables, often referred to as units. Some of these variables are observed, that is their state is known. These are often referred to as the ‘visible’ units and are used to represent the training data. For example in the image domain, the visible units correspond to the pixels of the image.

The variables that are not observed, are latent variables, often referred to as ‘hidden units’ as they are not observed.

Connections between units are used to encode relationships between the variables, where the relationship may be causal, such as in a Sigmoid Belief network or an encoding/representation in the Restricted Boltzmann Machine.

Collections of units, are often referred to as ‘patterns’ or ‘vectors’ in that they are represented by a vector or pattern of bits. For instance in the context of an image, the visible pattern would be the pixels of the image ravelled into a one dimensional vector.

### 2.2.2 PGMs as a tool reasoning about generative models

**TODO Draft Points Still to write TODO**

Probalistic Graphical Models or PGMs for short, are an expressive way to represent a collection of related, stochastic variables. If the graph is directed then the edges represent causation, this is also referred to a Bayesian network. Conversely, if the graph was undirected then edges represent a dependancy or mapping. Throughout this report, RBMs, Sigmoid Belief Networks and the proposed ORBM will be shown in this format.

Figure 2.1 shows an abstract example of a directed PGM, where B is the underlying cause of A, we cannot observe B directly, instead it’s state is represented as a ‘belief’ or a probability of being in a given state.

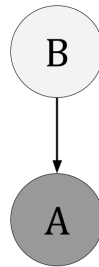


Figure 2.1: An example PGM, showing an observed variable ‘A’ and it’s hidden cause ‘B’.

## 2.3 Sampling and inverting the model

Sampling is the process of drawing samples from a distribution. It is used when the distrubution we want samples from is intractable to calculate analytically. Sampling is required to train generative models, as often the gradient to be climbed/descended involves calculating a probability in the generative model.

**TODO Draft Points Still to write TODO**

- Inference is the process of given reasoning about what we do not know given that of which we do know.
- In a Generative Model this amounts to the Posterior

### 2.3.1 Gibbs sampling, a subset of Markov Chain Monte Carlo

- The importance of Markov Chains and mixing time are crutial in this project

Gibbs sampling is a special case of Markov Chain Monte Carlo, a technique for drawing sampling from a complex distribution. The probability mass (the joint distribution) of a generative model is a common use case for Gibbs sampling.

Gibbs sampling explores the desired probability distribution, taking samples of that distributions state, allowing iterations of exploration between drawing of a sample to ensure

that the samples are independent. The process of taking a step between states is referred to as a Gibbs iteration.

Gibbs sampling is used for performing inference in the RBMs, Sigmoid Belief Networks and in the ORBM. The mixing time, that is how many Gibbs iterations are needed to reach a satisfactory sample is an important part issue in the ORBM, in that more than one may be required.

### Mixing Time

MCMC methods require a ‘mixing’ phase to ensure convergence, that is that the sample is being drawn from a representative part of the desired distribution. This is part of the trade off the ORBM attempts to make, as a mixing time is introduced that is not present in the RBM.

### 2.3.2 Reconstructions, visualising what the model has learnt

Generative Models can create an internal representation given an input. They can also generate a faux input given an internal representation. Performing one Gibbs iteration, that is sampling from the hidden units given an input  $P(\tilde{h}|\tilde{v})$  and then taking the generated hidden state and generating a faux input. The model tries to reconstruct the input.

### Fantasies of the model

TODO Draft Points Still to write TODO

In the same way that a generative model uses reconstructions to try and recreate the supplied input based purely on how it’s represented that input, performing many, many (greater than 100) Gibbs iterations with no input pattern clamped allows the reconstructions to explore the probability mass that the model has built up during training. Sampling from these wanderings creates what are referred to as ‘fantasies’ or ‘dreams’. These give a sense of what the model has learnt, and can act a smoke test for if the model has actually capture anything. (TODO-CITE-PAPER-WITH-MNIST-DREAM-EVALUATION, they were crappy).

## 2.4 An intractable model for causes

### 2.4.1 Sigmoid Belief Networks

TODO Draft Points Still to write TODO

The ORBM relies on the Sigmoid Belief Network to capture the causation. The Sigmoid Belief Network is composed of units with weights and a sigmoid activation function, akin to that of a perceptron linear threshold unit/Perceptron. The probability of a node being ‘on’ is found by taking the weighted sum of all input to that node and applying a Sigmoid function or another activation function that ensures a values between 0 and 1.

Belief Networks appear to be an intuitive way to model data in machine learning, as rich dependancies often present in real data can be expressed in it’s architecture. Nodes in the network represent binary variables which are dependent on ancestor nodes, the degree of which is encoded in a weight on a directed edge between them.

Performing inference in a Sigmoid Belief network would allow source separation in that each hidden unit could represent a cause. Meaning if a causes state could be inferred from an input item, individual causes could be examined for an input. For example if the input was an n by n image, the Sigmoid Belief Net makes the assumption that each pixel has an independent cause.

Despite the Sigmoid Belief Network being expressive and providing a succinct encoding of inter-variable dependencies, the expressiveness is too rich such that performing inference is intractable. There do exist algorithms for performing inference in Sigmoid Belief Networks. For instance, the Belief Propagation algorithm [TODO CITE: The paper where BP/Sum Prod proposed](#) operates on this encoding, calculating the probabilities of a given network state (i.e. the state of all the variables). Belief Propagation is intractable to use as the number of variables grow [TODO CITE: the paper explaining intractable for belief prop.](#)

This intractability arises from the Sigmoid Nets richness and the ‘explaining away effect’. Inference is required for training generative models making Sigmoid Belief Networks impractical to train. [TODO CITE: It has been done, link to paper where they do it.](#)

### 2.4.2 Explaining Away creates a trade off between richness and tractability

[TODO Draft Points Still to write TODO](#)

The power of the Belief Network is also it’s weakness, a rich structure that models a system of interest inherently has dependencies. In its minimal case explaining away can be seen in a 3 node network popularised by [TODO CITE: AI-A-MODERN-APPROACH-TODO. TODO-GRAPHIC](#) as shown in figure 2.2. Each of the nodes represents a binary state. For instance *Burglar* = 1 means that the person owning the Alarm has been burgled. Also note how the connections between the units have arrows, this is causal.

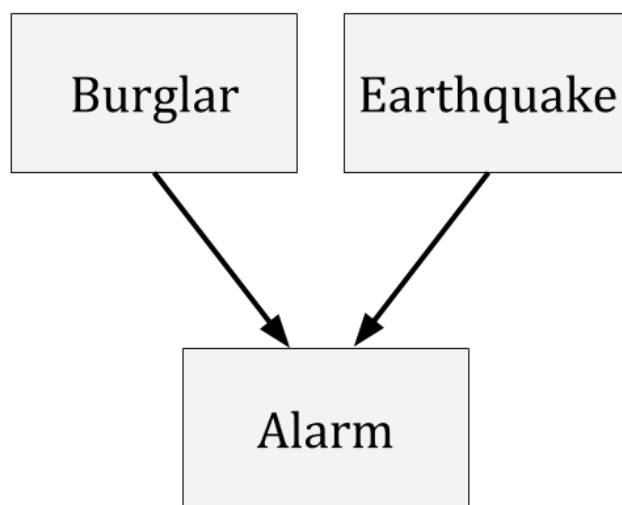


Figure 2.2: The famous Burglar, Earthquake, Alarm network showing a minimal case of explaining away.

In the network shown in figure 2.2, knowledge of the Alarm creates a dependence between Burglar and Earthquakes. For instance, say the Alarm has gone off and we know an earthquake has occurred, our belief in being burgled decreases. The dependence in belief networks means that sampling from the network requires a longer Markov Chain to mix, as changing the value of Earthquake, effects the value of Burglar. [TODO WORDING In a network with many connected nodes the dependence introduced makes sampling take longer. In the context of images, where there may be upwards of 1000 observable values, all with different dependencies this is intractable.](#)

### 2.4.3 Boltzmann Machines

A Boltzmann machine **TODO CITE: Cite the Harmonium, markov field** has qualities in common with Belief Networks. Both are generative models with their nodes having probabilities of being active based on neighboring nodes. Connections between nodes have associated weights as shown in figure 2.3. These weights are symmetric. **TODO WORDING** Unlike a Belief Network, a Boltzmann Machine is a undirected network that allows cycles and thus more complex data can be captured.

**TODO WORDING** Connections between nodes no longer encode causal information, instead a depedancy, the difference being that a connection encodes a relationship as they are not directed.

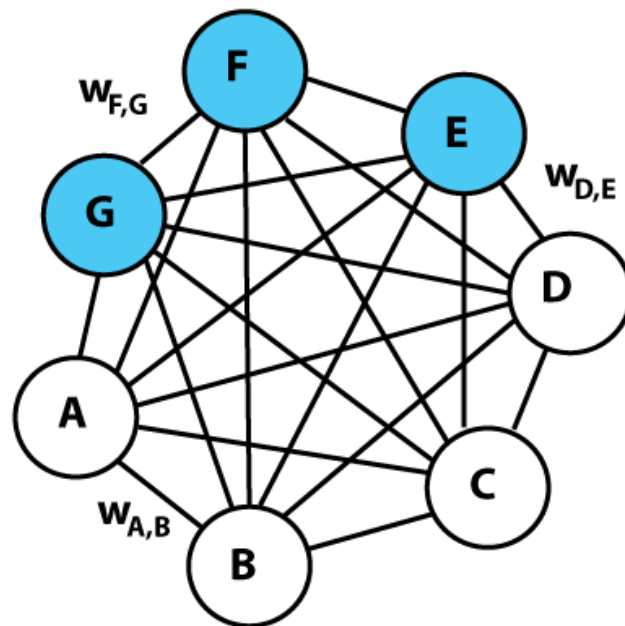


Figure 2.3: A Boltzmann Machine, the blue shaded nodes representing the observed variables, and the non-shaded nodes the latent variables.

Performing gibbs sampling appears trivial in a Boltzmann Machine, in that to find the probability of a given unit being active a weighted input to that node is passed through a sigmoid function. However, in practice the recurrent nature of Boltzmann Machines makes sampling intractable.

**TODO CITE: TODO-REFERENCE-PAPER-OF-THIS** The Boltzmann Machine was shown, given an unreasonable amount of time, to be able to perform better than the state of the art at the time.

## 2.5 The Current Approach: A Strong assumption

### 2.5.1 Restricted Boltzmann Machines

While Boltzmann Machines are impractical to train and sample from as networks grow in size **TODO CITE: Need to cite this...** their architecture can be altered to alieviate these shortcomings. The restriction, proposed by **TODO CITE: Hinton, a proper cite** requires the

network to be bipartite, where connections are forbidden between the layer of hidden units and the layer of visible units respectively.

An example Restricted Boltzmann Machine architecture is shown in figure 2.4. The affect of the restriction is that inference can be tractably computed, as the latent variables no longer become dependant given the observed variables. For example in figure 2.4 the hidden unit  $h_1$  is not dependant on  $h_2$  whether or not we know anything about the visible units. This is the opposite of a Sigmoid Belief Network where knowledge of the visible units makes the hidden units dependant. The RBMs nature removes the recurrence present in Boltzmann Machines. This reduces the expresiveness of the network but makes the RBM useable in practice. **TODO CITE: The paper about Boltzmann Machines being really good when actually left to find solution.**

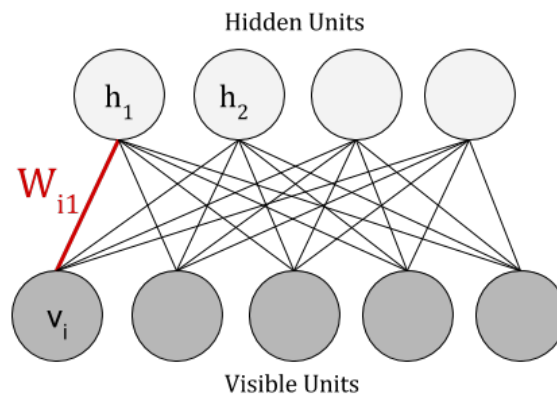


Figure 2.4: A Boltzmann Machine, the blue shaded nodes representing the observed variables, and the non-shaded nodes the latent variables.

**TODO Draft Points Still to write TODO**

Talk about how Gibbs sampling in RBMs allows us to approximately sample from the hidden (unknown/representation) given the visible (known/input data).

**TODO Draft Points Still to write TODO**

Hinton **TODO-REFERENCE-THE-PAPER** proposed a restriction by way of assumption to the Boltzmann Machine that makes it tractable to sample from and therefore train. Boltzmann Machines of this architecture are referred to as Restricted Boltzmann Machines, or RBMs for short.

The assumption being that the observable and latent variables are independant respectively, enforcing a two layer, fully connected bipartide structure.

### Tractable Training - Contrastive Divergence

Hinton **TODO-CITE-CLASSIC-PAPER** proposed Contrastive Divergence as a method for training RBMs efficiently. The algorithm leverages the now tractable wake phase because  $P(h|v)$  is efficeint to compute. However the free or sleep phase required another restriction where the network is only left to its own dynamics can be limited to only one iteration and still perform well. **TODO-CITE-CD-PAPER**

The observed variables are often referred to as the visible units, and will be so forth in this report. The latent variables are often referred to as the hidden units, and will be so forth in this report. Therefore the Restricted Boltzmann machine transforms some visible unit into a hidden representation. These two layers of units can be thought of as vectors of binary values, referred to as  $v$  and  $h$  for visible and hidden layers respectively.



This restriction allows an efficient calculation of the Wake Phase of generative model learning, as the  $P(h|v)$  can be calculated as a simple weighted sum passed through a sigmoid followed by a bernouli trial where the probability of being 1 is equal to the result of sigmoid.

### 2.5.2 Deep Learning

- Discuss deep learning as there are clear parallels to Deep Belief Networks and the new approach
- in particular how the deep networks have this process of freezing the weights and creating a sigmoid belief layer instead. There seem to be clear parallels between a deep network with one RBM to the ORBM.
- Unrolling the gibbs chain and we are in effect training an infinite depth sigmoid belief net (TODO-REFERENCE-HINTONS-PAPER-HERE)

### 2.5.3 Inference

One of the reasons the Restricted Boltzmann Machine is effective in practice is inference can be performed efficiently. Inference being computing the posterior.

### 2.5.4 Evaluating Restricted Boltzmann Machines

- Being unsupervised makes it difficult to evaluate RBMs. Often used as part of a deeper network, feature extractor, autoencoder
- Hinton Diagrams allow visualisation of hidden unit utilisation (TODO-SOME-SORT-OF-CITE). The weights out of a given hidden unit can be visualised in visible data space. The weights should exhibit some structure if they are being utilised. This is a good smoke tests for non-utilised hidden units will look very similar to units with random initial weights.
  - Small Cases
    - \* In trivial cases an RBM can inspected analytically. Reconstructions of the dataset should match the dataset with approximately the correct proportion. For instance training RBMs on 2 bit XOR should result in mostly [1,0] and [0,1] but not [1,1] and [0,0].
    - \* Hand craft weights can be used to perform inference in a 'perfect model'. For instance an RBM that can capture two bit, logical XOR can be represented as :TODO-INSERT-PIC
  - Large Cases
    - \* In non-trivial cases, with larger datasets, reconstructions can be compared to the dataset but given the unsupervised nature of RBMs empirically detecting if a model is trained is difficult.
    - \* The log likelihood of the RBMs generative model exhibiting the dataset is a good measure that can be approximated (because we have to sample).
    - \* We can train a classifier on the RBMs hidden representation. This can be compared for a ORBM and RBM.

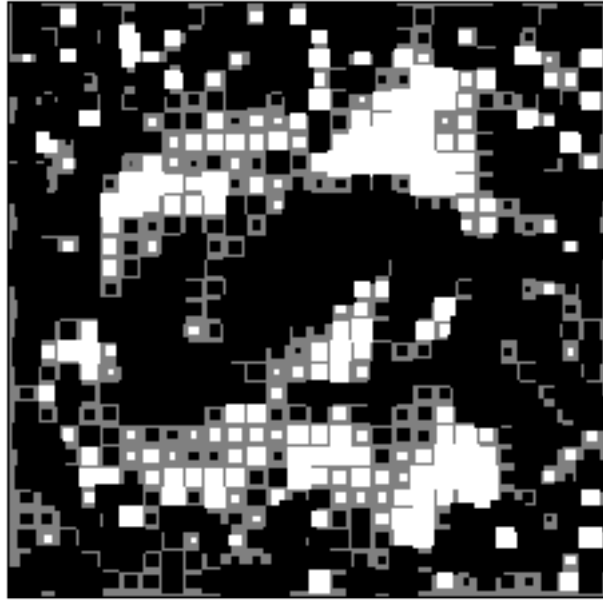


Figure 2.5: Good Hinton Diagram

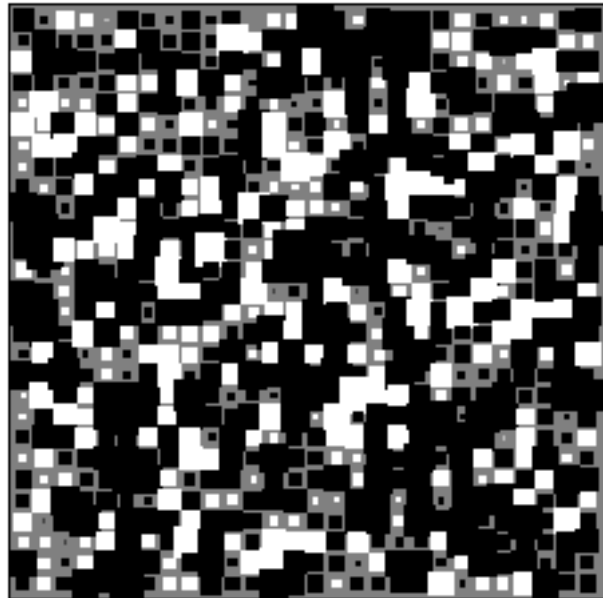


Figure 2.6: Bad Hinton Diagram

## 2.6 A New Approach - The ORBM

- Frean and Masland's new approach combines the Restricted Boltzmann Machine and the Sigmoid Belief Network, the RBMs allowing the rich complex causes to be encoded independently, and the Sigmoid Belief Network modelling the combination of the causes to form the observable data.'
- By building on these existing methods can leverage existing algorithms
- Can verify the inference algorithm with pre-trained RBMs for each cause.
- Like the RBM leveraged in the ORBM, it is difficult to evaluate, But similar techniques can be leveraged

### 2.6.1 Architecture

- Two RBMs, one for each cause, then they combine via a Sigmoid Belief Layer. Weights between the RBM and Sigmoid Layers are shared. (TODO-A-DIAGRAM)
- Diagram of the ORBM Architecture Including U Layer. Make sure I'm explaining the U layer.
- In fact  $U_a, U_b$  are like mirrors of the visible.

### 2.6.2 Inference In the ORBM

The difference in architecture from an RBM means that a slightly different inference algorithm is required, as the representations are represented separately for a composite input.

- Inference in this generative model is  $P(h_a, h_b \text{ given } v)$ , we want to represent the causes separately.
- Unfortunately,  $h_a$  and  $h_b$  are dependant given the visible  $v$ . meaning to perform inference, obtaining a hidden representation requires a Gibbs chain
- Diagram showing the inference gibbs chain.

### Calculating the Posterior

To find the  $P(h_a, h_b | v_{comp})$ , where  $h_a$  and  $h_b$  are the separate representations of the data caused  $a$  and  $b$ , and  $v_{comp}$  is the composed/composite input. We must sample from a Gibbs Chain, as  $h_a$  and  $h_b$  are dependant given  $v_{comp}$ . This ends up being almost identical to the RBM except we add a 'Correction', when computing the update for  $h_a$  and  $h_b$  respectively.

TODO-SHOW-THE-FULL-CORRECTION

### 2.6.3 Source Separation - Reconstructions in the ORBM

To actually perform source separation, one needs to simply take the internal representation generated by the inference step, and generate a visible pattern in the same way you would with a standalone RBM.



# Bibliography

[1] History of SciPy. [http://wiki.scipy.org/History\\_of\\_SciPy/](http://wiki.scipy.org/History_of_SciPy/). Accessed: 2015-07-20.