

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 [TODO WORDING I or We????](#)



# 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	A Problem . . . . .	1
1.2	Deep Belief Networks can achieve state of the art performance . . . . .	1
1.3	A Proposed Solution and Contributions . . . . .	1
<b>2</b>	<b>Backgroud</b>	<b>3</b>
2.1	Generative Models and Source Separation . . . . .	3
2.1.1	Terminology in Generative Models observable and hidden variables .	3
2.1.2	PGMs as a tool reasoning about generative models . . . . .	4
2.2	Sampling and inverting the model . . . . .	4
2.2.1	Gibbs sampling, a subset of Markov Chain Monte Carlo . . . . .	4
2.2.2	Mixing Time . . . . .	4
2.2.3	Reconstructions, visualising what the model has learnt . . . . .	5
2.2.4	Ancestral Sampling or fanstasies of the model . . . . .	5
2.3	An intractable model for causes . . . . .	5
2.3.1	Sigmoid Belief Networks . . . . .	5
2.3.2	Explaining Away — a trade off between modeling causes and tractability	6
2.3.3	Boltzmann Machines . . . . .	6
2.4	Restricted Boltzmann Machines: A Strong assumption . . . . .	7
2.4.1	Energy, and the log likelihood of the joint . . . . .	7
2.4.2	Gibbs Sampling in RBMs . . . . .	8
2.4.3	Evaluating Restricted Boltzmann Machines . . . . .	10
2.5	A network equivalent to an RBM . . . . .	12
2.5.1	Unrolling a Gibbs Chain, a different perspective on RBM inference . .	12
2.6	A New Approach, The ORBM . . . . .	14
2.6.1	Architecture . . . . .	14
2.6.2	Gibbs Sampling in the ORBM . . . . .	14
<b>3</b>	<b>Design and Implementation</b>	<b>19</b>
3.1	Evaluation Dataset Choice . . . . .	19
3.2	Implementation Design . . . . .	20
3.2.1	Language Choice . . . . .	20
3.2.2	Program Architecture . . . . .	20
<b>4</b>	<b>Evaluation</b>	<b>23</b>
4.1	Evaluation Design . . . . .	23
4.1.1	Mixing Time Evaluation . . . . .	23
4.1.2	Reconstructions As a measure of Performance . . . . .	23
4.1.3	Growing Complexity of tasks . . . . .	23
4.2	Evaluation Methods . . . . .	23

4.3	Starting from the minimal case, two bit XOR . . . . .	24
4.4	Y bit pattern, X neighbouring bits on . . . . .	27
4.5	2D Patterns in a small dimensional space . . . . .	28
4.6	2D Pattern, Different Rectangle Separation . . . . .	29
4.7	MNIST Digits Reconstructions . . . . .	29
4.7.1	MNIST Analysis . . . . .	30
<b>5</b>	<b>Conclusions and Future Work</b>	<b>33</b>

# Figures

1.1	An example composition of a handwritten 4 and 3, illustrating a non-trivial task where the ground truth is known. . . . .	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	An example Restricted Boltzmann Machine with four hidden units, and five visible units. Note that the edges between units are not directed - representing a dependency not a cause. . . . .	8
2.5	A figure illustrating a Gibbs chain where left to right indicates a Gibbs iteration. Note this is <i>not</i> a PGM. . . . .	9
2.6	A diagram showing $\psi_j$ , the weighted sum into the $j$ th hidden unit. Note that $W_{0j}$ is the hidden bias, represented as a unit that is always on with a weight into each hidden unit. . . . .	10
2.7	Hinton Diagrams illustrating a trained hidden unit, versus that of an untrained/unutilised hidden unit. This is with no measures to enforce sparsity. . . . .	10
2.8	A diagram illustrating 'unrolling' an RBM by one Gibbs iteration. Note the connections between $H$ and $V$ layers are now directed. . . . .	12
2.9	A diagram showing Ancestral sampling in the equivalent network, where normal sampling in the top 2 layers is performed until Gibbs iteration $t$ the hidden state is pushed through the bottom layers Sigmoid Network. . . . .	13
2.10	The full ORBM architecture, where $A$ and $B$ are the two causes that combine to form the data. . . . .	15
2.11	A diagram that shows the structure of the correction over the weighted sums into the hiddens of one of the RBMS versus both. . . . .	16
3.1	A figure showing the architecture for the implemented test suite in UML notation. . . . .	21
4.1	A figure illustrating two five by five pixel images combining to form a composite/multicause input. The images are binary. . . . .	23
4.2	The two bit RBM for modelling a single bit on at a time in a two bit pattern. The diagonal weights are negative and the non-diagonal weights are positive. . . . .	24
4.3	Example reconstructions for the Handcrafted RBM, For 10000 independant reconstructions given the input $[1, 0]$ . . . . .	24
4.4	Dreams in the XOR RBM, note how only the training data is present. . . . .	25
4.6	The process of generating reconstruction is shown in this diagram. . . . .	26
4.5	Figure showing how the XOR RBM fits into the ORBM architecture. . . . .	26

4.7	A figure with the results of running ORBM inference 1000 times, with the two different approaches to calculating the correction. The frequency for which a given reconstruction occurred is shown. . . . .	27
4.8	A figure showing the result of generating 1000 reconstructions with the ORBM with $[1,0]$ as input, again comparing the Approximated correction and Full correction. . . . .	27
4.9	A figure showing the result of generating 1000 reconstructions with an RBM. It is clear the RBM has no mechanism to separate the sources. Hence showing it $[1,1]$ when it has only been trained on XOR results in no separation, i.e. reconstructions including $[1,1]$ . . . . .	28
4.10	Dataset-wide MNIST Score Results . . . . .	31
4.11	Cosine Score breakdown for the highest and lowest performing datasets, 0 and 9. . . . .	31



# Chapter 1

## Introduction

### 1.1 A Problem

Consider an image of a face. At least two systems are at play in the image of a face, the face itself and the illumination. Both face and illumination are complex and can vary greatly, but they are fundamentally acting independently of each other up until they compose to form the image. A form of Deep Neural Networks, Deep Belief Networks, have achieved state of the art performance in facial recognition, but this is only possible with a large amount of training data. This suggests that these networks are missing a way to represent these *sources* independently. This project takes steps toward representing complex causes separately with the primary task of decomposing joint images (data).

Separating faces and illumination is too challenging for this project, there is no way to verify/evaluate that the new approach is working as the concept of a face without illumination cannot be visualised. Instead I will start with the modest task of working with images where the source images being combined are known.

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

Deep Belief networks (DBNs) are powerful models that have proven to achieve state of the art performance in tasks such as image classification, dimensionality reduction, natural language recognition, Document classification, Semantic Analysis. **TODO CITE:** Despite a DBN's expressiveness, there is no way to extract independent sources, the model instead learns how to represent the complex combination. The complex combination of sources is inherently richer than the individual sources acting alone. The DBN may learn features that correspond to each source during its training process, however the architecture or training algorithm make no attempt to enforce this.

DBNs are built of shallow networks called Restricted Boltzmann Machines.

### 1.3 A Proposed Solution and Contributions

Frean and Marsland propose an alternative architecture to that of an RBM, that aims to encode two complex sources independently. Frean and Marsland also propose an algorithm to put this alternative architecture to use.

This project contributes:

- The first articulation of the architecture, presenting the context needed to understand the new architecture.

- A suite of graded evaluations/tests that explore how the architecture and source separation algorithm work in practice.

The evaluations will not address separating faces and illumination, instead only performing tasks such as separating two handwritten digits composed on each other (see figure 1.1).

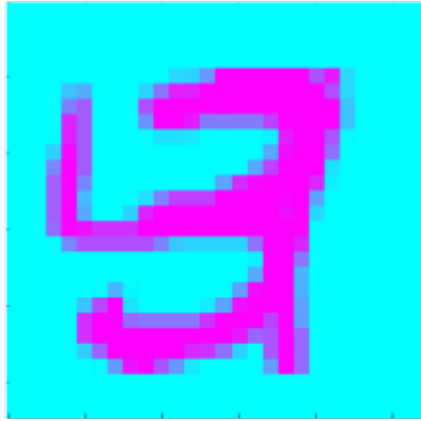


Figure 1.1: An example composition of a handwritten 4 and 3, illustrating a non-trivial task where the ground truth is known.

## 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 Generative Models and Source Separation

Generative models are a powerful way to model data. The rational behind them being that we aim to learn a model that can both create the training data and represent it. Input data, say images, could be mapped from raw values into some higher level features. Hinton gave a compelling argument why higher level features are desirable in the context of generative models [2].

Consider, for example, a set of images of a dog. Latent variables such as the position, size, shape and color of the dog are a good way of explaining the complicated, higher-order correlations between the individual pixel intensities, and some of these latent variables are very good predictors of the class label. In cases like this, it makes sense to start by using unsupervised learning to discover latent variables (i.e. features) that model the structure in the ensemble of training images.

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

#### 2.1.1 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.

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

Probabilistic 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 dependency or mapping.

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. Throughout this report, RBMs, Sigmoid Belief Networks and the proposed ORBM will be shown in this format.

## 2.2 Sampling and inverting the model

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

Inverting a generative model can be referred to as inference, the process of reasoning about what we do not know, given that of which we do know. In generative models this is the 'posterior', the probability distribution of the hidden (latent) variables, given we know the state of the observable (visible) units. The algorithm introduced and evaluated in this report aimed to invert the ORBM, generating two separate hidden representations given an input.

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

Gibbs sampling is a special case of Markov Chain Monte Carlo, a technique for drawing sampling from a complex distribution. Sampling from the probability mass (or '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 distribution's 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 or a Gibbs Step.

Gibbs sampling is used for performing inference in RBMs and as result also 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.

### 2.2.2 Mixing Time

MCMC methods aim to approximate a distribution, by exploring likely states. As we often start this process from a random state, it's important that enough Gibbs steps are taken before a sample is drawn. This is because the random state may not be close any part of the true distribution we want to sample from, so by running the chain for many iterations we increase the likelihood of samples being from the desired distribution.

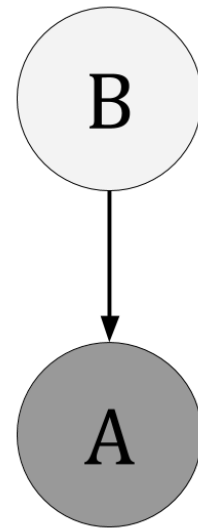


Figure 2.1: An example PGM, showing an observed variable 'A' and it's hidden cause 'B'.

This process of waiting for enough steps to before drawing samples is referred to as the Mixing Time. When Hinton when proposed a fast training for algorithm for RBMs and DBNs [3], one Gibbs step was sufficient in practice for training and using an RBM. The ORBM is not so fortunate.

### 2.2.3 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.

### 2.2.4 Ancestral Sampling or fantasies of the model

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: (TODO-CITE-PAPER-WITH-MNIST-DREAM-EVALUATION, they were crappy).**

## 2.3 An intractable model for causes

### 2.3.1 Sigmoid Belief Networks

The ORBM relies on the Sigmoid Belief Network to capture the causation. The Sigmoid Belief Network (SBN) 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.

SBNs offer a way to model causation in machine learning, as every hidden unit represent a single, simple cause. 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 simple cause. Meaning if a cause's 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 independent causes gave rise to a given pixel.

Despite the Sigmoid Belief Network being expressive and providing a succinct encoding of inter-variable dependancies. However, there is too much dependency in it's expressiveness, to the extent that performing inference is intractable for a network of practical size. There do exist algorithms for performing inference in Sigmoid Belief Networks. For instance, the Belief Propagation algorithm proposed by Judea Pearl [8] 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. This intractability arises from the SBNs dependencies and the 'explaining away effect' as Hinton touches on

in light of training stacked RBMs (DBNs) [3]. Inference is required for training generative models making Sigmoid Belief Networks impractical to train.

### 2.3.2 Explaining Away — a trade off between modeling causes and tractability

The benefit of the Belief Network is also it's weakness, a 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 David Barber in his book [1], 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 conveys causation.

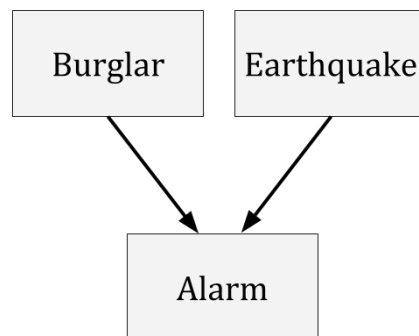


Figure 2.2: The famous Burglar, Earthquake, Alarm network showing a minimal case of explaining away.

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. 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 dependancies this becomes intractable. Neal showed this by comparing the number of gibbs iterations required for small enough error rates in [7].

### 2.3.3 Boltzmann Machines

A Boltzmann machine has proposed in various forms throughout the years from different backgrounds, for instance Paul Smolensky presented it in context of 'Harmony Theory' [9]. Where-as Hinton and Sejnowski also propped the Boltzmann machine in [4]. The Boltzmann Machine 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. Unlike a Belief Network, a Boltzmann Machine is a undirected network that allows cycles and thus more complex relationships can be captured. Also connections no longer encoe causal information, instead they encode a dependency.

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

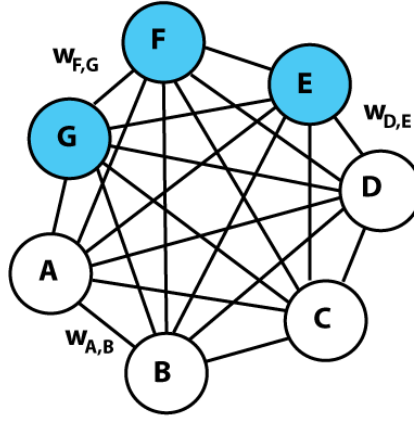


Figure 2.3: A Boltzmann Machine, the blue shaded nodes representing the observed variables, and the non-shaded nodes the latent variables.

sampling intractable as updating a node will change the probabilities of those connected. However, it was shown that given unlimited training time Boltzmann Machines could be trained, out performing the state of the art models of the time [TODO CITE: This](#).

## 2.4 Restricted Boltzmann Machines: A Strong assumption

While Boltzmann Machines are impractical to train and sample from as networks grow in size [TODO CITE](#): their architecture can be altered to alleviate these shortcomings. The restriction, originally proposed by [TODO CITE: Hinton, a proper cite](#), and then later revitalised with a deeper architecture training algorithm in Hinton [4]. The restriction requires the network to be a two layer bipartite network, each layer corresponding to the observed (visible) and latent (hidden) units. 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 collection of hidden units, forming a layer are referred to as the *hidden layer*. The collection of visible units are referred to as the *visible layer*.

The effect of this restriction is inference becomes tractable, as the latent variables no longer become dependant given the observed variables. This is illustrated 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. By removing the recurrence present in Boltzmann Machines, it reduces the expresiveness of the RBM network but makes the RBM useable in practice.

### 2.4.1 Energy, and the log likelihood of the joint

An RBM models the joint distribution of hidden and visible states. The RBM assigns to every configuration of  $\tilde{h}$  and  $\tilde{v}$  an Energy, where the lower the energy, the more likely the RBMs configuration is to *fall* into that state. Hopfield, in the context of what is now called the Boltzmann Machine [6], presented this energy as defined by the function as

$$E(\tilde{v}, \tilde{h}) = - \sum_{i \in \text{visible}} W_{0i} v_i - \sum_{j \in \text{hidden}} W_{0j} h_j - \sum_{i,j} v_i h_j W_{ji}$$

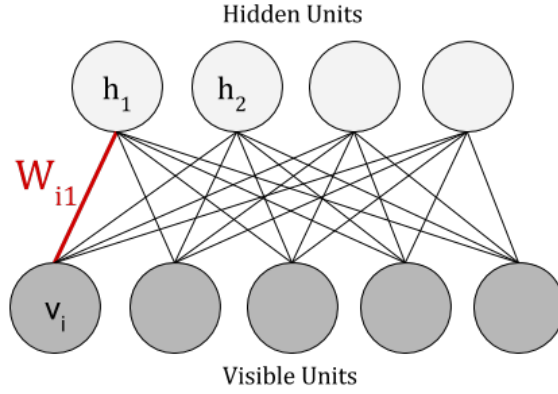


Figure 2.4: An example Restricted Boltzmann Machine with four hidden units, and five visible units. Note that the edges between units are not directed - representing a dependency not a cause.

The probability of the RBM being in a given configuration is the joint probability of  $\tilde{h}$  and  $\tilde{v}$ .

$$P(\tilde{h}, \tilde{v}) = \frac{1}{Z} \prod_{j,i} e^{h_j W_{ji} v_i}$$

If we move to log space that becomes

$$\log P(\tilde{h}, \tilde{v}) = \frac{1}{Z} \log \sum_{j,i} h_j W_{ji} v_i$$

$\frac{1}{Z}$  is the partition function, which normalises the probability of the joint. Calculating this would require summing over all possible configurations of  $h$  and  $v$ , which is intractable for practical numbers of units. For instance a 28 by 28 image corresponds to 784 visible units, and for say 10 hidden units this would amount to  $2^{784} * 2^{10}$  possible configurations. We opt to work in terms of  $P^*$  which is the non-normalised probability of the joint over  $h$  and  $v$ . So we arrive at

$$\log P^*(h, v) = \sum_i \sum_j e^{h_j v_i W_{ji}} \quad (2.1)$$

## 2.4.2 Gibbs Sampling in RBMs

In order to describe Gibbs sampling in the ORBM, it must be first introduced in the context of RBMs. For the following discussion, the hidden units are indexed by  $j$  and the visible units are indexed by  $i$ . To perform inference, create reconstructions, and train RBMs, Gibbs sampling is used to draw samples from  $P(h|v)$  (referred to as sampling from the posterior) and  $P(v|h)$  respectively. It is introduced here to be reference later in Gibbs in an ORBM. In an RBM the process of Gibbs sampling is as follows:

- One must sample from  $P(\tilde{h}|\tilde{v})$  giving a hidden state  $\tilde{h}'$
- Using this hidden state, a visible state is then generated,  $\tilde{v}'$ , by sampling from  $P(\tilde{v}'|\tilde{h}')$ . This process of generating a hidden pattern, and subsequent visible pattern is referred to as a Gibbs step.



- This chain of Gibbs steps between sampling from  $P(h|v)$  and  $P(v|h)$  can then be repeated as desired, the longer the chain the closer the samples will be to the true joint distribution that the model has learnt. For training an RBM Hinton showed that 1 step is often enough in practice, as one step is enough to infer a direction to adjust the weights in.

This process can be calculated in one step for all units in the target layer as they are independent of one another. The process of updating the hidden, then visible layers forms what is referred to as the Gibbs Chain. **TODO CITE: Feel like I could cite this again in like a CD paper or something.** This process is visualised at a layer level in figure 2.5

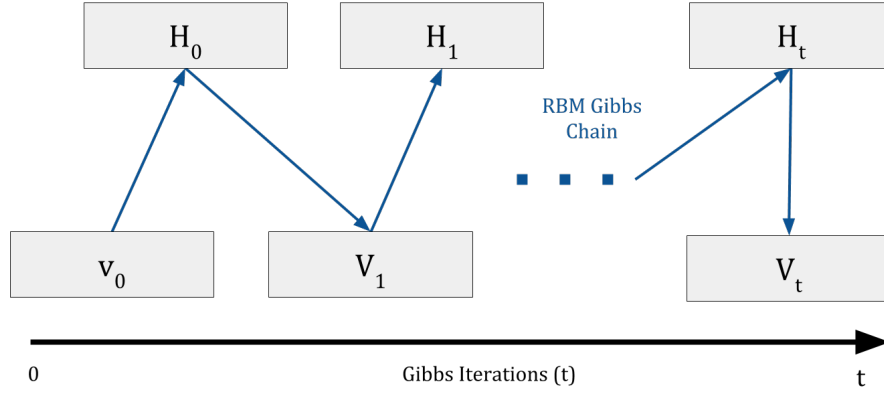


Figure 2.5: A figure illustrating a Gibbs chain where left to right indicates a Gibbs iteration. Note this is *not* a PGM.

### The Gibbs update

In a standard RBM, updating a hidden unit  $h_j$  when performing Gibbs sampling is calculated by finding  $P(h_j = 1|\tilde{v})$  where  $\tilde{v}$  is an input pattern. In the context of an image,  $\tilde{v}$  would be the pixel values where each pixel corresponds to a visible unit,  $v_i$ . The probability of a given hidden unit activating is: **TODO CITE: Gibbs sampling would be good here.**

$$P(h_j = 1|\tilde{v}) = \sigma(\psi_j) \quad (2.2)$$

Where  $\psi_j$  is the weighted sum into the  $j$ th hidden unit and  $\sigma()$  is the Sigmoid function, or it also known as the Logistic function  $\sigma(x) = 1/(1 + e^{-x})$ . Figure 2.6 illustrates  $\psi_j$  for an example RBM.

As the weights are symmetric, sampling from the visible layer, given a hidden state is similar. That is  $P(v_i = 1|\tilde{h})$ , where  $\tilde{h}$  is the entire hidden vector is given by:

$$P(v_i = 1|\tilde{h}) = \sigma(\phi_i) \quad (2.3)$$

Where  $\phi_i$  is the weighted sum into the  $i$ th visible unit, which is:  $\phi_i = \sum(W_{ji}h_j) + W_{0i}$ . Both  $\phi_j$  and  $\psi_i$  can be expressed in alternative, but useful way:

$$\phi_j = \log P^*(v, h|v_i = 1) - \log P^*(v, h|v_i = 0) \quad (2.4)$$

$$\psi_i = \log P^*(h, v|h_j = 1) - \log P^*(h, v|h_j = 0) \quad (2.5)$$

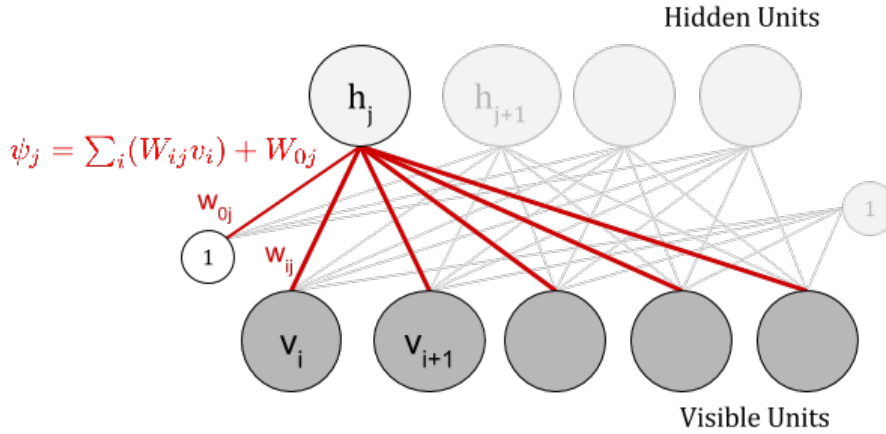
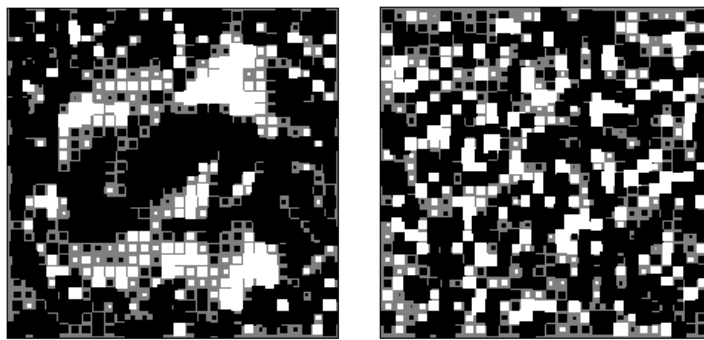


Figure 2.6: A diagram showing  $\psi_j$ , the weighted sum into the  $j$ th hidden unit. Note that  $W_{0j}$  is the hidden bias, represented as a unit that is always on with a weight into each hidden unit.

### 2.4.3 Evaluating Restricted Boltzmann Machines

A challenge faced by this project and work with RBMs is that they are non-trivial to evaluate **TODO CITE: Something please**. Being unsupervised black box, one cannot inspect the hidden units for an input and be sure that a good model has been learnt.

We can examine the weights into a given hidden unit in the shape of the visible vector. For instance in the context of images, Hinton Diagrams allow visualisation of what a given hidden unit is ‘doing’ by visualising the matrix. This was first used by Hinton in the context of Boltzmann Machines in [5]. They also give insight into hidden unit utilisation. Weights are often initialised to small random values, the Hinton diagram will have no visual structure to it. This is illustrated in figures 2.7a and 2.7b. The former showing a hidden units weights where some structure has been learnt, and the latter showing the opposite.



(a) Utilised hidden unit Hinton Diagram (b) Unutilised hidden unit Hinton Diagram

Figure 2.7: Hinton Diagrams illustrating a trained hidden unit, versus that of an untrained/unutilised hidden unit. This is with no measures to enforce sparsity.

Evaluating RBMs is also problem dependant, if split by dimensionality of the task we have some different approaches:

**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 two bit XOR should result in mostly  $[1,0]$  and  $[0,1]$  but not  $[1,1]$  and  $[0,0]$ . If the task has a small enough dimensionality, like two bit XOR, then a learned RBM can be compared to one constructed by hand. This approach is used later when the ORBM tackles two bit XOR.

**Large Cases** In non-trivial cases, with larger datasets, reconstructions can be compared to the training dataset, but empirically detecting if a model is trained is difficult given the unsupervised, black box nature of RBMs. We can train a classifier on the RBMs hidden representation, with the idea being that better hidden features should result in a better way of distinguishing classes, and therefore a better classification accuracy.

## 2.5 A network equivalent to an RBM

### 2.5.1 Unrolling a Gibbs Chain, a different perspective on RBM inference

Before the ORBMs architecture and inference algorithm can be introduced, Gibbs sampling in an RBM must be presented in a different, yet equivalent way. Hinton, when introducing the idea of training a Deep Belief Network showed that a Gibbs chain in an RBM is equivalent to a infinitely deep belief network, with an RBM on the top with tied weights. **TODO CITE: Hinton's paper on unrolling the Gibbs chain.** An unrolling of a single Gibbs iteration is illustrated in figure 2.8, the  $U$  layer corresponding to the  $V$  layer after one Gibbs iteration. The top two layers ( $U$  and  $H$ ) form a standard RBM, but the bottom connections ( $V$  and  $H$ ) form a Sigmoid Belief Network. To further clarify, the  $U$  layer corresponds to  $V_1$  in figure 2.5. Note in figure 2.8 the weights between the  $H$  and  $U$  layer are shared between the  $V$  and  $H$  layers.

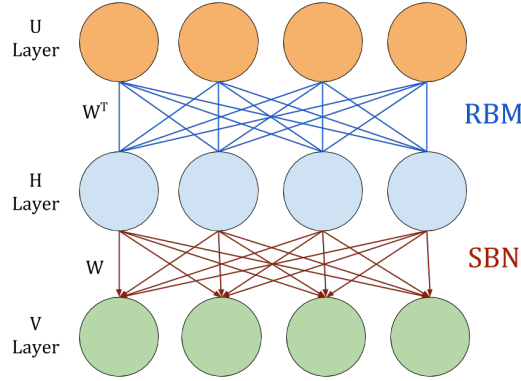


Figure 2.8: A diagram illustrating ‘unrolling’ an RBM by one Gibbs iteration. Note the connections between  $H$  and  $V$  layers are now directed.

We can now show that Gibbs sampling in this RBM unrolled with a Sigmoid belief network is equivalent to Gibbs Sampling in a standard RBM.

#### Ancestral Sampling in this equivalent network

Ancestral sampling in an RBM **TODO CITE: as described in the section where I talk about dreams....** This network behaves a similar way, where a Gibbs chain is run between the top two layers,  $U$  and  $H$ , until the last iteration. At the last iteration a hidden pattern is sampled and then is pushed through the Sigmoid belief network between  $H$  and  $V$ . This is illustrated in figure 2.9.

To generate a sample from  $h$  is equivalent we can use equation 2.2, by running the Gibbs chain between  $h$  and  $U$ . To sample from  $v$  in the SBN is by definition:

$$P(v_i = 1|h) = \sigma_i(h)$$

Thus Gibbs sampling from a simple RBM ending in a sample for  $v$  is the same as sampling  $h$  from the same RBM and using a SBN for the last step.

There is however another useful way to draw samples in such a network that we will leverage in the ORBM. The log likelihood of the joint can be written as:

$$\log P^*(h, v) = \log P^*(h) + \log P(v|h) \quad (2.6)$$

We have the second term already:

$$\log P(v|h) = \sum_i v_i \log \sigma_i(h) + (1 - v_i) \log(1 - \sigma_i(h))$$

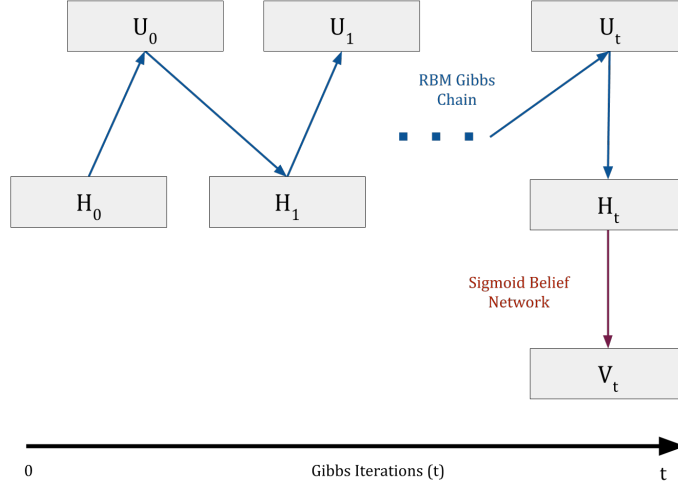


Figure 2.9: A diagram showing Ancestral sampling in the equivalent network, where normal sampling in the top 2 layers is performed until Gibbs iteration  $t$  the hidden state is pushed through the bottom layers Sigmoid Network.

To find the remaining term,  $P^*(h)$  we need to marginalise the joint (eq 2.6) over all  $\mathbf{v}$  configurations:

$$\begin{aligned}
 P^*(h) &= \sum_{v_1=0}^1 \cdots \sum_{v_n=0}^1 \exp \left[ \log P^*(h, v) \right] \\
 &= \sum_{v_1=0}^1 \cdots \sum_{v_n=0}^1 \exp \left[ \sum_i \sum_j h_j W_{ji} v_i + \sum_i W_{0i} v_i + \sum_j W_{j0} h_j \right] \\
 &= \sum_{v_1=0}^1 \cdots \sum_{v_n=0}^1 \exp \left[ \sum_i v_i \phi_i(h) + \sum_j W_{j0} h_j \right] \\
 \text{where } \phi_i(h) &= \sum_j W_{ji} h_j + W_{0i} \\
 &= \exp \left[ \sum_j h_j W_{j0} \right] \times \sum_{v_1=0}^1 \cdots \sum_{v_n=0}^1 \prod_i \exp \left[ v_i \phi_i(h) \right] \\
 &= \exp \left[ \sum_j h_j W_{j0} \right] \times \prod_i \left( 1 + e^{\phi_i(h)} \right) \\
 \text{and so } \log P^*(h) &= \sum_j h_j W_{j0} + \sum_i \log \left( 1 + e^{\phi_i(h)} \right) \\
 &= \sum_j h_j W_{j0} + \sum_i \phi_i(h) - \sum_i \log \sigma_i(h)
 \end{aligned}$$

So far we've figured out  $\log P^*(h)$  for the RBM that is the 'top layer'.

Therefore another way to write  $\log P^*(h, v)$  is

$$\log P^*(h, v) = \underbrace{\sum_j h_j W_{j0} + \sum_i \phi_i(h) - \sum_i \log \sigma_i(h)}_{\log P^*(h)} + \underbrace{\sum_i v_i \log \sigma_i(h) + (1 - v_i) \log(1 - \sigma_i(h))}_{\log P(v|h)}$$

By collecting terms and simplifying one can readily see that this matches the earlier form in equation 2.1.

## 2.6 A New Approach, The ORBM

### 2.6.1 Architecture

Frean and Marsland extend on this idea of a one Gibbs iteration unrolled RBM. They propose adding another  $U$  and  $H$  layers to represent another rich source of data, which then combines with the original  $U$  and  $H$  layer via a SBN to form  $V$ . This architecture is illustrated in figure 2.10, where A and B are used to denote the two independent causes we are modeling. To clarify the ORBMs architecture, there are two RBMs, one for each cause. These RBMs are connected to a SBN into the visible layer, the weights of which are tied to the respective RBMs. This architecture is simpler in practice as the  $U$  layers can be captured in the inference algorithm.

By building on RBMs and SBN we can leverage existing algorithms and work on RBMs and also the potential for extending this work to deeper architectures is possible. Also the architecture supports ‘plug and play’ with the models in that existing RBMs can be plugged into the architecture. The implications of this are exciting in that an RBM that has been trained on hand written digits could be plugged into an RBM that has been trained on a paper texture (the bumps and ridges of paper). Then using the ORBM a ‘clean’ representation of the text and paper could be separated out given an image of a digit on paper.

The ORBM, like the RBM is difficult to evaluate empirically, but the same techniques that can be used to evaluate RBMs can be applied to the ORBM.

### 2.6.2 Gibbs Sampling in the ORBM

#### Ancestral Sampling in the ORBM

Ancestral sampling in an RBM involved running a Gibbs chain and then taking the last visible in that chain. To perform ancestral sampling in the ORBM we can leverage the generative power of the RBM and then combine their inputs in a simple way. We know that joint in the unrolled RBM is expressed as eq 2.6. That is, the ORBM probability of  $v$  given  $h^A$  and  $h^B$  is defined by:

$$\log P(v|h^A, h^B) = \sum_i v_i \log \sigma(\phi_i^A + \phi_i^B) + (1 - v_i) \log(1 - \sigma(\phi_i^A + \phi_i^B))$$

Where  $\phi_i^A$  and  $\phi_i^B$  are the weighted sums into the  $i$ th visible units from the models A and B respectively. In this generative model a visible unit is created by taking the weighted sum from both sources, adding their contribution and then passing through a Sigmoid function to give a probability.

#### Inference in the ORBM

In a standard RBM sampling from  $P(h_j)$  is given by  $P(h_j) = \sigma(\psi_j)$ , where  $\psi_j$  is defined in equation 2.5. In an ORBM we cannot consider the single RBM alone when trying to find

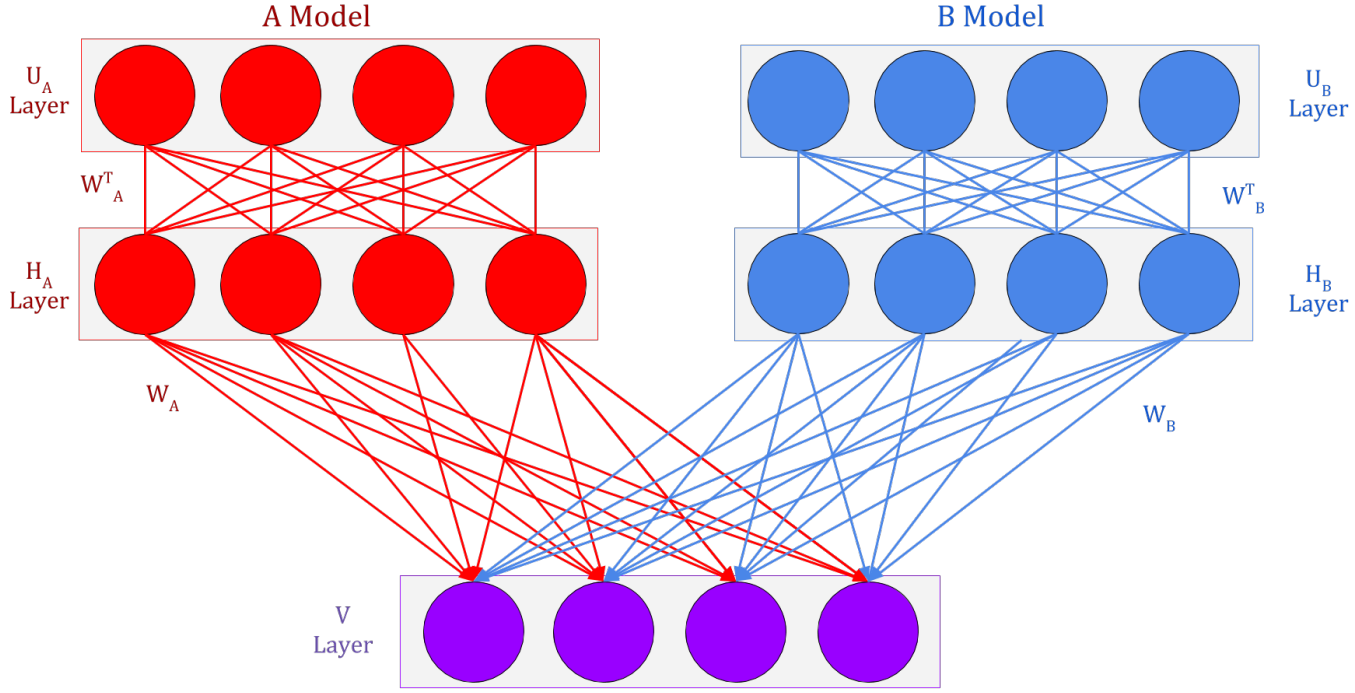


Figure 2.10: The full ORBM architecture, where  $A$  and  $B$  are the two causes that combine to form the data.

$P(h_j)$ , as the hidden layers of the two RBMs in the ORBM are dependent given a visible layer  $v$ . This amounts to explaining away as described in section 2.3.2. There is a nice feature of the ORBM in that the hidden representations ( $h^A$  and  $h^B$ ) we extract from the visible layer require no interaction with the  $U$  layers to sample from. They are only needed to generate a composite  $V$  pattern (or independent reconstructions).

We aim to use Gibbs sampling to generate  $h^A$  and  $h^B$  given a  $v$ : We need to calculate

$$\psi_j^A = \log P^*(h, v | h_j^A = 1) - \log P^*(h, v | h_j^A = 0)$$

We will use that fact that the weighted sum into a visible unit  $i$  where some hidden unit  $h_j$  is on, is equivalent to the same configuration except  $h_j$  is off plus the weight between these two units. This is by definition true, and expressed below:

$$\phi_i^A(h | h_j^A = 1) = \phi_i^A(h | h_j^A = 0) + W_{ji}^A$$

We will abbreviate  $\phi_i^A(h | h_j = 0)$  to  $\phi_i^{Aj0}$ . Given these we obtain:

$$\psi_j^A = \sum_i v_i \log \left( \frac{1 + e^{-\phi_i^{Aj0} - \phi_i^B}}{1 + e^{-\phi_i^{Aj0} - W_{ji} - \phi_i^B}} \frac{1 + e^{\phi_i^{Aj0} + W_{ji} + \phi_i^B}}{1 + e^{\phi_i^{Aj0} + \phi_i^B}} \right) + \sum_i \log \left( \frac{1 + e^{\phi_i^{Aj0} + W_{ji}}}{1 + e^{\phi_i^{Aj0}}} \frac{1 + e^{\phi_i^{Aj0} + \phi_i^B}}{1 + e^{\phi_i^{Aj0} + W_{ji} + \phi_i^B}} \right)$$

Now  $\phi = \log \frac{1+e^\phi}{1+e^{-\phi}}$  (Marcus' magic identity), which is  $= \log \frac{\sigma(\phi)}{\sigma(-\phi)}$ . So the first term simplifies to  $\sum_i v_i W_{ji}$ , which is the same as that in an RBM. The second term can also be simplified, using the identity  $\log(1 - \sigma(\phi)) = \phi - \log(1 + e^\phi)$ . This leads to the following Gibbs Sampler probability of the  $j$ -th hidden unit in network  $A$  being 1:  $p_j = \sigma(\psi_j^A)$  with

$$\psi_j^A = \underbrace{\sum_i W_{ji}^A v_i}_{\text{vanilla RBM}} + \underbrace{\sum_i C_{ji}^A}_{\text{correction}}$$

where the correction is:

$$\begin{aligned} C_{ji}^A &= \log \left[ \frac{\sigma(\phi_i^{Aj0})}{\sigma(\phi_i^{Aj0} + W_{ji}^A)} \cdot \frac{\sigma(\phi_i^{Aj0} + W_{ji}^A + \phi_i^B)}{\sigma(\phi_i^{Aj0} + \phi_i^B)} \right] \\ &= \log \sigma(\phi_i^{Aj0}) + \log \sigma(\phi_i^{Aj0} + W_{ji}^A + \phi_i^B) - \log \sigma(\phi_i^{Aj0} + W_{ji}^A) - \log \sigma(\phi_i^{Aj0} + \phi_i^B) \\ &= \log \left[ \frac{\sigma(\phi_i^A - h_i^A W_{ji}^A)}{\sigma(\phi_i^A + (1 - h_i^A) W_{ji}^A)} \right] - \log \left[ \frac{\sigma(\phi_i^{AB} - h_i^A W_{ji}^A)}{\sigma(\phi_i^{AB} + (1 - h_i^A) W_{ji}^A)} \right] \end{aligned} \quad (2.7)$$

where  $\phi_i^{AB} = \phi_i^A + \phi_i^B$ . Note that  $v$  plays no role in the correction. It is clear that adding  $\phi_i^B$  has introduced a dependency between the entirety of  $h$ . This means that a Gibbs chain will need to be run, in practice this chain proves to be short, even in the larger dimension tasks like MNIST.

### Examining and approximating the Correction

This correction (eq 2.7) is a non-trivial computation to make in that for every weight, a series of semi-large matrix operations have to be performed. As a result, Freat and Marsland propose an approximation for this correction.

It is useful to see the correction contours on axes  $\phi_i^A$  versus  $\phi_i^{AB}$ , for a positive weight  $W_{ij}^A$ . There are two plots for the two cases  $h_j^A = 0, 1$  These are plotted in figure 2.11.

The proposed approximation uses a Sigmoid to approximate [TODO WORDING Marcus help! I don't understand where the approximation comes from. Do you have to justify it or can I just drop it in and be like, proposed approximation will also be tested?](#)

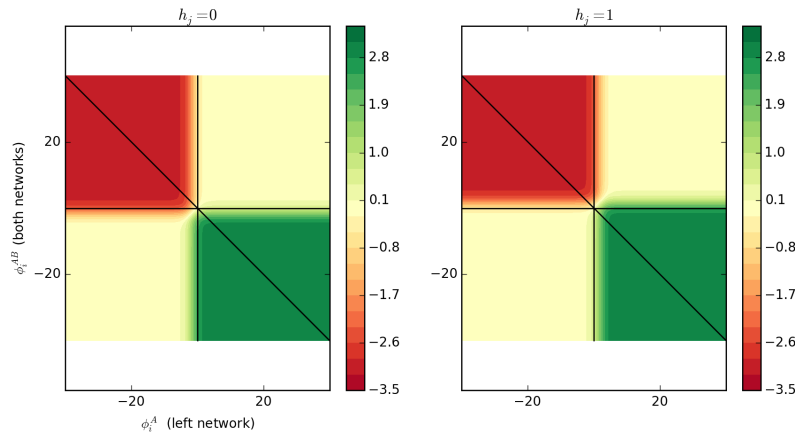


Figure 2.11: A diagram that shows the structure of the correction over the weighted sums into the hidden units of one of the RBMs versus both.



### What happened to the biases?

The ORBM utilises the hidden biases present on the RBMs when calculating the hidden states  $h^A$  and  $h^B$ . However, the visible biases are not used in sampling from the visible. The reasoning behind this being that the RBMs visible bias acts like the 'background rate', i.e. what visible reconstruction would we see from an all zero hidden activation. As visible bias are not captured in the ORBM, it is important that RBMs plugged into it's structure do not use a visible bias.

### Reconstructions and Dreams in the ORBM

Reconstructions in the ORBM are a natural extension of the inference algorithm.

- First hidden representations  $h^A$  and  $h^B$  are generated given an input  $v$ .
- The RBMs (RBM A and B) uses their respective hidden representations to generate visible patterns independently. That is, we can use the same way of performing Gibbs sampling we saw in equation 2.3.

This means that for a visible pattern there are potentially two reconstructions, one from each model.



## Chapter 3

# Design and Implementation

A large portion of the project was gaining enough understanding of the existing work on RBMs to be able to implement the ORBMs algorithm and architecture. This was crucial as bugs in the implementation are a threat to validity.

*TODO WORDING Make sure you talk about why the sigmoid is the most appropriate activation function! Maybe this should go in background..*

### 3.1 Evaluation Dataset Choice

- 2 Bit XOR - It's the minimal case. Can examine reconstructions and dreams/fantasy empirically. Also can check that code performs as expected by manually calculating expected values for the algorithm.
- X Neighbouring Bits On in Y Bit Pattern - A natural next step from 2 bit XOR. Still trivial to train an RBM to represent, quick feedback to ensure the algorithm works. Allows comparing the different approximations tractable as number of hidden units is small enough.
- Square Patterns in a 5 by 5 Pixel Image - Another natural step from neighbouring bits, trivial to construct a dataset. Very easy to compose. Can use the same model for both models in the ORBM Architecture.
- Rectangles in a 5 by 5 Pixel Image - Next step, builds on the previous test but adds a different model for each source. Also different shapes of model can be used.
- Continuous Rectangles (Pixel values 0 to 1) in a 5 by 5 Image - Another variation working towards the MNIST dataset, same as the previous but has continuous visible pixels. Allows to different intensities, interesting cases where the images overlap and increase in intensity.
- MNIST Handwritten Digit Dataset - Pixel values from 0 to 1, 28 by 28 pixel image (784 pixel features.) Used extensively in previous studies. Gives me confidence that RBMs can learn these representations. Also allows metaparameters don't have to be tweaked as much as studies already have found reasonable values. Non-trivial to evaluate, but more of a real world case.

## 3.2 Implementation Design

### 3.2.1 Language Choice

- Library choice meant efficiency with benefits of python for quick dev
- Easy to deploy to grid (versus java higher risk, lack of experience)
- Ordering of evaluation tasks made unit testing, with hand made test cases possible. Also less risk, remove the uncertainty around the RBM.
- testing approach

### 3.2.2 Program Architecture

The program was implemented with unit testing and composability in mind. It was important that the design supported comparing the Full and Approximated correction. By using Python, multiple inheritance could be leveraged. For instance adding continuous pixel value support to the Full Correction Sampler, one simply needed to extend Continuous Approximate Correction Sampler and the Full Correction Sampler, with no code actually in the new class. The architecture is pictured in the class diagram 3.1. There were three main roles these classes filled:

1. The Trainer was used to train the weights of a supplied RBM. It would do so using a supplied sampler, decoupling how samples were generated from the training process.
2. The RBM was the model of an RBM, storing the weight matrix, as well as parameter information. Also this supported conversion of the SciPy Sklearn libraries' RBM implementation into the implementation used in this project. Decoupling the concept of an RBM from the Sampler and the Trainer meant that RBMs could be 'plugged' into the ORBM architecture with ease.
3. The Sampler defined how to perform Gibbs sampling, the subclasses defining whether it is standard RBM sampling or ORBM sampling. This made it trivial to compose samplers, which was required for the ORBM samplers.

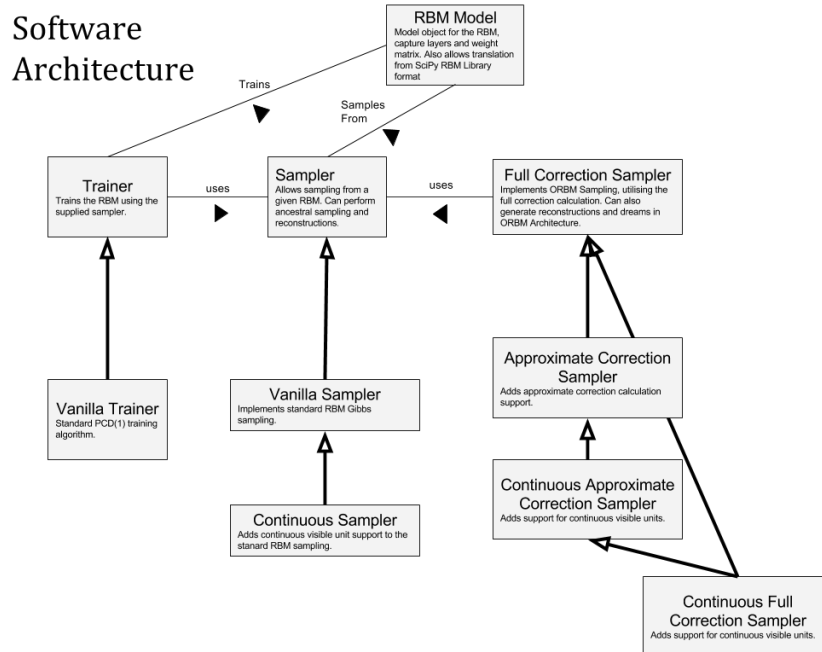


Figure 3.1: A figure showing the architecture for the implemented test suite in UML notation.



# Chapter 4

## Evaluation

### 4.1 Evaluation Design

#### 4.1.1 Mixing Time Evaluation

- List the types Traceplot, Densityplots, Gelbin and Rubin multiple sequence diagnostics
- Performed Geweke Diagnostic

#### 4.1.2 Reconstructions As a measure of Performance

- Cite literature where reconstructions are used, starting from early RBM training papers by Hinton, up to more recent in deeper networks - the goal being to show the reader that reconstructions are used well in practice.

#### 4.1.3 Growing Complexity of tasks

- Same models
- Growing dimensionality, bit strings to 2d to MNIST
- Allows comparison of Approx corrections in smaller scale - fail there don't continue using them...

### 4.2 Evaluation Methods

My [TODO WORDING x experiements](#) all followed the same high level process, working from trivial cases to more challenging tasks. The reasoning being that trivial cases make unit testing feasible, therefore ensuring conclusions can be drawn with regard to the algorithm and model, not an incorrect implementation.

The following evaluations aimed to evaluate the ORBM and it's inference algorithm by examining reconstructions given a multi-cause input. An example of a multi-cause input with two quadrilaterals is shown in figure 4.1. The optimal reconstructions are equivalent to the two images that combined form the input.

In the smaller dimensional cases the reconstructions are inspected by hand, manually plotting the reconstructions and the

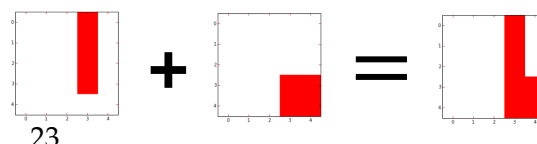


Figure 4.1: A figure illustrating two five by five pixel images combining to form a composite/multicause in-

h

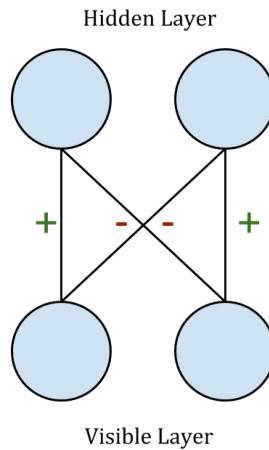


Figure 4.2: The two bit RBM for modelling a single bit on at a time in a two bit pattern. The diagonal weights are negative and the non-diagonal weights are positive.

frequency with which they occurred after a large amount of repetitions. In the larger dimensional tasks, two 'scores' were used to evaluate reconstructions against the training set.

To be able to use this reconstruction based evaluation, knowledge of the ground truth was required. This was so:

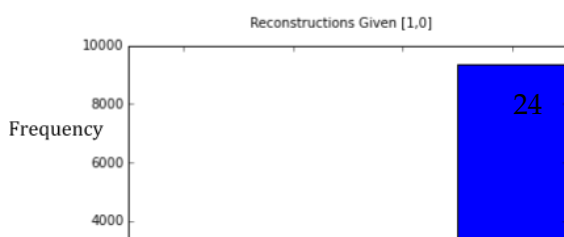
- RBMs could be trained and then plugged into the ORBM network.
- Reconstructions could be compared directly (by score) to the ground truth.

### 4.3 Starting from the minimal case, two bit XOR

#### Description

The starting point for the evaluation was examining the ORBM reconstructions in a two bit pattern, where two of the same model were combining to form the data. This model made one bit on in a two bit pattern at a time. That is it could either make  $[1, 0]$  or  $[0, 1]$ . The training set is the two bit XOR truth table. Two bit case also provides the minimal case to compare the full correction calculation, versus that of the proposed approximation (see equation 2.7).

#### Architecture and Parameters



Being a trivial dimensionality, the RBMs weights were constructed by hand, meaning that only the ORBMs inference algorithm was



Visible Input	Expected Reconstructions	
[1, 1]	[1, 0]	[0, 1]
[1, 0]	[1, 0]	[1, 0]
[0, 1]	[0, 1]	[0, 1]

Table 4.1: A Table showing the expected reconstructions from performing ORBM inference with various input patterns. The left and right hand column of the Expected Reconstructions column indicate the reconstructions from the left and right RBMs in the ORBM.

being evaluated and not the training of the RBMs plugged into it. This network is picture in figure 4.2, [TODO WORDING with weights greater than 5](#) the networks stable state results in either [1, 0] or [0, 1]. The RBMs had 2 hidden units and 2 visible units.

## Method

This RBM was checked to ensure that it behaved in practice by ensuring it's reconstructions matched the input for a large frequency of reconstructions. The results of trying this with the input

[1, 0] are showing in figure 4.3. The most common reconstruction is [1, 0], which matches the input. We see that 500 reconstructions out of the 10,000 reconstructions were incorrect, corresponding to [0, 1]. This is due to the stochastic nature of the RBM.

In a similar way to the reconstructions, ancestral samples can be taken from the model and evaluated. [TODO WORDING In larger dimensions Hinton has shown RBMs to be poor generative models without the extra layers above...](#) however in the small dimensions of this task the dreams should match the training set:

A histogram of frequencies of reconstructions given the RBM is in the free phase, is shown in figure 4.4. Again the model behaves as expected, generating dream patterns that match the training dataset the majority of the time. Given this good model, the XOR RBM was duplicated and plugged into the ORBM structure as picture in figure 4.5.

The inference algorithm was run in the ORBM architecture for various visible inputs, giving two hidden vectors for the two representations  $h^A$  and  $h^B$ . For each pair of hidden vectors, a reconstruction was created, the process illustrated in figure 4.6. This process was repeated in a similar way to how the reconstructions were evaluated in the lone RBM, counting the frequency each reconstruction occurred over 1500 runs.

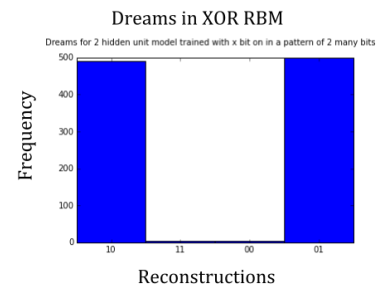


Figure 4.4: Dreams in the XOR RBM, note how only the training data is present.

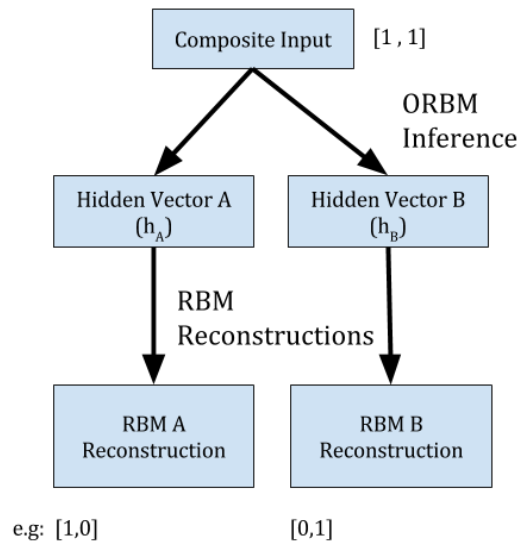


Figure 4.6: The process of generating reconstruction is shown in this diagram.

### XOR ORBM Analysis

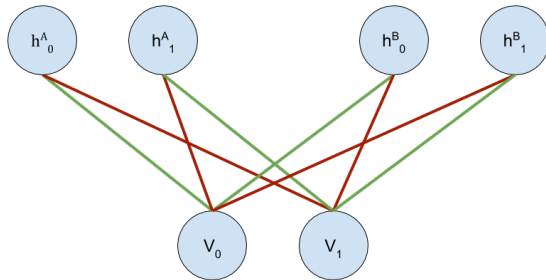


Figure 4.5: Figure showing how the XOR RBM fits into the ORBM architecture. These reconstructions were compared to one of the RBMs trained to recognise a single pattern being on in two bits. As expected a machine that has been trained to recognise one bit, has no concept of two bits being on and hence the reconstructions show no mechanism for separating the sources. This is illustrated in figure 4.9.

The results of repeating this process with the input  $[1,0]$  yielded successful results. We would hope that ORBM architecture could also handle inference with just a single subject. The results of this are shown in figure 4.8. **TODO WORDING Do the markers care that it works in this case... it's so trivial does it even mean anything for larger images..**

### Two bit XOR Conclusion

The ORBM was able to extract the target patterns  $[1,0]$  and  $[0,1]$  given the input  $[1,1]$ . This was the case for both the Approximated and Full corrections, which gave confidence that the more computationally efficient Approximated correction could be relied on going forward — as for larger datasets it was a lot faster in practice. The generative model also copes with

The results of this process are shown in figure 4.7. Where the x-axis shows the reconstructions of the form A reconstruction, B reconstruction. The results show how applying the full correction compares to the approximated correction.

The ORBM is able to separate the causes with both types of correction, as the model is being duplicated, it produces  $[1,0]$  and  $[0,1]$  approximately half the time and symmetrically  $[0,1]$   $[1,0]$  the other half of the time.

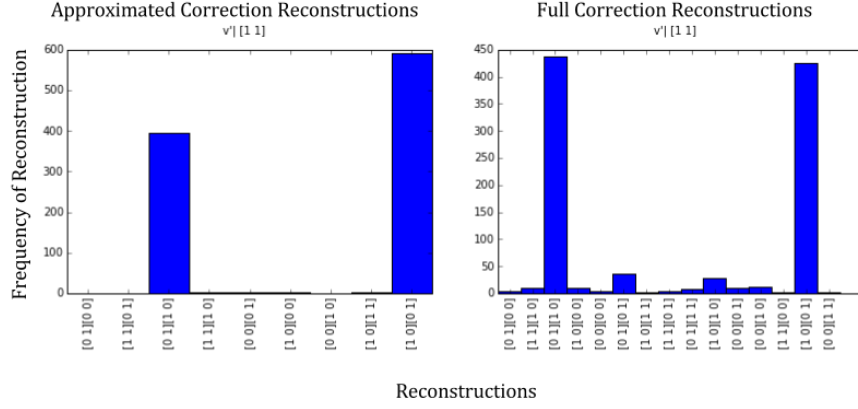


Figure 4.7: A figure with the results of running ORBM inference 1000 times, with the two different approaches to calculating the correction. The frequency for which a given reconstruction occurred is shown.

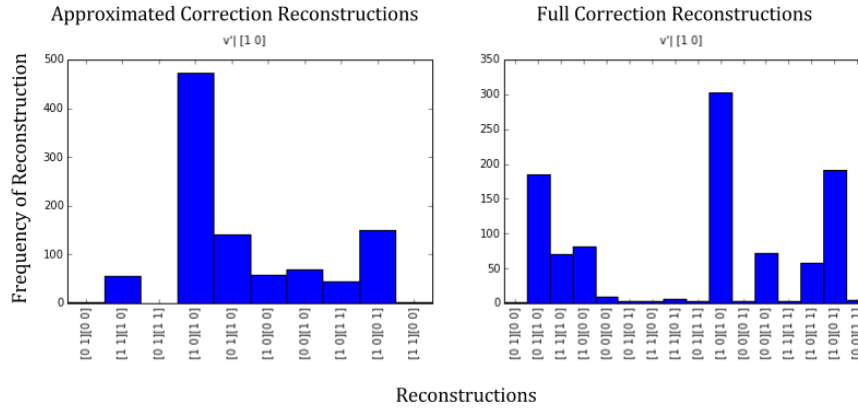


Figure 4.8: A figure showing the result of generating 1000 reconstructions with the ORBM with  $[1, 0]$  as input, again comparing the Approximated correction and Full correction.

the subject overlapping, which in the two bit case arised when  $[1, 0]$  is composed with  $[1, 0]$ . In both the Approximated and Full corrections the highest occuring reconstruction is the target  $[1, 0]$ , however there appears to be a lot less confidence.

#### 4.4 $Y$ bit pattern, $X$ neighbouring bits on

##### Description

A natural next step from a single bit on in a 2 bit pattern, is moving up to  $X$  bits side by side on in a  $Y$  bit pattern is a next step. For example if  $Y = 4$  and  $X = 2$  then a valid inputs are the rows of the following:

$$dataset = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & 0 & 1 & 1 \end{bmatrix}$$

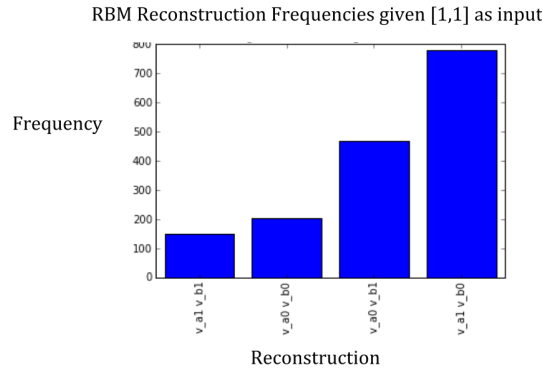


Figure 4.9: A figure showing the result of generating 1000 reconstructions with an RBM. It is clear the RBM has no mechanism to separate the sources. Hence showing it  $[1,1]$  when it has only been trained on XOR results in no separation, i.e. reconstructions including  $[1,1]$ .

This allows for some interesting cases, for instance using the same example as above of  $Y = 4$  and  $X = 2$ , we can see how the ORBM is able to separate interesting patterns such as  $[1, 1, 1, 0]$ , which is a composition of  $[1, 1, 0, 0]$  and  $[0, 1, 1, 0]$ .

### Architecture and Parameters

- $|Y|$  Hidden units per RBM, any less and we were unable to train the RBM to make the target reconstructions or dreams.
- Random normally distributed, mean and standard deviation of one, starting weights. The RBM was trained with 1000 epochs and a learning rate of 0.002.
- A sample of 10,000 dreams were sampled from the RBMs, and then plotted on a histogram in a similar fashion to figure ???. The frequency of produced dreams was ensured to be approximately equal and that the dreams should directly match the training set. This is feasible as the ground truth patterns are known.

### Method

For values of  $Y$  where  $y \in Y | 2 < y < 7$  and all values of  $X$  where  $x \in X | 2 < x < y < 7$  The process below is repeated.

The inference algorithm was run for 1000 reconstructions in the ORBM architecture for all unique compositions of the training set. The result was 1000 reconstructions for both the Full Correction and Approximate Correction, for some interesting compositions of the dataset.

**TODO CITE: Need to grab the images out of the ipython notebook.**

## 4.5 2D Patterns in a small dimensional space

- Dataset:
  - 2x2 square in a 5x5 image. Dataset of every possible configuration of said square.
- Method:

- train a single RBM to represent 2x2 squares in 5x5 space. Nice and small, can still inspect reconstructions, and dreams.
- Compose two square images with each other, can the ORBM architecture separate the images.

## 4.6 2D Pattern, Different Rectangle Separation

- Dataset
  - n by m rectangles in 5x5 Images (TODO-IMAGE-IN-HERE) where n and m are not always equal.
- Method
  - Superimpose two images as the same way as before. How well can they separate.

## 4.7 MNIST Digits Reconstructions

### Description

MNIST is a widely used dataset of handwritten digits (0 – 9). This task explored the non-trivial task of given two handwritten digits composited to form one image, how effectively can the ORBM separate the sources compared to the naive RBM? An example input is illustrated in figure ??.

### Architecture and Parameters

- MNIST Digit images are 28 by 28 pixel images, with pixel values between 0 – 1.
- 10, 100 Hidden unit RBMs trained on 500 examples of each digit respectively.
- Reconstructions and dreams of these RBMs were inspected by hand to ensure that they resembled the dataset.

### Method

For every digit dataset (of size 500), each digit in each dataset was composed with every other digit in every other dataset. Given this set of composite datasets, the corresponding RBMs were plugged in the ORBM architecture and used to create reconstructions. 100 Gibbs iterations were used when calculating the correction (generating the hidden states ( $h^A$  and  $h^B$ ) for a given input). These RBMs were also used to create standard RBM reconstructions. Given the ORBMs reconstructions (two per image) and the RBM reconstructions (also two per image) two scores were applied to compare these reconstructions to the ground truth:

**Cross Entropy** The cross entropy was calculated between the reconstruction and the ground truth.

**Cosine Angle** The angle between the flattened reconstruction vectors was computed, then negated to give a 'score'. The higher the score the smaller the angle between the reconstruction vector and the ground truth.

These scores can then be summed over each digit and over the entire dataset to find a single score for each digit composed with every other digit forming a 9 by 9 matrix. A cell of this matrix (say  $j, i$ ) corresponds to the difference between the ORBM score and RBM scores for digits  $j$  and  $i$  composited together. This entire process was repeated 10 times gain certainty in the results — given the stochasticity of the RBM and ORBM. This process is shown in Algorithm 1.

**Data:** MNIST Digit datasets, each with 500 examples, 10 RBMs trained on MNIST data 0–9 respectively

**Result:** Composite Datasets for every digit

**for** 10 repetitions to gain confidence in results **do**

**for** All MNIST digits datasets **do**

**for** Every digit example in the MNIST digit **do**

            composite the current digit dataset with every other dataset;

            Generate reconstructions on that dataset using the ORBM, RBM;

            Calculate the Cosine Angle score and Cross Entropy for both the RBM's and ORBM's reconstructions versus the ground truth;

**end**

        Sum the scores over every digit example;

**end**

    Tabulate the summed scores in a several matrices indexed by the digits being composited;

**end**

**Algorithm 1:** Algorithm explaining how the Scores matrices were constructed.

By finding the difference between the ORBMs reconstructions and the RBMs reconstruction we can create another matrix, in which a positive value represents a 'win' for the ORBMs reconstructions and a negative value represents a 'loss' (where the RBM outperformed the ORBM).

## Results

These score difference matrices are seen in for the Cross Entropy and Cosine Angle scores in 4.10a and 4.10b respectively, where the colour encodes the score. For the Cosine score, using the approximate correction, zero composited with every other digit yielded the best results for ORBM relative to the RBM, and nine the converse yielding particularly bad results (worse by a factor of 10 compared to zero). Given this, I plotted the RBMs Score versus the ORBM Score as a scatter plot, where each point corresponds to a score. This allows us to examine if the ORBM is performing worse as a whole, or if in some cases it is performing better than the RBM. These plots for zero and nine are shown in 4.11.

### 4.7.1 MNIST Analysis

It is clear that RBM has better scores than the ORBM when summed over the dataset, and surprising the Approximated correction results in better scores than the Full Correction calculation. For the dataset wide scores 4.10 I would expect to see symmetry in that, for example four composited with five, will be an equivalent dataset to five composited with four. There is a symmetry present in the Full Correction, but this is not the case for the Approximate Correction.

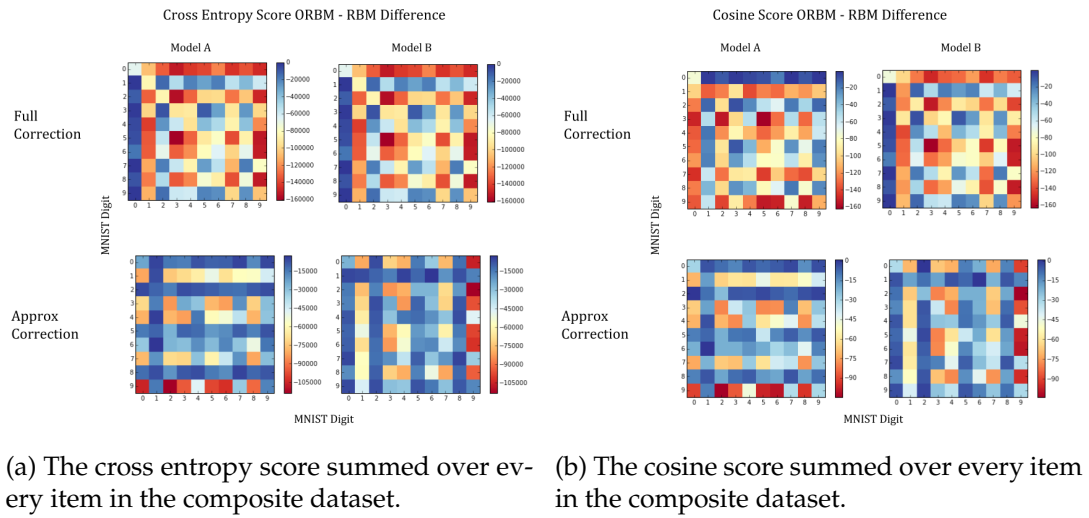


Figure 4.10: Dataset-wide MNIST Score Results

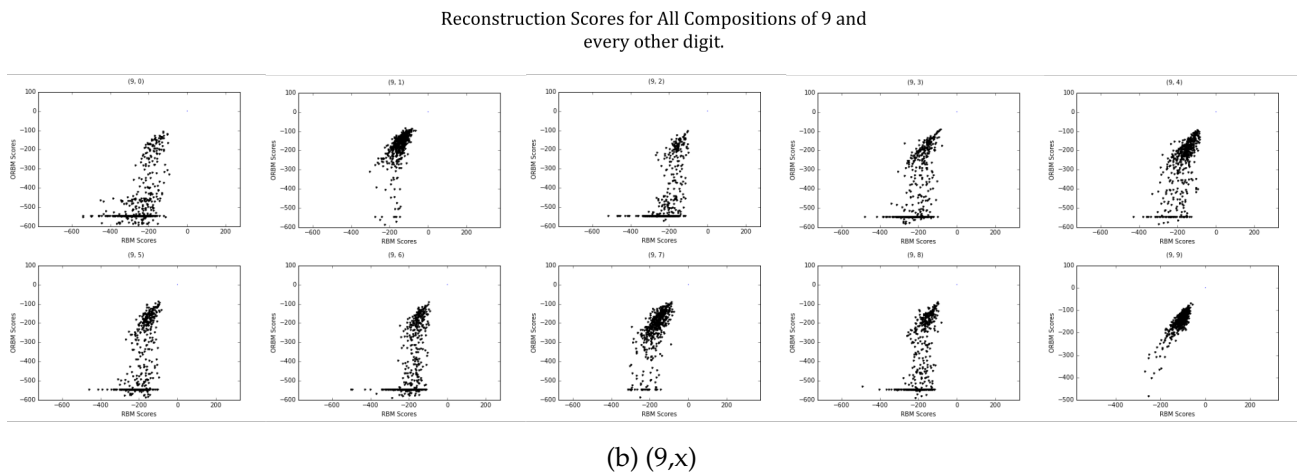
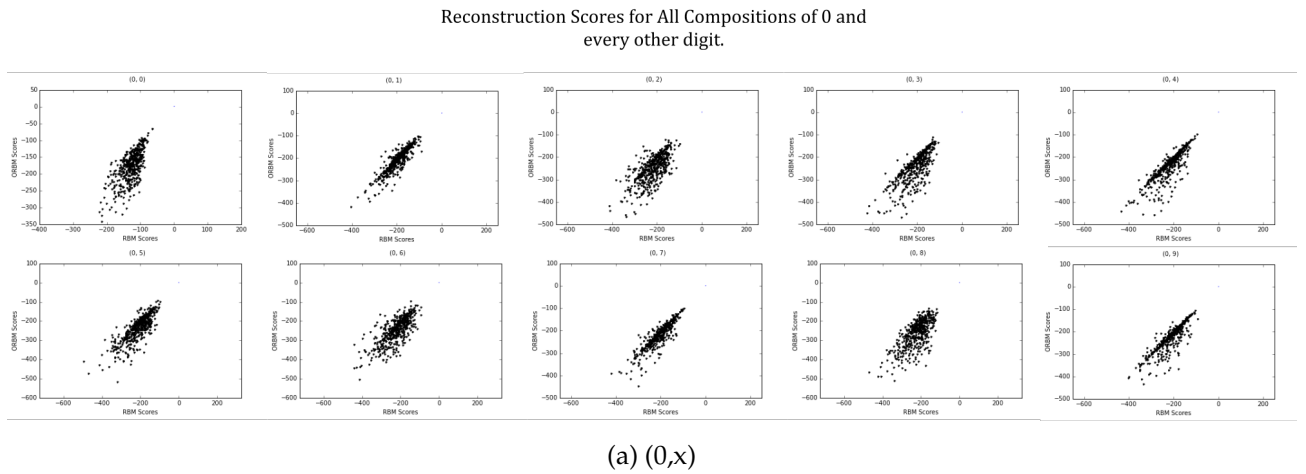


Figure 4.11: Cosine Score breakdown for the highest and lowest performing datasets, 0 and 9.





## Chapter 5

# Conclusions and Future Work

- ORBM inference is working in smaller cases really well. Even better than the RBM!
- Training algorithm for this generative model has been proposed. Haven't been able to get it going but it could be really promising - Amounts to blind source separation.



# Bibliography

- [1] BARBER, D. *Bayesian Reasoning and Machine Learning*. Cambridge University Press, New York, NY, USA, 2012.
- [2] HINTON, G. E. To recognize shapes, first learn to generate images. *Progress in brain research* 165 (2007), 535–547.
- [3] HINTON, G. E., OSINDERO, S., AND TEH, Y.-W. A fast learning algorithm for deep belief nets. *Neural Comput.* 18, 7 (July 2006), 1527–1554.
- [4] HINTON, G. E., AND SEJNOWSKI, T. J. Analyzing cooperative computation. *Proceedings of the 5th Annual Congress of the Cognitive Science Society* (may 1983).
- [5] HINTON, G. E., AND SEJNOWSKI, T. J. Parallel distributed processing: Explorations in the microstructure of cognition, vol. 1. MIT Press, Cambridge, MA, USA, 1986, ch. Learning and Relearning in Boltzmann Machines, pp. 282–317.
- [6] HOPFIELD, J. J. Neural networks and physical systems with emergent collective computational abilities. *Proceedings of the National Academy of Sciences* 79, 8 (1982), 2554–2558.
- [7] NEAL, R. M. Connectionist learning of belief networks. *Artificial intelligence* 56, 1 (1992), 71–113.
- [8] PEARL, J. Reverend Bayes on inference engines: A distributed hierarchical approach. In *Proceedings of the American Association of Artificial Intelligence National Conference on AI* (Pittsburgh, PA, 1982), pp. 133–136.
- [9] SMOLENSKY, P. *Foundations of harmony theory: Cognitive dynamical systems and the sub-symbolic theory of information processing*. Parallel distributed processing: Explorations in the ..., 1986.