

Article

The Partial Information Decomposition of Generative Neural Network Models

Tycho M.S. Tax ^{1,†}, Pedro A.M. Mediano ^{2,†} and Murray Shanahan ²

¹ Corti; tt@cortilabs.com

² Imperial College London; {pmediano, m.shanahan}@imperial.ac.uk

† These authors contributed equally to this work.

Academic Editor: name

Version August 13, 2017 submitted to Entropy

Abstract: In this work we study the distributed representations learnt by generative neural network models. In particular, we investigate the properties of redundant and synergistic information that groups of hidden neurons contain about the target variable. To this end, we use an emerging branch of information theory called Partial Information Decomposition (PID) and track the informational properties of the neurons through training. We find two differentiated phases during the training process: a first, short phase in which the neurons learn redundant information about the target, and a second phase in which neurons start specialising and each of them learns unique information about the target. We also find that in smaller networks individual neurons learn more specific information about certain features of the input, suggesting that learning pressure can encourage disentangled representations.

Keywords: Partial Information Decomposition; Neural Networks; Information Theory

1. Introduction

Neural networks are famously known for their excellent performance, yet infamously known for their thin theoretical grounding. While common deep learning “tricks” that are empirically proven successful tend to be later discovered to have a theoretical justification (e.g. the Bayesian interpretation of dropout [1,2]), deep learning research still operates “in the dark” and is guided almost exclusively by empirical performance.

One common topic in learning theory is the study of data representations, and in the case of deep learning it is the hierarchy of such representations that is often hailed as the key to neural networks’ success [3]. More specifically, *disentangled representations* have received increased attention recently [4–6] and are particularly interesting given their reusability and their potential for transfer learning [7,8]. A representation can be said to be disentangled if it has factorisable or compositional structure, and has consistent semantics associated to different generating factors of the underlying data generation process.

In this article we explore the evolution of learnt representations in the hidden layer of a restricted Boltzmann machine as it is being trained. Are groups of neurons correlated or independent? To what extent do neurons learn the same information or specialise during training? If they do so, when? To answer these questions, we need to know how multiple sources of information (the neurons) contribute to the correct prediction of a target variable – which is known as a multivariate information problem.

To this end, the Partial Information Decomposition (PID) framework by Williams and Beer, which seeks a rigorous mathematical generalisation of mutual information to the multivariate setting, provides an excellent foundation for this study [9]. In PID, the information that multiple sources

33 contain about a target is decomposed into unique non-negative *information atoms*, the distribution of
34 which gives insight into the interactions between the sources.

35 1.1. Why information theory?

36 Information theory was developed to optimise communication through noisy channels, and it
37 quickly found other areas of application in the mathematical and computer sciences. Nevertheless, it is
38 not commonly linked to machine learning and it is not part of the standard deep learning engineer's
39 toolkit or training. So why, then, is information theory the right tool to study neural networks?

40 To answer that question, we must first consider some of the outstanding theoretical problems in
41 deep learning: what kind of stimuli do certain neurons *encode*; how do different layers *compress* certain
42 features of an input image; or how can we *transfer* learnt information from one dataset to another?

43 These problems (encodings, compression, transfer) are precisely among the problems information
44 theory was made to solve. Casting these questions within the established framework of information
45 theory gives us a solid language to reason about these systems and a comprehensive set of quantitative
46 methods to study them.

47 We can also motivate this choice from a different perspective: in the same way as neuroscientists
48 have been using information theory to study computation in biological brains, here we try to
49 understand an artificially developed *neural code* [10]. Although the code used by artificial neural
50 networks is most likely much simpler than the one used by biological brains, deep learning researchers
51 can benefit from the neuroscientists' set of tools.

52 1.2. Related work

53 When it comes to representations, the conventional way of obtaining insights about a network
54 has typically been through visualisation. Famously, [Le et al.](#) trained a neural network on web-scraped
55 images and reported finding neural receptive fields consisting mostly of human faces, human bodies
56 and cat faces [11]. Later, [Zeiler and Fergus](#) devised a technique to visualise the features learnt by
57 neurons in hidden layers, and provided good qualitative evidence to support the long-standing claim
58 that deeper layers learn increasingly abstract features of the input [12].

59 While visualisation is a great exploration tool, it provides only qualitative insights and is therefore
60 unable to make strong statements about the learning dynamics. Furthermore, as later work showed,
61 the specific values of weights are highly sensitive to the details of the optimisation algorithm used,
62 and therefore cannot be used to make definite judgements about the network's behaviour [13,14].

63 More recently there is a small line of emerging work investigating the behaviour of neural
64 networks from an information-theoretic perspective [15–20], with some work going as far back as
65 [21]. The most relevant of these is the work by [Schwartz-Ziv and Tishby](#), who show that feed-forward
66 deep neural networks undergo a dynamic transition between drift- and diffusion-like regimes during
67 training.

68 The main contribution of this article is to show how PID can be used for the analysis of learning
69 algorithms, and its application to neural generative models. The results of our PID analysis show two
70 distinct learning phases during the optimisation of the network, and a decrease in the specialisation of
71 single neurons in bigger networks.

72 2. Methods

73 2.1. Restricted Boltzmann Machines

74 We deal with the problem of multiclass classification, in which we have a dataset \mathcal{D} of (\mathbf{x}, y)
75 tuples, where y is a discrete *label* (also called the *target* variable) and \mathbf{x} is a vector of predictor variables.
76 The goal is to learn an approximation to the joint distribution of the predictors and the labels, $p(\mathbf{x}, y)$.
77 We will use a class of neural generative models known as Boltzmann Machines.

78 Boltzmann Machines (BM) are energy-based probabilistic graphical models, the origin of which
 79 goes as far back as Paul Smolensky's Harmonium [22]. Of particular interest are the so-called Restricted
 80 Boltzmann Machines (RBM). These are called *restricted* because all the nodes in the model are separated
 81 in two layers, and intra-layer connections are prohibited. These typically receive the names of *visible*
 82 and *hidden* layers.

In this article we follow [23] and perform classification with a *discriminative* RBM (DRBM). To do this we introduce the vector of target classes y as part of the visible layer, such that the DRBM represents the joint distribution over the hidden, visible, and target class variables. The distribution parametrized by the DRBM is:

$$p(y, \mathbf{x}, \mathbf{h}) = \frac{1}{Z} e^{-E(y, \mathbf{x}, \mathbf{h})}, \quad (1)$$

where $E(y, \mathbf{x}, \mathbf{h})$ is the DRBM *energy function*, given by

$$E(y, \mathbf{x}, \mathbf{h}) = -\mathbf{h}^T \mathbf{W} \mathbf{x} - \mathbf{b}^T \mathbf{x} - \mathbf{c}^T \mathbf{h} - \mathbf{h}^T \mathbf{U} \vec{y} - \mathbf{d}^T \vec{y}, \quad (2)$$

83 where $\vec{y} = (1_{y=i})_{i=1}^C$ for the C different classes. For comparison, the energy function for a standard
 84 RBM is the same but with the last two terms removed. Figure 1 shows a schematic diagram of a DRBM
 85 and the variables involved.

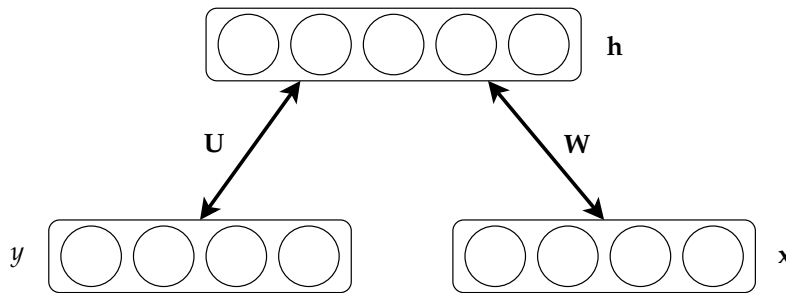


Figure 1. Graphical representation of the Discriminative Restricted Boltzmann Machine and its components. Vectors \mathbf{x} and \mathbf{y} correspond to the training input and label respectively, \mathbf{h} is the activation of the hidden neurons, and \mathbf{U} and \mathbf{W} are the weight matrices to be learnt. (Adapted from [23].)

Now the model is specified, we calculate the predictive posterior density $p(y|\mathbf{x}, \theta)$ given DRBM parameters $\theta = \{\mathbf{W}, \mathbf{b}, \mathbf{c}, \mathbf{U}, \mathbf{d}\}$. At this point the restricted connectivity of RBMs comes into play – this connectivity induces conditional independence between all neurons in one layer given the other layer. This resulting intra-layer conditional independence allows us to factorise $p(y_i, x_i|\theta)$ and to write the following conditional distributions [24]:

$$\begin{aligned} p(X_i = 1|\mathbf{h}) &= \sigma(b_i + \sum_j W_{ji} h_j) \\ p(H_j = 1|y, \mathbf{x}) &= \sigma(c_j + U_{jy} + \sum_i W_{ji} x_i) \\ p(y|\mathbf{h}) &= \frac{e^{d_y + \sum_j U_{jy} h_j}}{\sum_{y^*} e^{d_{y^*} + \sum_j U_{jy^*} h_j}}, \end{aligned} \quad (3)$$

86 where $\sigma(t) = (1 + e^{-t})^{-1}$ is the standard sigmoid function. With these equations at hand we can
 87 classify a new input vector \mathbf{x}^* by sampling from the predictive posterior $p(y|\mathbf{x}^*, \theta)$, or we can sample
 88 from the joint distribution $p(y, \mathbf{x}|\theta)$ via Gibbs sampling.

Finally, the network needs to be trained to find the right parameters θ that approximate the distribution of the data. We use a standard maximum likelihood objective function,

$$\mathcal{L}(\theta) = - \sum_i \log p(y_i, x_i | \theta) . \quad (4)$$

89 Gradients of this objective cannot be obtained in closed form, and we must resort to contrastive
 90 sampling techniques. In particular, we use the Contrastive Divergence (CD) algorithm [24] to estimate
 91 the gradient and we apply fixed-step size stochastic gradient updates to all parameters in the network.
 92 The technical details of CD and other contrastive sampling estimators are outside the scope of this
 93 paper, and the interested reader is referred to the original publications for more information [24,25].

94 2.2. Information Theory

95 In this section we introduce a few relevant tools from Information Theory (IT) that we will use to
 96 analyse the networks trained as explained in the previous section. For a broader introduction to IT and
 97 more rigorous mathematical detail we refer the reader to [26].

98 We focus on systems of discrete variables with a finite number of states. Throughout the paper
 99 we will deal with the scenario in which we have one *target* variable and a number of *source* variables.
 100 We refer to the target variable as Y (matching the nomenclature in Sec. 2.1), to the source variables as
 101 Z_i and let Z denote generically any nonempty subset of the set of all sources. Summations always run
 102 over all possible states of the variables considered.

Mutual Information (MI) is a fundamental quantity in IT that quantifies how much information is shared between two variables Z and Y , and is given by

$$I(Y; Z) = \sum_{y,z} p(y, z) \log \left(\frac{p(y, z)}{p(y)p(z)} \right) . \quad (5)$$

MI can be thought of as a generalised (non-linear) correlation, which is higher the more a given value of Z constrains possible values of Y . Note that this is an average measure – it quantifies the information about Y gained when observing Z *on average*. In a similar fashion, *specific information* [27] quantifies the information contained in Z associated with a particular outcome y of Y , and is given by

$$I(Y = y; Z) = \sum_z p(z|y) \log \left(\frac{p(y|z)}{p(y)} \right) . \quad (6)$$

Specific information quantifies to what extent the observation of Z makes outcome y more likely than otherwise expected based on the prior $p(y)$. Conveniently, MI can easily be written in terms of specific information as

$$I(Y; Z) = \sum_y p(y) I(Y = y; Z) .$$

103 2.2.1. Non-negative Decomposition of Multivariate Information

104 In this section we discuss the main principles of the PID framework proposed by [Williams and Beer](#). Technical details will not be covered and the interested reader is referred to the original paper [9].

105
 106 The goal of PID is to decompose the joint mutual information that two or more sources have about the target, $I(Y; \{Z_1, Z_2, \dots, Z_n\})$, into interpretable non-negative terms. For simplicity, we present the two-variable case here, although the framework applies to an arbitrary number of sources. In the two-variable PID (or PI-2), there are three types of contributions to the total information of $\{Z_1, Z_2\}$
 107
 108
 109
 110 about Y which form the basic atoms of multivariate information:

- 111 • *Unique* information U one of the sources provides and the other does not.
- 112 • *Redundant* information R both sources provide.

- 113 • *Synergistic* information S the sources provide jointly, which is not known when either of them is
114 considered separately.

115 There is a very intuitive analogy between this decomposition and set theory – in fact the decomposition
116 for any number of variables can be shown to have a formal lattice structure if R is mapped to the
117 set intersection operation. This mapping corresponds to the intuitive notion that R should quantify
118 the *overlapping information* of Z_1 and Z_2 . Consequently, these quantities can be represented in a Venn
diagram called the *PI diagram*, shown in Figure 2.

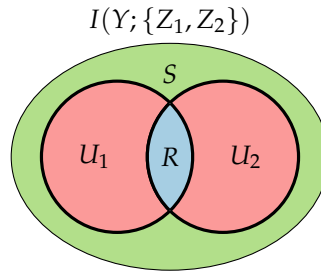


Figure 2. Partial information diagram for 2 source variables and a target. Outer ellipse corresponds to the MI between both sources and the target, $I(Y; \{Z_1, Z_2\})$, and both inner circles (highlighted in black) to the MI between each source and the target, $I(Y; Z_i)$. Coloured areas represent the PI terms described in the text.

119 Mathematically, the relation between MI and S , R and U (which we refer to jointly as *PI terms*)
can be written as follows:

$$\begin{aligned} I(Y; Z_1) &= R(Y; \{Z_1, Z_2\}) + U(Y; Z_1), \\ I(Y; Z_2) &= R(Y; \{Z_1, Z_2\}) + U(Y; Z_2), \\ I(Y; \{Z_1, Z_2\}) &= R(Y; \{Z_1, Z_2\}) + U(Y; Z_1) + U(Y; Z_2) + S(Y; \{Z_1, Z_2\}). \end{aligned} \quad (7)$$

120 This is an underdetermined system of three equations with four unknowns, which means the PI
121 decomposition in itself does not provide a method to calculate the PI terms. To do that we need to
122 specify one of the four variables in the system, typically by providing an expression to calculate either
123 R or S . There are a number of proposals in the literature [28–31], but at the time of writing there is no
124 consensus on any one single candidate.

In this study we follow the original proposal by [Williams and Beer](#) and use I_{\min} as a measure of redundancy, defined as

$$R(Y; \{Z_1, Z_2\}) = I_{\min}(Y; \{Z_1, Z_2\}) = \sum_y p(y) \min\{I(Y = y; Z_1), I(Y = y; Z_2)\}, \quad (8)$$

125 where $I(Y = y; Z_i)$ is the specific information¹ in Eq. 6. Despite the known flaws of I_{\min} , we choose it
126 for its extensibility, tractability and inclusion-exclusion properties. With this definition of redundancy
127 and the standard MI expression in Eq. 5 we can go back to system 7 and calculate the rest of the terms.

While all the terms in PI-2 can be readily calculated with the procedure above, with more sources the number of terms explodes very quickly – to the point that the computation of all PI terms is intractable even for very small networks. Conveniently, with I_{\min} we can compute the overall redundancy, synergy, and unique information terms for arbitrarily many sources – restricted only by

¹ In their original article, [DeWeese and Meister](#) propose two quantities to measure “the information gained from one symbol:” specific information and specific surprise. Confusingly, the quantity that [Williams and Beer](#) chose and named specific information is actually the specific surprise of [DeWeese and Meister](#).

the computational cost and amount of data necessary to construct large joint probability tables. We write here the overall redundancy, synergy and unique information equations for completeness, but the interested reader is referred to [32] for a full derivation:

$$\begin{aligned}
 R(Y; \{Z_1, \dots, Z_n\}) &= I_{\min}(Y; \{Z_1, \dots, Z_n\}) \\
 &= \sum_y p(y) \min\{I(Y = y; Z_1), \dots, I(Y = y; Z_n)\} \\
 S(Y; \{Z_1, \dots, Z_n\}) &= I(Y; \{Z_1, \dots, Z_n\}) - I_{\max}(Y; \{\mathbf{A} \in \{Z_1, \dots, Z_n\} : |\mathbf{A}| = n - 1\}) \\
 U(Y; Z_i) &= I(Y; Z_i) - I_{\min}(Y; \{Z_i, \{Z_1, \dots, Z_n\} \setminus Z_i\}),
 \end{aligned} \tag{9}$$

128 where I_{\max} is defined exactly the same as I_{\min} except substituting max for min [33].

129 3. Results

130 Instead of generating a synthetic dataset, we opt for a more realistic benchmark and use the
 131 MNIST dataset of hand-written digits. We use a stochastic binarised version of MNIST – every time an
 132 image is fed as input to the network, the value of each pixel is sampled from a binomial distribution
 133 with a probability equal to the normalised intensity of that pixel. Then we use Eqs. 3 to sample the
 134 state of the network, and repeat this process to build the probability distributions of interest.

135 For training, the gradients are estimated with contrastive divergence [24] and the weights are
 136 optimised with vanilla stochastic gradient descent with fixed learning rate (0.01). We did not make
 137 strong efforts to optimise the hyperparameters used during training.

138 To produce the results below we train an ensemble of 100 RBMs and take snapshots of these
 139 networks during training. Each RBM in the ensemble is initialised and trained separately using a
 140 different PRNG seed. All information-theoretic measures are reported in bits and debiased with
 141 random permutation tests – to debias the estimation of any measure on a given set of data we generate
 142 many surrogate data sets by randomly permuting the original data, calculate the mean of the measure
 143 across all surrogates and subtract this from the original estimation on the unshuffled data [34].

144 3.1. Classification error and mutual information

145 First, we train a small RBM with 20 hidden neurons and inspect its learning curve during training.
 146 In Fig. 3 we show the classification error and the mutual information between the predicted labels \hat{Y}
 147 and the real labels Y during training, averaged for the ensemble of 100 RBMs.

148 As expected, classification error decreases and MI increases during training, the relationship
 149 between the two being an almost perfect line. This gives us an intuitive correspondence between a
 150 relatively abstract measure like bits and a more easily interpretable measure like error rate. We note
 151 that a perfect classifier with 0 error rate would have $I(\hat{Y}, Y) = H(Y) = \log_2(10) \approx 3.32$ bit.

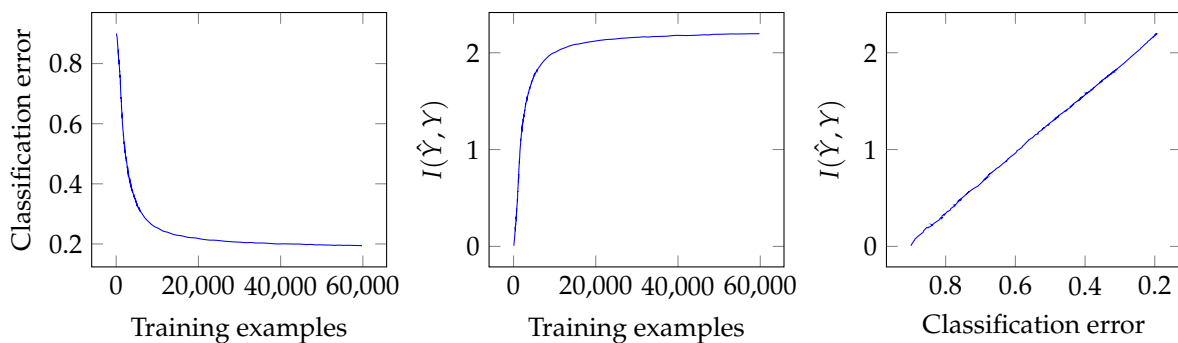


Figure 3. Classification error and mutual information between real and predicted labels, $I(\hat{Y}, Y)$, calculated through training. Note: X-axis in the rightmost plot is reversed for illustration purposes, so that training time goes from left to right.

152 As should be apparent to any occasional reader of the ML literature, the classification error
 153 presented in Fig. 3 is worse than the authors reported originally in [23], and significantly worse than
 154 the state of the art on this dataset. The main reason for this is that we are restricting our network to a
 155 very small size to obtain a better resolution of the phenomena of interest.

156 3.2. Phases of learning

157 In this section we investigate the evolution of the network through training and show three
 158 complementary pieces of evidence for the existence of two separate learning phases. We describe the
 159 main results illustrated in Figs. 4 and 5.

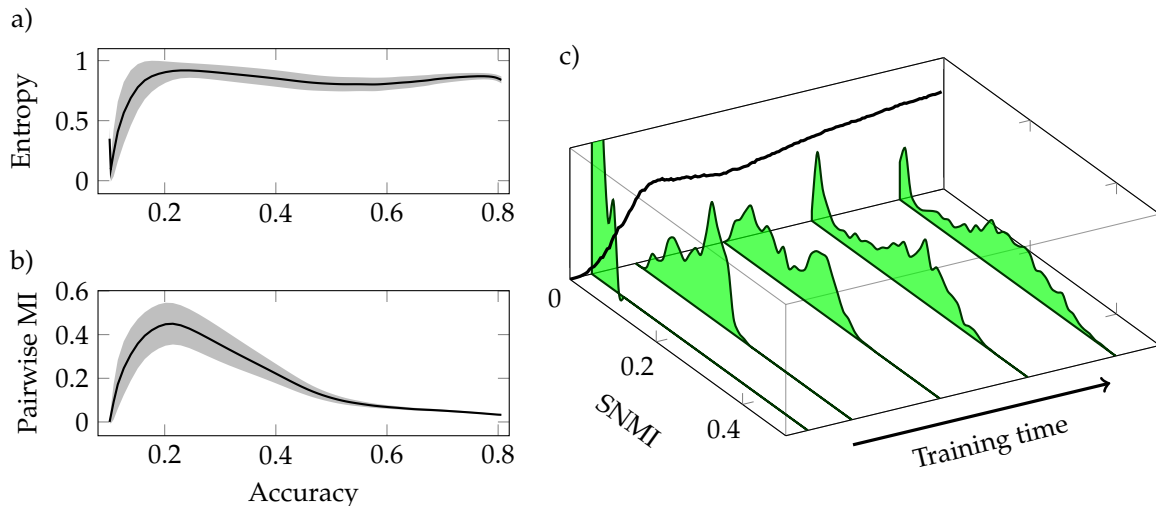


Figure 4. Single-neuron entropy and mutual information follow non-trivial patterns during training. (a) Entropy quickly rises up to close to its maximum value of 1 bit. (b) Inter-neuron correlation as measured by pairwise MI peaks midway through training. (c) Histograms of Single-Neuron MI split midway through training, implying that some neurons actually lose information. Average SNMI is shown in black projected on the frame box.

160 First, in Fig. 4a, we show the evolution of the average entropy of single neurons in the hidden
 161 layer, where the average is taken over all neurons in the same network. Entropy increases rapidly
 162 at the start of training until it settles around the 0.8 to 0.9 bit range, relatively close to its maximum
 163 possible value of 1 bit. This means that throughout most of the training, including the final state,
 164 most of the informational capacity of the neurons is being actively used – if this were not the case, in
 165 a network with low entropy in which most neurons do not change their states much, the encoding
 166 capability of the network would be heavily reduced.

As a measure of inter-neuron correlation, we calculate the average Pairwise Mutual Information (PWMI) between hidden neurons H_i , defined as

$$\text{PWMI} = \langle I(H_i; H_j) \rangle_{ij}.$$

167 PWMI is shown in Fig. 4b and is the first sign of the transition mentioned above – it increases rapidly
 168 at the start, it reaches a peak at an intermediate point during training and then decays back to near
 169 zero.

170 Next we calculate the average MI between a single hidden neuron and the target, $I(Y; H_i)$, which
 171 we refer to as Single-Neuron Mutual Information (SNMI), and show the results in Fig. 4c. As expected,
 172 at first neurons barely have any information about the target and early in training we see a quick
 173 uniform increase in SNMI.

174 Remarkably, at the transition point there is a split in the SNMI histogram, with around half of
 175 the neurons reverting back to low values of SNMI and the other half continuing to increase. At the
 176 whole network level we do not find any sign of this split, as shown by the monotonically decreasing
 177 error rate in Fig. 3. This is a seemingly counterintuitive finding – some neurons actually get *worse*
 178 at predicting the target as the network learns. We currently do not have a solid explanation for this
 179 phenomenon, although we believe it could be due to the effects of local minima or to the neurons
 180 relying more on synergistic interactions at the cost of SNMI, as suggested by the results below.

181 After exploring the behaviour of individual neurons, we now turn to PID and study the
 182 interactions between them when predicting the target. Since a full PID analysis of the whole network is
 183 intractable, we follow a procedure inspired by [35] to estimate the PI terms of the learnt representation:
 184 we sample pairs of neurons, calculate the PI terms for each of them, apply random permutation
 185 correction to each pair separately, and finally compute averages over all pairs. We present results
 186 obtained with I_{\min} following Sec. 2.2.1, but qualitatively identical results are obtained if the more
 187 modern measures in [28,36] are used.

188 We calculate synergy S , redundancy R and total unique information $U = U(Y; Z_1) + U(Y; Z_2)$,
 189 as well as their normalised versions calculated by dividing S , R or U by the joint mutual information
 190 $I(Y; \{H_1, H_2\})$. Results are depicted in Fig. 5, and error intervals shown correspond to two standard
 191 deviations across pairs.²

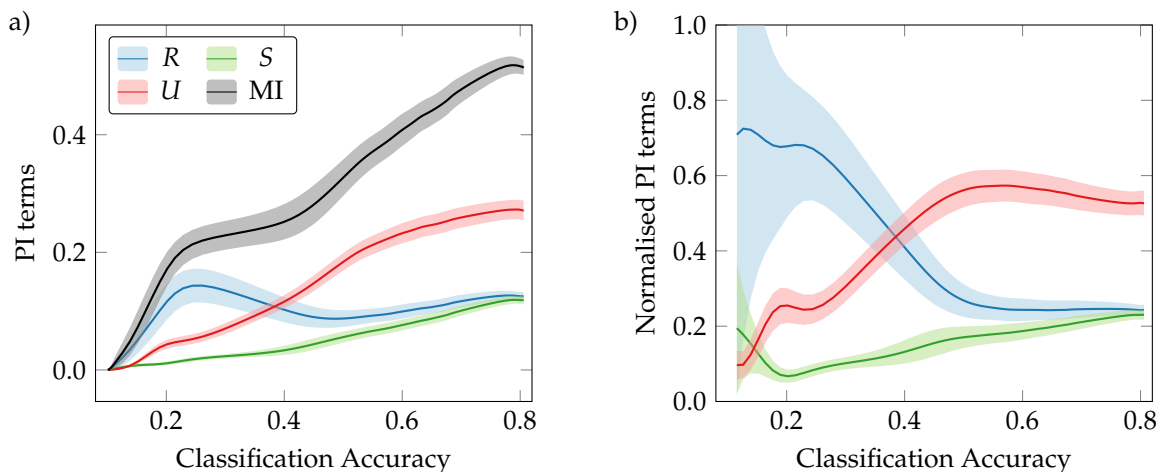


Figure 5. PI terms (left) and PI terms normalised by joint mutual information (right). Mutual information in black, redundancy in blue, synergy in green and unique information in red. MI increases consistently during training, but the PI terms reveal a transition between a redundancy-dominated phase and a unique-information-dominated phase.

192 Here we see again a transition between two phases of learning. Although synergy and joint
 193 MI increase steadily at all times, there is a clear distinction between a first phase dominated by
 194 redundancy and a second one dominated by unique information. It is at this point that neurons
 195 specialise, suggesting that this is when disentangled representations emerge.

196 These three phenomena (peak in PWMI, split SNMI histogram and redundant-unique information
 197 transition) do not happen at the same time. In the figures shown, the peak in PWMI marks the onset
 198 of the decline of redundancy, and the split in SNMI happens between then and the point when
 199 redundancy is overtaken by unique information. This is, however, a consistent pattern we have
 200 observed in networks of multiple sizes, and in bigger networks these three events tend to come closer
 201 in time (results not shown).

² The number of surrogates was increased until their SEM was much lower than STD across samples.

202 We note that there is a relation between PWMI and R and between SNMI and U . As indicated in
 203 Eq. 7, SNMI is an upper bound on that neuron's unique information; and usually higher PWMI comes
 204 with higher redundancy between the neurons. However, although they follow similar shapes, these
 205 magnitudes do not quantify the same thing. Take the OR logical gate as an example – if we feed it a
 206 uniform distribution of all possible inputs (00, 01, 10, 11) both input bits will be perfectly uncorrelated,
 207 yet their redundancy (according to I_{\min}) will be nonzero.

208 These findings are in line with those of [Schwartz-Ziv and Tishby](#), who observe a similar transition
 209 in a feed-forward neural network classifier. One of the pieces of evidence for [Schwartz-Ziv and](#)
 210 [Tishby](#)'s claim is in the change of gradient signal dynamics from a drift to a diffusion regime. We
 211 did not analyse gradient dynamics as part of this study, but investigating the relationship between
 212 informational and dynamical accounts of learning is certainly a promising topic.

213 3.3. Neural interactions

214 In this last set of experiments we examine the representations learnt jointly by larger groups of
 215 neurons. Due to computational constraints, we run the analyses only on fully trained networks instead
 216 of at multiple points during training. We train networks of different sizes, ranging from 20 to 500
 217 hidden neurons (using the same algorithm, but allowing each network to train for more epochs until
 218 convergence), and consider larger groups of neurons for the PID analysis. We use a procedure similar
 219 to the one used in the previous section, but this time sampling tuples of K neurons, and calculating
 220 their overall synergy following Eq. 9. We refer to this as the PI- K synergy. Results are shown in Fig. 6.

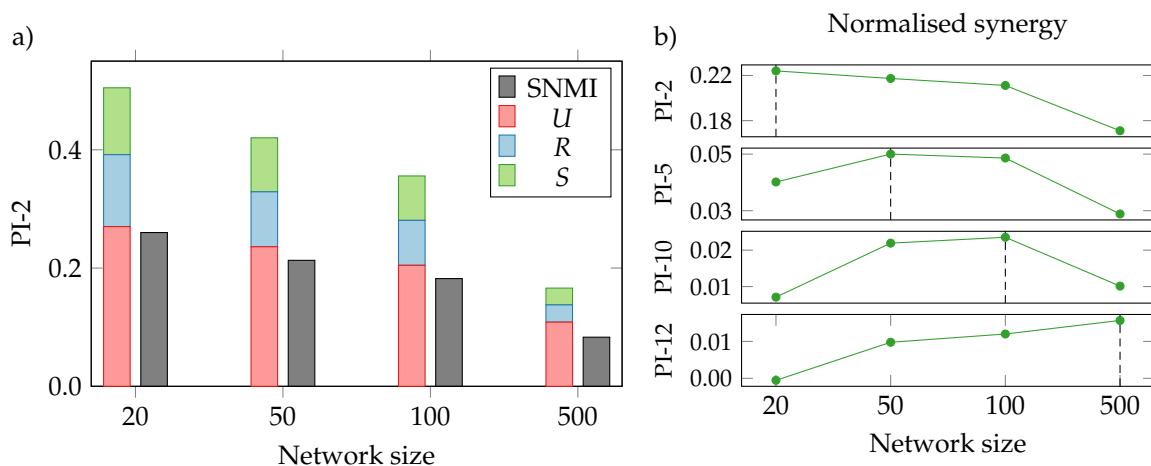


Figure 6. PID analysis of larger groups of neurons in networks of different sizes. (a) Single-neuron MI is consistently smaller in bigger networks, indicating that, although the network as a whole is a better classifier, each individual neuron has a less efficient encoding. (b) Normalised PI- K synergy, with network size increasing from left to right and K from top to bottom. Network with maximum synergy for each PI- K highlighted with a vertical dashed line. The PI group size with the highest synergy becomes larger in larger networks, indicating that in bigger networks one needs to consider larger groups to capture strong cooperative interactions.

221 The first result in Fig. 6a is that SNMI decreases consistently with network size. This represents
 222 reduced efficiency in the neurons' compression – despite the overall accuracy of the network being
 223 significantly higher for bigger networks ($\sim 20\%$ error rate for a network with 20 hidden neurons, vs.
 224 $\sim 5\%$ for a network with 500), each individual neuron contains less information about the target. This
 225 suggests that the representation is more distributed in bigger networks, as emphasised below.

226 What is somewhat counterintuitive is that normalised unique information actually grows in
 227 bigger networks, which is apparently in contradiction with more distributed representations. These

two, however, are perfectly compatible – bigger networks have more and less correlated neurons, and despite U growing relative to S and R , it still decreases significantly in absolute terms³

Interestingly, Fig. 6b shows that the network size that achieves maximum normalised synergy shifts to the right as we inspect larger groups of neurons. For bigger networks, bigger groups carry more synergistic information, meaning that representations also become more distributed. There is a consistent pattern that in bigger networks we need to explore increasingly high neural groups to see any meaningful PI values, which means that perhaps part of the success of bigger networks is that they make better use of higher-order correlations between hidden neurons. This can be seen as a signature of bigger networks achieving richer and more complex representations [37].

3.4. Limitations

The main limitation of the vanilla PID formulation is that the number of PI terms scales very rapidly with bigger group sizes – the number of terms in the PI decomposition of a system with n sources is the $(n - 1)$ -st Dedekind number, which is 7579 terms for a 5-variable system and has not been computed yet for systems of more than 8 variables. For this reason we have restricted our analyses mostly to pairs of neurons, although in practice we expect larger groups of neurons to have strong effects on the prediction. Potentially some approximation to the whole PID or a reasonable grouping of PI terms could help scale this type of analysis to larger systems.

On a separate topic, some of the phenomena of interest we have described in this article (two phases of learning, peak in correlation between the neurons) happen very early on during training, in practice. In a real-world ML setting, most of the time is spent in the last phase where error decreases very slowly; and so far we haven't seen any unusual behaviour in that region. Future work should focus on this second phase and try to characterise it in more detail, with the aim of improving performance or speeding convergence.

4. Conclusion

In this article we have used Information Theory, and in particular the Partial Information Decomposition framework, to explore the latent representations learnt by a restricted Boltzmann machine. We have found that the learning process of neural generative models has two distinct phases: a first phase dominated by redundant information about the target, and another phase in which neurons specialise and each of them learns unique information about the target and synergy. This is in line with the findings of [Schwartz-Ziv and Tishby](#) in feed-forward networks, and we believe further research should explore the differences between generative and discriminative models in this regard.

Additionally, we found that representations learnt by bigger networks are more distributed, yet significantly less efficient at the single-neuron level. This suggests that the learning pressure of having fewer neurons encourages those neurons to specialise more, and therefore yields more disentangled representations. The interesting challenge is to find a principled way of encouraging networks towards disentangled representations while preserving performance.

An interesting piece of follow-up work would be to investigate whether these findings generalise to other deep generative models, most notably variational autoencoders [38]. We believe that further theoretical study of these learning systems is necessary to help us understand, interpret and improve them.

Acknowledgments: The authors would like to thank Raúl Vicente for insightful discussions, and Pietro Marchesi for valuable feedback in the early stages of this project.

³ Note that the U term plotted in Fig. 6 is the sum of the unique information of both neurons in the pair. Naturally, for one neuron $U(Y; H_i) \leq I(Y; H_i)$ as per Eq. 7.

270 **References**

- 271 1. Srivastava, N.; Hinton, G.; Krizhevsky, A.; Sutskever, I.; Salakhutdinov, R. Dropout: A Simple Way to
272 Prevent Neural Networks from Overfitting. *Journal of Machine Learning Research* **2014**, *15*, 1929–1958.
- 273 2. Gal, Y.; Ghahramani, Z. Dropout as a Bayesian Approximation: Representing Model Uncertainty in Deep
274 Learning **2015**. [[1506.02142](#)].
- 275 3. Bengio, Y.; Courville, A.; Vincent, P. Representation Learning: A Review and New Perspectives **2012**.
276 [[1206.5538](#)].
- 277 4. Higgins, I.; Matthey, L.; Glorot, X.; Pal, A.; Uria, B.; Blundell, C.; Mohamed, S.; Lerchner, A. Early Visual
278 Concept Learning with Unsupervised Deep Learning **2016**. [[1606.05579](#)].
- 279 5. Mathieu, M.; Zhao, J.; Sprechmann, P.; Ramesh, A.; LeCun, Y. Disentangling Factors of Variation in Deep
280 Representations Using Adversarial Training **2016**. [[1611.03383](#)].
- 281 6. Siddharth, N.; Paige, B.; Van de Meent, J.W.; Desmaison, A.; Wood, F.; Goodman, N.D.; Kohli, P.; Torr, P.H.S.
282 Learning Disentangled Representations with Semi-Supervised Deep Generative Models **2017**. [[1706.00400](#)].
- 283 7. Lake, B.M.; Ullman, T.D.; Tenenbaum, J.B.; Gershman, S.J. Building Machines That Learn and Think Like
284 People **2016**. [[1604.00289](#)].
- 285 8. Garnelo, M.; Arulkumaran, K.; Shanahan, M. Towards Deep Symbolic Reinforcement Learning **2016**.
286 [[1609.05518](#)].
- 287 9. Williams, P.L.; Beer, R.D. Nonnegative Decomposition of Multivariate Information. *arXiv preprint*
288 *arXiv:1004.2515* **2010**.
- 289 10. Rieke, F.; Bialek, W.; Warland, D.; de Ruyter van Steveninck, R. *Spikes: Exploring the Neural Code*; MIT Press,
290 1997; p. 395.
- 291 11. Le, Q.V.; Ranzato, M.; Monga, R.; Devin, M.; Chen, K.; Corrado, G.S.; Dean, J.; Ng, A.Y. Building High-Level
292 Features Using Large Scale Unsupervised Learning **2011**. [[1112.6209](#)].
- 293 12. Zeiler, M.D.; Fergus, R. Visualizing and Understanding Convolutional Networks. European Conference
294 on Computer Vision. Springer, 2014, pp. 818–833.
- 295 13. Choromanska, A.; Henaff, M.; Mathieu, M.; Arous, G.B.; LeCun, Y. The Loss Surfaces of Multilayer
296 Networks **2014**. [[1412.0233](#)].
- 297 14. Kawaguchi, K. Deep Learning Without Poor Local Minima **2016**. [[1605.07110](#)].
- 298 15. Sørngård, B. Information Theory for Analyzing Neural Networks. Master’s thesis, Norwegian University
299 of Science and Technology, 2014.
- 300 16. Schwartz-Ziv, R.; Tishby, N. Opening the Black Box of Deep Neural Networks via Information **2017**.
301 [[1703.00810](#)].
- 302 17. Achille, A.; Soatto, S. On the Emergence of Invariance and Disentangling in Deep Representations **2017**.
303 [[1706.01350](#)].
- 304 18. Tishby, N.; Zaslavsky, N. Deep Learning and the Information Bottleneck Principle **2015**. [[1503.02406](#)].
- 305 19. Berglund, M.; Raiko, T.; Cho, K. Measuring the Usefulness of Hidden Units in Boltzmann Machines With
306 Mutual Information. *Neural Networks* **2015**, *64*, 12–18.
- 307 20. Balduzzi, D.; Frean, M.; Leary, L.; Lewis, J.; Ma, K.W.D.; McWilliams, B. The Shattered Gradients Problem:
308 If Resnets are the Answer, Then What is the Question? **2017**. [[1702.08591](#)].
- 309 21. Hinton, G.E.; van Camp, D. Keeping the Neural Networks Simple by Minimizing the Description Length
310 of the Weights. Proceedings of the Sixth Annual Conference on Computational Learning Theory; ACM:
311 New York, NY, USA, 1993; COLT ’93, pp. 5–13.
- 312 22. Smolensky, P. Information Processing in Dynamical Systems: Foundations of Harmony Theory. Technical
313 report, DTIC Document, 1986.
- 314 23. Larochelle, H.; Bengio, Y. Classification Using Discriminative Restricted Boltzmann Machines. Proceedings
315 of the 25th International Conference on Machine Learning, 2008, pp. 536–543.
- 316 24. Hinton, G.E.; Osindero, S.; Teh, Y.W. A Fast Learning Algorithm for Deep Belief Nets. *Neural computation*
317 **2006**, *18*, 1527–1554.
- 318 25. Tieleman, T. Training Restricted Boltzmann Machines Using Approximations to the Likelihood Gradient.
319 Proceedings of the 25th International Conference on Machine Learning; ACM Press: New York, New York,
320 USA, 2008; pp. 1064–1071.
- 321 26. Cover, T.M.; Thomas, J.A. *Elements of Information Theory*; Wiley: Hoboken, 2006.

- 322 27. DeWeese, M.R.; Meister, M. How to Measure the Information Gained from one Symbol. *Network*
323 *Computation in Neural Systems* **1999**, *12*, 325–340.
- 324 28. Ince, R.A.A. Measuring Multivariate Redundant Information with Pointwise Common Change in Surprisal.
325 *Entropy* **2017**, *19*, 318, [1602.05063].
- 326 29. Griffith, V.; Ho, T. Quantifying Redundant Information in Predicting a Target Random Variable. *Entropy*
327 **2015**, *17*, 4644–4653.
- 328 30. Harder, M.; Salge, C.; Polani, D. Bivariate Measure of Redundant Information. *Physical Review E* **2013**,
329 *87*, 012130.
- 330 31. Bertschinger, N.; Rauh, J.; Olbrich, E.; Jost, J. Shared Information—New Insights and Problems in
331 Decomposing Information in Complex Systems. Proceedings of the European Conference on Complex
332 Systems 2012. Springer, 2013, pp. 251–269.
- 333 32. Williams, P.L. Information Dynamics: Its Theory and Application to Embodied Cognitive Systems. PhD
334 thesis, Indiana University, 2013.
- 335 33. Williams, P.L. Information Dynamics: Its Theory and Application to Embodied Cognitive Systems. Phd
336 thesis, Indiana University, 2011.
- 337 34. Lizier, J.T. *The Local Information Dynamics of Distributed Computation in Complex Systems*; Springer Theses,
338 Springer Berlin Heidelberg: Berlin, Heidelberg, 2010.
- 339 35. Timme, N.; Alford, W.; Flecker, B.; Beggs, J.M. Synergy, Redundancy, and Multivariate Information
340 Measures: an Experimentalist's Perspective. *Journal of computational neuroscience* **2014**, *36*, 119–140.
- 341 36. Bertschinger, N.; Rauh, J.; Olbrich, E.; Jost, J.; Ay, N. Quantifying Unique Information. *Entropy* **2014**,
342 *16*, 2161–2183.
- 343 37. Montúfar, G.; Ay, N.; Ghazi-Zahedi, K. Geometry and Expressive Power of Conditional Restricted
344 Boltzmann Machines. *Journal of Machine Learning Research* **2015**, *16*, 2405–2436.
- 345 38. Kingma, D.P.; Welling, M. Auto-Encoding Variational Bayes **2013**. [1312.6114].

346 © 2017 by the authors. Submitted to *Entropy* for possible open access publication under the terms and conditions
347 of the Creative Commons Attribution (CC BY) license (<http://creativecommons.org/licenses/by/4.0/>).