

Andrius Grabauskas

Measuring mutual information in Neural Networks

Computer Science Tripos – Part II

Robinson College

Wednesday 15th May, 2019

Declaration

I, Andrius Grabauskas of Robinson College, being a candidate for Part II of the Computer Science Tripos, hereby declare that this dissertation and the work described in it are my own work, unaided except as may be specified below, and that the dissertation does not contain material that has already been used to any substantial extent for a comparable purpose.

Signed

Date

Proforma

Name: **Andrius Grabauskas**
College: **Robinson College**
Project Title: **Measuring mutual information in Neural Networks**
Examination: **Computer Science Tripos – Part II, July 2001**
Word Count: **6208¹**
Line Count: **1303²**
Project Originator: **Dr. Damon Wischik**
Supervisor: **Prof. Alan Mycroft**

Original Aims of the Project

Work Completed

Special Difficulties

¹This word count was computed by `detex *.tex | tr -cd '0-9A-Za-z \n' | wc -w`

²This line count was computed by `wc -l **/*.py`

Contents

1	Introduction	11
1.1	The Information plane	12
2	Preparation	15
2.1	Entropy and Mutual Information	15
2.1.1	Entropy	15
2.1.2	Conditional Entropy	16
2.1.3	Mutual Information	16
2.2	Neural Networks	17
2.2.1	The Prediction problem	17
2.2.2	Machine Learning Frameworks	18
2.2.3	Neural Networks	18
2.2.4	Abstracting the Neural Network	19
2.3	The Information plane	20
2.3.1	Setup	21
2.3.2	Visualization	22
2.3.3	Interpretation of the Information Plane	23
2.4	Testability	24
2.5	Software Engineering	24
2.6	Starting Point	24
3	Implementation	25
3.1	The Compression Rebuttal	25
3.1.1	Viability of compression	25
3.1.2	Determinism of the Transition Function	26
3.1.3	Recap	28
3.2	Measuring Mutual Information within NN	28
3.2.1	General Algorithm	28
3.2.2	Hyperparameters	29
3.2.3	Mutual Information Estimators	30
3.2.4	MIE - Binning	31
3.2.5	MIE - Kernel Density Estimation	34
3.2.6	Discrete vs Continuous Distributions	35
3.2.7	MIE - As If Random	37
3.2.8	AIR	37
3.2.9	Binning method	38
3.2.10	Kernel Density Estimation	39

3.2.11	Advanced methods	39
3.3	Implementation Optimizations	40
3.3.1	Maximising Resource Utilization	40
3.3.2	Minimising the Workload	40
3.4	Repository Structure	44
4	Evaluation	45
4.0.1	Deterministic networks	46
4.0.2	Why Randomness is hard to capture	46
5	Conclusion	47
	Bibliography	47
A	Project Proposal	51

List of Figures

1.1	information plane for a neural network with 5 layers, which was only trained for one epoch.	12
1.2	Information plane for a neural network with 4 layers, which was trained for approximately 10 000 epochs	13
2.1	Source: Wikimedia Commons	19
2.2	Visualization of a neural networks structure. x here is any input to the network from the set $\{x_i, \dots, x_N\}$. \hat{y} is the prediction of the network for the input x which may or may not be to the correct label. The values t_1, \dots, t_L here are activations of layers $1, \dots, L$	20
2.3	Definition of correlated random variables $X, T_{e,l}$ and, Y . The Probability distributions are generated from the dataset $D = \{(x_i, y_i) i = 1, \dots, N\}$. $F_{\theta(e)}^l$ is defined by Equation 2.18	22
2.4	explaining IP	23
3.1	The general algorithm for calculating mutual information inside a neural network.	29
3.2	Updated definition of correlated random variables $X, T_{e,l}$ and, Y . The Probability distributions are generated from the dataset $D = \{(x_i, y_i) i = 1, \dots, N\}$. $F_{\theta(e)}^l$ is defined by Equation 2.18	30
3.3	prob	31
3.4	32
3.5	32
3.6	Implementation of binning MIE. Functions estimateX and estimateY are called by Algorithm from Figure 3.1.	33
3.7	Implementation of KDE MIE. Functions estimateX and estimateY are called by Algorithm from Figure 3.1.	36
3.8	Pseudo code for computing mutual information refer to Figure 3.9 and Figure 3.10 for entropy computation.	38
3.9	Algorithm for computing entropy – Binning method	38
3.10	Algorithm for computing conditional entropy	38
3.11	Algorithm for computing entropy – Kernel Density Estimation method. The same algorithm as used by Saxe.	39
3.12	Delta Skip Exact. The skip value is assumed to be global it specifies how many epochs to skip until we measure again	42
3.13	Backtrack Algorithm	42
3.14	Example how distance between two epochs is measured.	43
3.15	Delta Skip Approximate	44

Chapter 1

Introduction

Deep Neural Networks (DNNs) are an extremely successful tool, they are widely adopted commercially and closely studied academically, however even given the attention they have gathered there is no comprehensive understanding of how these models generalize data and provide such impressive performance - in fact very little is known about how DNNs learn or about their inner workings. Recently Prof. Tishby produced a paper claiming to understand the basic principle of how DNNs work. He suggested that there are two phases that the network goes through while being trained - the fitting phase and the compression phase. During the fitting phase the network memorizes the training data and makes predictions based on what it has seen before, during the compression phase the network generalizes, it forgets the unnecessary information from the training data. Tishby suggested that the incredible performance that DNNs are able to achieve is due to this compression phase, and that this process of compression is a result of randomness inherent in Stochastic gradient descent. Tishby showed this by looking at DNNs through the information domain, most notably he used what is now called the information plane. The information plane summarizes how the information is flowing through the DNN, for every neuron layer the plane shows mutual information it has with the input data and the label data. In his experiments Tishby has concluded that every layer loses unnecessary information from the input data and tries to keep information of the label. Tishby made some interesting and significant claims about how DNNs work, however he did not provide a formal proof, his conclusions are based only on experimental evidence.

In our work we look at Tishby's claim that DNNs compress data and throw away unnecessary information about the input. We reimplement his experiments as a form of independent verification, showing that the results Tishby got are robust and are stable to parameter changes. We also take a look at a paper produced by Saxe that provides an opposing view to that of Tishby's. Saxe showed that compression that Tishby showed is only a result of Tishby's choice of activation function for the neural network. He showed that compression only happens when Tanh activation function is used and does not happen when ReLU is used.

Lastly, we think that the experiments presented in both papers don't fully align with the ideas that Tishby presented to us, specifically his idea that weights should be treated 'as if' they are random variables. Tishby's and Saxe's experiment deal with this 'as if' random notion quite crudely or try to sidestep it completely. As a result we devised an experiment that tries to capture this idea more explicitly, although it is still relatively crude and more work should be put into it in the future.

1.1 The Information plane

If we look at a neural network through the lens of information we can think of it as a form of a Markov chain. Where every node takes in a piece of data processes it and passes to the next node. For every such node $t \in T$ we can measure its mutual information with the input patterns $x \in X$ and labels $y \in Y$.

$$Y - X \rightarrow T_1 \rightarrow \dots \rightarrow T_i \rightarrow \hat{Y}$$

In this analogy a node corresponds to a layer in a non-convolutional neural network. When looking at a neural network in this way we can say that its job is to preserve as much information as possible about the label, or to minimise information loss about the label during the transitions from node to node.

The Information plane summarizes this Markov chain view and applies it to the entirety of the training process.

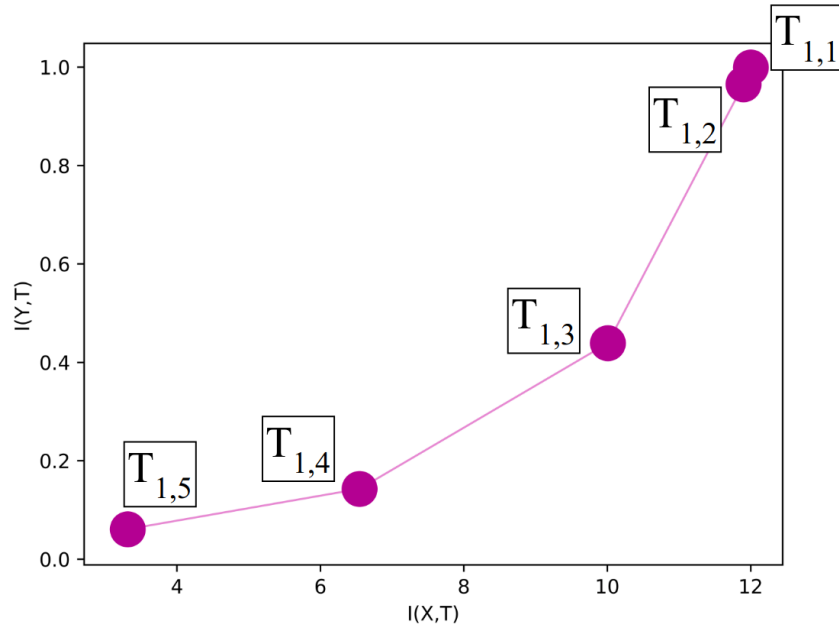


Figure 1.1: information plane for a neural network with 5 layers, which was only trained for one epoch.

Figure 1.1 shows an example of an information plane - in this case the network has 5 layers and has only been trained for one epoch. Every node corresponds to a layer, and the lines between them just help us distinguish the order of the layers. The x -axis shows mutual information between any layer T and the input X , while the y -axis shows mutual information with the label Y .

The figure below corresponds to a network that takes in a 12bit number and maps it to a 1bit number hence the mutual information ranges from 0 to 12 for the x -axis and from 0 to 1 for the y -axis.

The upper right most node corresponds with the very first layer before any data processing occurs, therefore it has maximum mutual information with the input and the label. As the network was only trained for one epoch the weights are essentially random and we see a steep drop in every subsequent layer, with the last layer of the network having close to 0 mutual

information with X and Y , which is what we expect if we assume that the network guesses at random before any training is done.

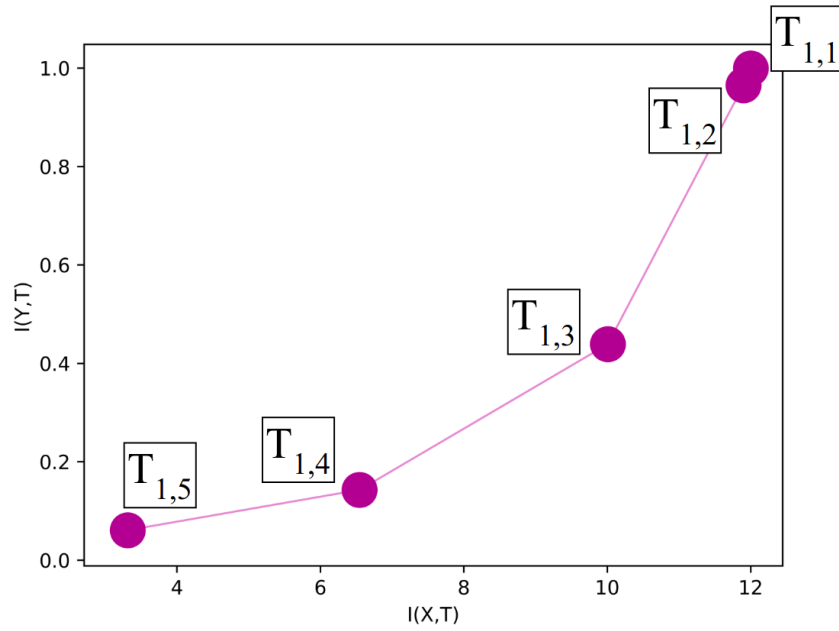


Figure 1.2: Information plane for a neural network with 4 layers, which was trained for approximately 10 000 epochs

Figure 1.2 shows an example of a full information plane. The gradient maps the colour to the epoch. From this figure we can clearly see the fitting and the compression phases Tishby was describing.

- The Fitting phase is the rapid improvement phase at the start of the training period, where information about the label is rapidly increasing while information about the input remains approximately the same. In this case the training is in the fitting phase for less than 1000 epochs.
- The Compression phase begins when the after the fitting phase ends and the rapid improvement stops. We can see during this phase we start losing information about the input while still gaining information about the label. This phase the was majority of the training time.

Chapter 2

Preparation

Neural Networks (NN) are an extremely successful tool, they are widely adopted commercially and closely studied academically, however even given the attention they have there is no comprehensive understanding of how these models generalize data and provide such impressive performance – in fact very little is known about how NN learn or about their inner workings. Recently Prof. Tishby [8] produced a paper claiming to understand the basic principles of how NN work. He decided to examine NNs through the information domain. Tishby made the claim that the incredible performance of NNs is due to their ability to compress information. Compressing data means the network is only able to keep relevant input features and it must discard the irrelevant bits of information, leading to ability to generalize.

Tishby made interesting claims and provided experimental evidence to support his claims, however he did not provide a formal proof leaving his results up for debate. A paper release by Saxe [4] has contested the claim made by Tishby arguing that compression cannot happen in Neural Networks and the results are a consequence of the hyper parameters Tishby used. However, Saxe’s paper suffers from the same problem as Tishby’s as it does not provide a formal proof only experimental evidence – as such it does not settle the rebuttal. To fully understand the discussion we need to understand the following topics Entropy and Mutual Information, Neural Networks, and Information Plane, described in section 2.1, section 2.2, and section 2.3.

2.1 Entropy and Mutual Information

2.1.1 Entropy

Entropy – quantifies information content of a random variable. It is generally measured in bits and can be thought of as the expected information content when we sample a random variable once. Let X be a discrete random variable that can take values in $\{x_1, \dots, x_n\}$. $H(X)$, the entropy of X , is defined by Equation 2.1.

$$H(X) = - \sum_{i=1}^n P(x_i) \log P(x_i) \quad (2.1)$$

Consider a random variable Y s.t $P(Y = 1) = P(Y = 0) = 0.5$, Equation 2.1 defines $H(Y)$ to be 1.

Similarly for a random variable Y s.t $P(Y = 0) = 0.5, P(Y = 1) = P(Y = 2) = 0.25$, we have $H(Y) = 1.5$.

2.1.2 Conditional Entropy

Conditional Entropy – quantifies the amount of information needed to describe an outcome of variable Y given that value of another random variable X is already known. Conditional Entropy of Y given X is written as $H(Y|X)$. Let X be defined as before, Let Y be a discrete random variable that can take values in $\{y_i, \dots, y_m\}$. Equations 2.2 and 2.3 are two examples how we can compute conditional entropies – using explicit probabilities or entropies of or entropies of random variables conditioned on an event.

$$H(Y|X) = - \sum_{i=1}^n \sum_{j=1}^m P(x_i, y_j) \log \frac{P(x_i, y_j)}{P(x_i)} \quad (2.2)$$

$$H(Y|X) = - \sum_{i=1}^n P(x_i) H(Y|X = x_i) \quad (2.3)$$

$$\text{where } H(Y|X = x) = - \sum_{i=1}^m P(y_i|X = x) \log P(y_i|X = x) \quad (2.4)$$

Let the correlated variables X and Y be defined by Table 2.1.

X \ Y	0	1
	0	1
0	0.25	0.25
1	0.5	0

Table 2.1: Joint probability distribution for X and Y

Equation 2.1 and Equation 2.2 defines entropy values to be:

$$\begin{aligned} H(Y|X) &= 0.5 \\ H(X|Y) &\approx 0.6887 \\ H(X) &= 1 \\ H(Y) &\approx 0.8112 \end{aligned} \quad (2.5)$$

2.1.3 Mutual Information

Mutual Information (MI) – measures how much information two random variables have in common. It quantifies information gained about one variable when observing the other. Equations 2.6 and 2.7 are two examples of how we can compute MI – using explicit probability computations or entropies of the random variables respectively, here X and Y are as previously defined.

$$I(X, Y) = \sum_{i=1}^n \sum_{j=1}^m P(x_i, y_j) \log \left(\frac{P(x_i, y_j)}{P(x_i) P(y_j)} \right) \quad (2.6)$$

$$I(X, Y) = H(X) - H(X|Y) \quad (2.7)$$

For example of MI consider the random variables X and Y as before in the conditional entropy section – defined by Table 2.1.

We computed the entropy values in Equation 2.5, we will use them in Equation 2.7 to compute $I(X, Y)$.

$$I(X, Y) = H(X) - H(X|Y) \approx 1 - 0.6887 = 0.3113 \quad (2.8)$$

Properties of Mutual Information

There are some important properties of MI that we need to take note of. Let X and Y be any probability distributions then:

Commutative

$$I(X, Y) = I(Y, X) \quad (2.9)$$

Information Loss MI of two random variables cannot exceed entropy of either of them.

$$\begin{aligned} H(X) &\geq I(X, Y) \\ H(Y) &\geq I(X, Y) \end{aligned} \quad (2.10)$$

Data Processing Inequality Let u be some function; then,

$$I(X, Y) \geq I(X, u(Y)) \quad (2.11)$$

Invertible Transformation Let u be some invertible function; then,

$$I(X, Y) = I(X, u(Y)) \quad (2.12)$$

2.2 Neural Networks

Before we understand neural networks we need to understand The Prediction Problem and the purpose of Machine Learning Frameworks.

2.2.1 The Prediction problem

Suppose we have some dataset D defined as (x_i, y_i) for $i = 1, \dots, N$. The prediction problem is finding a function f s.t Equation 2.13 is satisfied.

$$f(x_i) = y_i \text{ for } i = 1, \dots, N \quad (2.13)$$

Prediction task is a common task that involves having input data $\{x_i, \dots, x_N\}$ and finding the label, some desirable feature, $\{y_i, \dots, y_C\}$.

The prediction problem could be simple to extract: for example if our input is a natural number $x_i \in \mathbb{N}$, and our label is either *true* or *false* depending if x is even or odd – in which case function defined by Equation 2.14 satisfies the problem.

$$g(x) = \begin{cases} True, & \text{if } \exists n \in \mathbb{N}. x = 2n, \\ False, & \text{otherwise.} \end{cases} \quad (2.14)$$

The prediction problem also can be impossible to solve: for example the halting problem, if our x 's are programs and y 's boolean values corresponding if the program halts or not.

Of course the prediction problem can be hard or impossible to solve as is the case for problems:

input data	label	difficulty
medical symptoms	diagnosis	intractable
picture	object in the picture	
face photograph	identity	
stock market history	future stock prices	
program	does the program halt	proved to be unsolvable
boolean equation	is the equation satisfiable	expensive to compute

Table 2.2: Example of specific prediction problems

Problems listed in Table 2.2 are either intractable, unsolvable or too expensive to compute – hence we cannot produce an algorithm that always give the correct answer and runs in a reasonable time.

2.2.2 Machine Learning Frameworks

If the Prediction problem is too difficult and we are tolerant to errors in our labels we may want to use a supervised machine learning framework¹ to tackle the problem.

Every machine learning framework requires that we have some subset $\hat{X} \subseteq \{x_i, \dots, x_N\}$ s.t that $\forall x \in \hat{X}$ we know label the y . A framework uses the data in order to reach some goal – such as minimizing the prediction error. The way any machine learning framework learns from data varies, but generally more data means an increase in prediction performance.

2.2.3 Neural Networks

Neural Networks (NN) are an example of a machine learning framework. They learn from data and attempt to solve the prediction problem.

The structure of a NN consists of layers of nodes, where every consecutive layer is fully connected as in Figure 2.1a. When we try to predict a label of a specific input every node gets assigned a value in the real number space \mathbb{R} . Values for the Input nodes are provided, whereas values for the nodes in non-input layers are generated based on nodes from the preceding layer. Let $n_{l,i}$ be the value that the i 'th node in layer l takes. Let $w_{l,j,i}$ be a parameter that influences how relevant the j 'th node in layer l is to the i 'th node in layer $l + 1$.

$$n_{l+1,i} = g\left(\sum_{j=0}^{\text{layer } l \text{ size}} w_{l,j,i} * n_{l,j}\right) \quad (2.15)$$

g in Equation 2.15 is called the activation function and is generally taken to be:

- Leaky ReLu:

$$g(x) = \begin{cases} x, & \text{if } x > 0 \\ \alpha x, & \text{otherwise.} \end{cases}$$

α here is some small value such as $\alpha = 0.01$

¹In this thesis we are exclusively talking about supervised machine learning techniques – the word "supervised" will be omitted for brevity

- Sigmoid: $g(x) = \frac{1}{1+e^x}$
- Tanh: $g(x) = \tanh(x)$

Stochastic Gradient Descent (SGD) The NNs training algorithm is called Stochastic Gradient Descent. SGD periodically adjusts the parameters w according to some goal such as minimising the prediction error.

SGD is an iterative process – it introduces a notion of epochs where one iteration of the algorithm advances the epochs by one. Notion of epochs implies that the parameters w depends on which epoch we are currently at – thus it must take the epoch number e as an argument – $w(e)$.

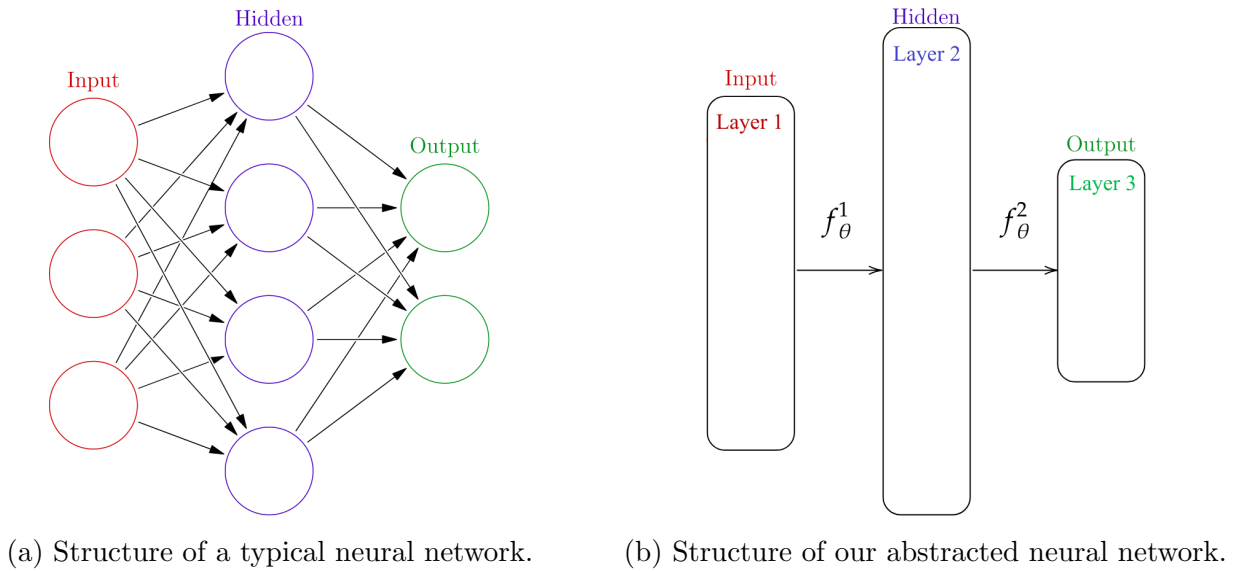


Figure 2.1: Source: Wikimedia Commons

2.2.4 Abstracting the Neural Network

For our purposes we will abstract away the individual nodes in the NN and only consider the layer structure – as in Figure 2.1b.

In our abstraction the NN is structured in layers where every layer holds an intermediate representation of the final prediction output - let us call this intermediate representation an **activation** of that specific layer. We will formally define a neural network with L layers to be a sequence of L functions $f_{\theta(e)}^1, f_{\theta(e)}^2, \dots, f_{\theta(e)}^L$ that are parameterized by θ s.t. Equations 2.16 hold. In our abstract the parameters θ correspond to the parameters w , therefore θ depends on the current epoch – $\theta(e)$. Let us also define an **activation** of layer l to be output of the function f_{θ}^l .

$$\begin{aligned}
& \text{let } t_0 = x, \\
& t_0 \rightarrow f_\theta^1(t_0) = t_1, \\
& t_1 \rightarrow f_\theta^2(t_1) = t_2, \\
& \dots \\
& t_{L-1} \rightarrow f_\theta^L(t_{L-1}) = t_L, \\
& \text{let } \hat{y} = t_L
\end{aligned} \tag{2.16}$$

values t_1, t_2, \dots, t_L here are **activations** of layers 1, 2, ..., L ,

x is any input to the NN from the set $\{x_i, \dots, x_N\}$,

\hat{y} is the prediction of the NN for the input, x which may or not be correct label,

the arrow ‘ \rightarrow ’ signifies a transition from one NN layer to another.

We now have one function for each layer transition of the NN this allows us to extract the activation of a specific layer as in Equation 2.17.

$$t_l = f_\theta^l(f_\theta^{l-1}(\dots(f_\theta^1(x)))) \tag{2.17}$$

activtion of later l for input x

Let us define F_θ^l to be the activation of layer l given input x i.e

$$F_\theta^l(x) = f_\theta^l(f_\theta^{l-1}(\dots(f_\theta^1(x)))) \tag{2.18}$$

Figure 2.2 summarizes our NN definition.

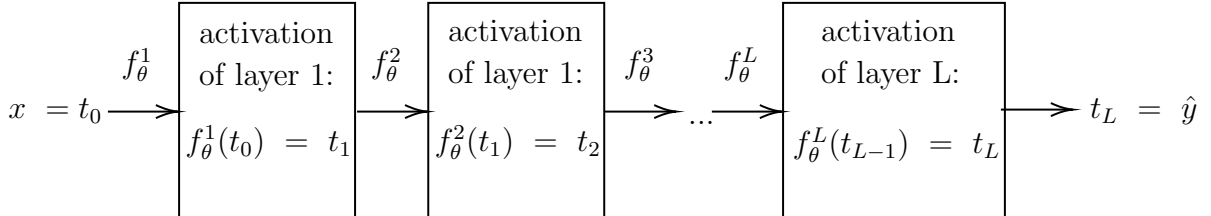


Figure 2.2: Visualization of a neural networks structure. x here is any input to the network from the set $\{x_i, \dots, x_N\}$. \hat{y} is the prediction of the network for the input x which may or may not be to the correct label. The values t_1, \dots, t_L here are **activations** of layers 1, ..., L .

2.3 The Information plane

The information plane is a way of visualizing the Neural Network's training process through the information domain. Meaning we are looking what information the NN is retaining throughout the layers, and how the information retained changes throughout the training process.

We are interested see what information about the input and the label is retained – as having this information would help us understand how the neural network learns. Suppose we have this data available then we can examine this data with respect to:

- Epochs – lets us ask questions such as:
 - Is the network getting rid of noise in the network.
 - Does the network ever retain all the information about the label.
 - If the network does retain all the information about the label does training it after the point lead to performance increase.
- Layer – lets us ask questions such as:
 - Is there a difference between how much information the layers are discarding.
 - How does the information dynamics change if we add more layers.

Suppose we are examining a specific NN. At the start of the training process we expect the network to perform poorly – meaning it does not retain information about the label. Through the training process we expect:

- The information about the label to rise until it saturates. If we train past this saturation point we don't expect the information about the label to change much.
- The information about the input to rise until we have saturated the information about the label. At this point we either expect:
 - The information about the input to keep rising if we believe the network starts to learn noise about the data.
 - The information about the input to start reducing until the networks does not retain any noise about the input if we believe the network generalizes the data.

Ideally a neural network would retain all the information relevant to the label and none of the noise inherent in the input.

2.3.1 Setup

We want to compute what information about the label and the input the network is retaining – that is we want to compute the mutual information:

- Between the input distribution and all of the network layers,
- Between the label distribution and all of the network layers.

In order to compute the needed MI values we need to define the probability distributions:

- X - probability distribution of the input,
- Y - probability distribution of the label,
- $T_{e,l}$ - probability distribution of the layer l for the epoch e .

The routine in Figure 2.3 defines the correlated random variables: $X, T_{e,l}, Y$. The probability distributions define the needed MI values: $I(X, T_{e,l})$ and $I(Y, T_{e,l})$.

```

1  def rxyty(e, l):
2      pick i ~ Uniform {1...N}
3      return (xi, Fθ(e)l(xi), yi)

```

Figure 2.3: Definition of correlated random variables $X, T_{e,l}$ and, Y . The Probability distributions are generated from the dataset $D = \{(x_i, y_i) | i = 1, \dots, N\}$. $F_{\theta(e)}^l$ is defined by Equation 2.18

Input Probability Distribution The probability distribution X is generated from the input dataset D . Any input dataset D has an inherent probability distribution but most of the time we don't know it exactly. Figure 2.3 makes the assumption that D is uniformly distributed and thus takes every input value to be equally likely. If we had a specific dataset we could adjust our assumption and assign a different probability distribution to the input dataset.

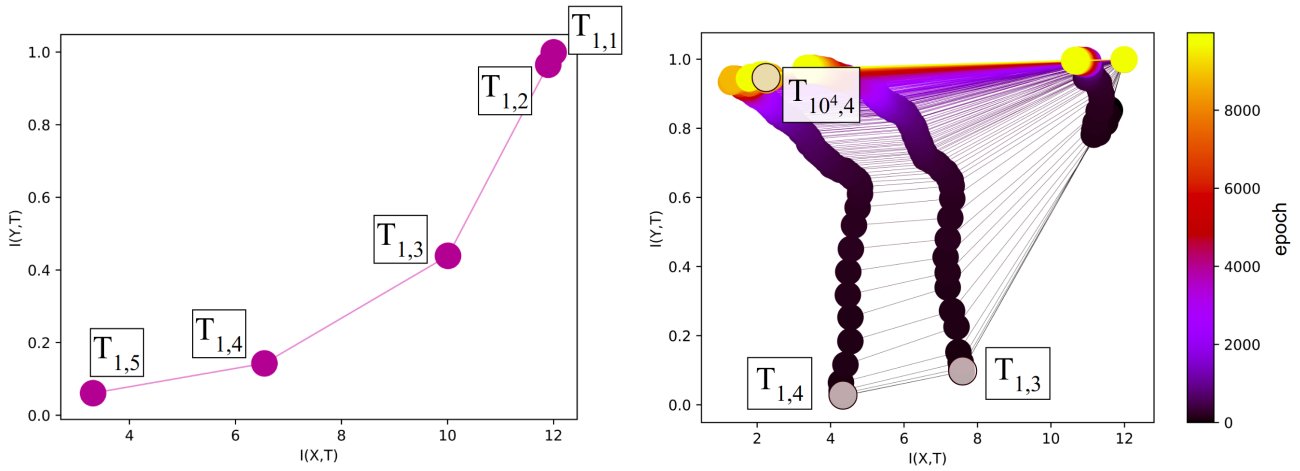
2.3.2 Visualization

The Information Plane visualizes how the information retained in the NN changes throughout the training period. In order to do this we need values $I(X, T_{e,l})$ and $I(Y, T_{e,l})$, for every epoch e and every layer l – we show how we obtain them in subsection 2.3.1.

To better understand how the Information Plane present the data let us consider **Figure 2.4a** – it shows the Information Plane for a NN for the first epoch before any training has been done. This implies that parameters $\theta(e)$ have not been affected by SGD and are thus random. If the network is random we would expect it not be able to retain much information. In the image we see that the network loses almost all information about the label $I(Y, T_{1,5}) \approx 0$, but retains some information about the input $I(X, T_{1,5}) \approx 3$ – we can see that this random NN only retains noise.

Let us now bring our attention to **Figure 2.4b** – it shows a NN that was trained for $\approx 10^4$. The NN was initialized with random parameters $\theta(1)$, then trained via the SGD algorithm which altered the parameters θ . The NN was trained so we expect to see a rise of mutual information with the label $I(Y, T_{e,l})$, throughout the training period – and we do see this in the Figure. However, we also see some other interesting features:

- Notice that the increase in $I(Y, T_{e,l})$ is very rapid at the start of the training period, but slows down considerably after ≈ 1000 epochs. The rapid increase might be the NN learning features that are highly correlated with the label, hence easy to learn.
- Notice how when increase in $I(Y, T_{e,l})$ slows down we start to see a decrease in the information about the input – $I(X, T_{e,l})$. This might be the point when the networks starts to remove the noise from the input.
- Notice how in both figures the NNs retain all the information about the label – This suggest that in order to achieve high performance retaining all the information is not enough. If all we needed to do is retain all the label information there would be no reason to add more layers to the neural network.



(a) information plane for a neural network with 5 layers, which was only trained for one epoch.

(b) Information plane for a neural network with 4 layers, which was trained for approximately 10 000 epochs.

Figure 2.4: The nodes in the figures correspond to information content of layers in the NN. The lines between the nodes help us distinguish individual epochs if more than one is plotted on the screen and lets us see the order of layers within a single epoch.

Notice that some nodes are labeled $T_{e,l}$, where e is the epoch number and l is the layer the node belongs to. Consider $T_{1,3}$ from Figure 2.4b – the node corresponds to the information content of layer 3 for the 1'st epoch, the x coordinate of the node is the value $I(X, T_{1,3}) \approx 8$, the y coordinate is the value $I(Y, T_{1,3}) \approx 0.1$

The Neural Networks in both figures have been trained on the same dataset as used by Tishby[8], The datasets input entropy, $H(X)$, is 12 and label entropy, $H(Y)$, is 1. Notice the x and y axis don't exceed the entropies $H(X)$ and $H(Y)$ respectively – this is due to the information loss property of MI Equation 2.10.

2.3.3 Interpretation of the Information Plane

Let us once again consider Figure 2.4b – we can see two general phases in the neural network. Tishby has named them *The Fitting Phase* and *The Compression Phase*.

These phases are observations made by Tishby and seem to be reproducible in by experiments. Even though Tishby gave no clear definition of how to identify the phases, they seem to follow the properties outlined below.

The Fitting Phase In Figure 2.4b the neural network is in the fitting phase from the start of the training up until epoch ~ 1500 . The duration of the fitting phase varies heavily on the training parameters and is most influenced by the size of our input dataset. The fitting phase is characterized by:

- A rapid increase in $I(Y, T_{e,l})$, the information about the label, as we advance through the epochs e . The increase is especially visible in the later layers, in our case layers 3 and 4.
- Either an increase or no change in $I(X, T_{e,l})$, the information about the input, as we advance through the epochs e , in our case we see very little change in $I(X, T_{e,l})$.

Tishby’s understanding of the Fitting phase is that the network tries to memorize the data and makes predictions based on the observations, this means that the network may learn features that only superficially correlate with the correct label.

The Compression Phase In Figure 2.4b the neural network enters the compression phase when the fitting phase ends around epoch ~ 1500 and lasts until we finish the training process. The compression phase is characterized by:

- A slowdown of how fast $I(Y, T_{e,l})$ is increasing with respect to epoch e .
- A slow decrease of $I(X, T_{e,l})$ with respect to epoch e .

Tishby’s understanding of the Compression phase is that the network learns how to compress the representation of the input. This means the network has to discard features that do not help with predicting the correct labels. Discarding irrelevant features helps the neural network to generalize and produce better prediction for new data.

2.4 Testability

Software related to NN tends to be very compact as it heavily relies on external libraries that are extensively tested. In our case the most effective way to test the software is to provide tools that help examine the data. I have provided such tools in form of data visualization – such as a tool to generate a video of the training process, which lets a human detect anomalies.

2.5 Software Engineering

The project was build using the Spiral model. I needed to build a working version as fast as possible in order to verify the validity of Tishby’s experiments. Once the Minimal Viable Product was build I added features as needed. Features such as different experiments, performance optimizations, visualization tools.

2.6 Starting Point

In order to understand the theory required I had learn the ideas presented in papers by Tishby[8][7] and Saxe[4]. I had to have a strong understanding of Neural Networks and Information Theory.

The starting for implementation was knowledge of *Python 3*. For this project I had to learn the Machine Learning Frameworks: *Keras* and *Tensorflow*.

Chapter 3

Implementation

3.1 The Compression Rebuttal

This section talks about if a NN can compress¹ data – i.e is it possible to lose information from layer to layer. There is an argument to be made for both sides. We can easily construct a NN that loses information. If we take the layer transition function f_{θ}^l to be

$$f_{\theta}^l(x) = 0$$

This function would lose information as we are not able to recover the input given the output, hence the NN must have discarded some information. Therefore compression is definitely possible. However, it is equally possible to create a function that preserves all the information about the input. Every layer in a NN is just a number \mathbb{R}^{d_l} , where d_l is the dimension of layer l . We know there is a one to one mapping between \mathbb{R} and \mathbb{R}^d for any dimension d . Therefore we can construct a neural network that loses no information if we take the layer transition function f_{θ}^l to be

$$f_{\theta}^l(x) = \text{map_}\mathbb{R}_\text{to_}\mathbb{R}^{d_l}(\text{map_}\mathbb{R}^{d_{l-1}}_\text{to_}\mathbb{R}(x))$$

That is mapping the value of previous layer to \mathbb{R} and mapping it back to \mathbb{R}^{d_l} . It is important to understand if NN can compress data if we are to study how NN manipulate information.

3.1.1 Viability of compression

Viability of compression inside a NN. Asking if a NN compresses data is equivalent to asking if Equation 3.1 holds.

$$I(X, T_{e,l}) < H(X) \tag{3.1}$$

Lets examine the value $I(X, T_{e,l})$ more closely. From Equation 2.7 we know,

$$I(X, T_{e,l}) = H(X) - H(X|T_{e,l}) \tag{3.2}$$

Form Equation 3.1 and Equation 3.2 we have,

$$I(X, T_{e,l}) < H(X) \Leftrightarrow H(X|T_{e,l}) > 0 \tag{3.3}$$

¹ N.B This section does not talk about the compression phase, just about compression in the network for a single epoch

Consider the value $H(X|T_{e,l})$. $H(X|T_{e,l})$ is 0 iff the value of $T_{e,l}$ uniquely identifies the value of X i.e

$$H(X|T_{e,l}) = 0 \Leftrightarrow P(X = x|T_{e,t} = t) = \begin{cases} 1, & \text{if } t = F_{\theta(e)}^l(x). \\ 0, & \text{otherwise.} \end{cases} \quad (3.4)$$

Equation 3.4 holds iff $F_{\theta}^l(x)$ is an invertible function. Recall the definition of $F_{\theta}^l(x)$

$$F_{\theta}^l(x) = f_{\theta}^l(f_{\theta}^{l-1}(\dots(f_{\theta}^1(x)))) \quad (3.5)$$

Equation 3.5 implies

$$F_{\theta}^l \text{ is invertible} \Leftrightarrow \forall i \in \{1 \dots l\}. f_{\theta}^i \text{ is invertible}$$

Putting it all together we get the bi implications

$$\begin{aligned} I(X, T_{e,l}) < H(X) &\Leftrightarrow \neg(\forall i \in \{1 \dots l\}. f_{\theta}^i \text{ is invertible}) \\ &\Leftrightarrow \exists i \in \{1 \dots l\}. f_{\theta}^i \text{ is not invertible} \\ &\Leftrightarrow \exists i \in \{1 \dots l\} \exists u, v \in \{x_1, \dots, x_N\}. f_{\theta}^i(u) = f_{\theta}^i(v) \wedge u \neq v \end{aligned} \quad (3.6)$$

and

This implies that compression can only happen if at least on transition function f_{θ}^l is not invertible.

3.1.2 Determinism of the Transition Function

Decomposing the Transition Function Lets break away from our NN abstraction and consider the transition function f_{θ}^l and how it is defined in concrete implementations of NNs – recall the Equation 2.15.

$$n_{l+1,i} = g\left(\sum_{j=0}^{\text{layer } l \text{ size}} w_{l,j,i} * n_{l,j}\right)$$

Where g is an invertible function such as Leaky ReLu, Sigmoid, Tanh. This means we can consider the function f_{θ}^l to be a composition of two different functions:

- A matrix multiplication – which is responsible for the weighted sum of previous layers activations $\sum_{j=0}^{\text{layer } l \text{ size}} w_{l,j,i} * n_{l,j}$. Let us call this matrix m_{θ}^l .
- An application of the invertible function g to every element of the vector produced by the matrix multiplication.

This function decomposition implies that

$$f_{\theta}^l \text{ is invertible} \Leftrightarrow m_{\theta}^l \text{ is invertible} \quad (3.7)$$

A Note on Invertible Matrices

Let M be a matrix and M^{-1} be the inverse, then;

$$M^{-1} \text{ is defined} \Leftrightarrow \det(M) \neq 0 \quad (3.8)$$

Let M be a random matrix, then;

$$P(\det(M) = 0) = 0 \quad (3.9)$$

Let M be random matrix, then from Equation 3.8 and Equation 3.9 we have

$$P(M^{-1} \text{ is defined}) = 1 \quad (3.10)$$

i.e every random matrix has an inverse.

Randomness in SGD Consider the matrix m_θ^l . The Matrix is defined by the parameters θ which are controlled by Stochastic Gradient Descent (SGD). At the start of the training period the parameters θ are initialized to random values. Every iteration of SGD we update the parameters – this update can also be considered random due to the Stochastic nature of the algorithm. From this we can treat the parameters θ and consequently the matrix m_θ^l as random throughout the training process.

Having m_θ^l be an instance of a random matrix would mean that Equation 3.10 holds – which would imply that every transition function is invertible, which in turn means that there cannot be any compression in a NN.

This would mean having discussion about Information in NN is moot. However, we can still have a meaningful discussion about neural networks if consider the parameters θ to be a random distribution rather than an instance of a concrete value. This is not agreed upon fully within the scientific community – Tishby assumes this is inherent in SGD, while Saxe has contested the claim.

Let us refer to the parameters as probability distribution with notation $\hat{\theta}$. $\hat{\theta}$ is a sequence of probability distributions that depend s.t $\hat{\theta}(e+1)$ depends on $\hat{\theta}(e)$, where e is the epoch.

Let us formally define $\hat{\theta}$. Since at the start of the training period SGD initializes parameters at random – let the start of the sequence $\hat{\theta}(1)$ be defined by Equation 3.11, where: \mathcal{N}_k – is the Multivariate Normal distribution with dimension k , μ – is a vector in k 'th dimension defining the mean of the distribution, Σ – is a $k * k$ matrix defining the variance of the distribution.

$$\hat{\theta}(1) \sim \mathcal{N}_k(\mu, \Sigma) \quad (3.11)$$

Let the sequence be defined as;

$$P(\hat{\theta}(e+1) = t | \hat{\theta}(e) = \hat{t}) = P(\phi = t), \quad (3.12)$$

$$\text{where } \phi \text{ is } \hat{\theta}(e) \text{ with one SGD iteration applied, we can assume } \phi \sim \mathcal{N}_k \quad (3.13)$$

We have shown that parameters of a NN can be thought as being a probability distribution.

Impact of $\hat{\theta}$ – parameters as probability distributions With the idea that parameters $\hat{\theta}$ are a probability distribution let us again consider the question of compression in neural networks. From subsection 3.1.1 we know that no compression in NN is equivalent to

$$I(X, T_{e,l}) = H(X)$$

which is equivalent

$$H(X|T_{e,l}) = 0$$

$H(X|T_{e,l}) = 0$ means by observing $T_{e,l}$ we can deduce the exact value of X . i.e

$$\begin{aligned} (\forall t. P(T_{e,l} = t) > 0) &\implies P(X = x|T_{e,l} = t) = 1) \\ \implies (\forall t. P(T_{e,l} = t) > 0) &\implies P(T_{e,l} = t|X \neq x) = 0) \end{aligned} \quad (3.14)$$

Notice from Equation 3.11 and Equation 3.13 that

$$\forall e. \hat{\theta}(e) \sim \mathcal{N}_k, \text{ where } e \text{ is an epoch} \quad (3.15)$$

Equations 3.14 and 3.15 are unsatisfiable together. Proof: Consider x and \hat{x} s.t $\hat{x} \neq x$. Let ϕ be s.t. $f_\phi^1(x) = t$

$$\begin{aligned} f_\phi^1(x) &= t, \\ \implies P(T_{e,1} = t) &> 0, \\ \implies P(T_{e,1} = t|X \neq x) &= 0, \text{ by Equation 3.14} \\ \implies P(T_{e,1} = t|X = \hat{x}) &= 0 \end{aligned} \quad (3.16)$$

however, we can construct $\hat{\phi}$ s.t $f_{\hat{\phi}}^1(\hat{x}) = t$

$$P(T_{e,1} = t|X = \hat{x}) = P(\hat{\theta}(e) = \hat{\phi}) > 0, \text{ by Equation 3.15} \quad (3.17)$$

hence, we get a contradiction

$$P(T_{e,1} = t|X = \hat{x}) = 0 \wedge P(T_{e,1} = t|X = \hat{x}) > 0 \quad (3.18)$$

and Equation 3.14 is unsatisfiable, this implies

$$H(X|T_{e,l}) > 0 \quad (3.19)$$

which finally implies

$$I(X, T_{e,l}) < H(X)$$

and proves that if we consider parameters to be random distributions the compression happens inside NN.

3.1.3 Recap

There is contention if NN are actually capable of compressing information. I have shown that if we assume the parameters of a NN to be concrete values compression is not possible and discussion about Information within them is moot. However, if we make the assumptions that parameters are actually random variables then compression does happen.

3.2 Measuring Mutual Information within NN

3.2.1 General Algorithm

We need a way to measure information content in NN throughout the training period – Figure 3.1 shows an algorithm that achieves this. It does this by explicitly running the SGD algorithm iteration by iteration and estimating information for every epoch and layer.

```

1  def informationNN(Data, Hyper, MIE):
2  # Measures how much information is retained in the NN. Aggregates the data
3  # over epochs and layers.
4  # Input:
5  # Data - Training data  $D = \{(x_i, y_i) | i = 1, \dots, N\}$ ,
6  # Hyper - Hyper-parameters - refer to subsection 3.2.2.
7  # MIE - Mutual Information Estimator - refer to subsection 3.2.3.
8  # Output:
9  #  $I_x(e, l)$  - Information content about the input in layer  $l$  and epoch  $e$ ,
10 #  $I_y(e, l)$  - Information content about the label in layer  $l$  and epoch  $e$ .
11  $I_x, I_y = \{\}, \{\}$ 
12  $N = \text{len}(\text{Data})$  # Number of data points
13  $\text{NN} = \text{Instantiate\_Neural\_Network}(\text{Hyper})$ 
14  $L = \text{NN.layer\_count}()$ 
15 for  $e$  in range(0, N):
16      $\text{NN.run\_SGD\_once}(\text{Data})$ 
17     for  $l$  in range(0, L):
18         data = []
19         for  $x \in \text{Data.X}$ :
20              $\hat{l} = \text{NN.layer}(l).predict(x)$ 
21             data.append( $\hat{l}$ )
22          $I_x(e, l) = \text{MIE.estimateX}(\text{data}, \text{Data.X})$  #  $\text{Data.X} = \{x_1, \dots, x_N\}$ 
23          $I_y(e, l) = \text{MIE.estimateY}(\text{data}, \text{Data.Y})$  #  $\text{Data.Y} = \{y_1, \dots, y_N\}$ 
24 return  $I_x(e, l), I_y(e, l)$ 

```

Figure 3.1: The general algorithm for calculating mutual information inside a neural network.

3.2.2 Hyperparameters

Figure 3.1 mentions Hyperparameters – these are parameters that do not change over the course of the training period and are usually set by Person conducting the experiment.

Hyperparameters can have a large impact on our results. It is important to vary our hyperparameters – by doing so we make sure that the results we see are robust and not just an epiphenomenon of our specific hyperparameters.

Examples of hyperparameters:

- Data Set Used – for example if we use Tishby’s[1] or the MNIST[3] dataset.
- Size of the training set – how much of the dataset should we use for the training of the neural network.
- Activation Function of the Network – activation function between layers, examples given previously: Leaky ReLu, Sigmoid, Tanh.
- Batch size for SGD – At every iteration SGD only consider a subset of the dataset the size of this subset is determined by the Batch Size.

- Number of training Epochs – How long we train the neural network, can be considered a hyperparameter. However, in our case it is less important as we measure information for individual epochs.
- Network shape – Network shape refers to number of layers and dimension of every layer.
- We can also consider the MIE and their hyperparameters to be hyper parameter as it affects our results. As for example of MIEs hyper parameters:
 - Binning subsection 3.2.4 – defines *bins*: which assigns activations of individual nodes to one of these *bins*.
 - Kernel Density Estimator subsection 3.2.5 – defines *Noise variance*: which tunes the added Gaussian Noise
 - As If Random subsection 3.2.7 – defines *Batches*: which defines how many epochs will be aggregated together to perform Mutual Information Estimation. It also defines *bins* as in the Binning MIE.

3.2.3 Mutual Information Estimators

Recall subsection 2.3.1 – which talks about how we can measure information within a NN. We note that in order to measure MI we need random variables, however we only have the empirical dataset $D = \{(x_i, y_i) | i = 1, \dots, N\}$. We would like to use the random variables defined by the routine in Figure 2.3. However, in subsection 3.1.2 we prove that parameters θ have to be treated as random variables – a property which is missing from the current definition. Figure 3.2 updates the definition of our random variable to capture the missing property.

Figure 3.2 defines the correlated random variables: $X, T_{e,l}, Y$ – which specify the MI values: $I(X, T_{e,l})$ and $I(Y, T_{e,l})$. However, computing the MI values exactly is a computationally infeasible task – as we cannot capture the full randomness produced by `line 2`. As such we need a way to estimate the MI values. Fundamentally we need to make assumptions about the input distribution – which would allow us to efficiently compute the values. We will consider 3 different Mutual Information Estimators (**MIE**), all of which simulate randomness in a different way: Binning, KDE, and AIR.

```

1  def rxy( $\hat{\theta}$ , l):
2      Let  $\phi$  be one update step of SGD on  $\hat{\theta}$ 
3      pick  $i \sim \text{Uniform}\{1 \dots N\}$ 
4      return  $(x_i, F_{\phi}^l(x_i), y_i)$ 

```

Figure 3.2: Updated definition of correlated random variables $X, T_{e,l}$ and, Y . The Probability distributions are generated from the dataset $D = \{(x_i, y_i) | i = 1, \dots, N\}$. $F_{\theta(e)}^l$ is defined by Equation 2.18

A Note on The Input Dataset

We have previously defined the data set $D = \{(x_i, y_i) | i = 1, \dots, N\}$ – it is the sequence of pairs where x is the input and y is the label. This is the dataset available to us when we are training

the NN. However, often D is a subset of D_{global} , where D_{global} is all valid possible (x, y) pairs, which in most interesting cases is infinite. Take for example the MNIST dataset. The MNIST dataset is a collection of labeled handwritten digits – where every input x is an image and every label y is digit in the image. In this case D is a set containing about 70,000 (image, label) pairs and D_{global} is the set of all images containing a handwritten image and the labels associated with the image.

The task of a NN is to predict the labels given any $x \in D_{global}$, regardless if $x \in D$ holds or not. Hence, we are fundamentally interested in how much information is retained in the NN about the inputs from the dataset D_{global} . However, we can only make an attempt to measure information retained about inputs from the dataset D . Therefore, even if we were able to exactly compute MI values defined by Figure 3.2, they would still be only estimates.

Auxiliary functions

```

1  def getProbabilitiesOfData(A):
2      # A = [a_1, ..., a_N]
3      # Unique(A) - returns A with duplicate elements removed
4      for  $\alpha$  in Unique(A):
5          count = 0
6          for a in A:
7              if  $\alpha == a$ :
8                  count = count + 1
9           $P_\alpha = \text{count} / N$ 
10     return P

```

Figure 3.3: The routines defines a random variable \hat{A} for any dataset A s.t. the equation $\forall a \in A. P(\hat{A} = a) = \frac{\text{Number of occurrences of } a \text{ in } A}{\text{size of } A}$ holds.

3.2.4 MIE - Binning

The Assumptions

- assumes X is uniform
- simulates randomness by binning. What is binning?? takes in hyperparameter $bins$
- low and high are usually decided by the activation function Sigmoid 0 to 1, Tanh -1 to 1, as for Leaky ReLu we first need to aggregate we take it to be the minimum and maximum values for the observed layer activations
- Consider value \hat{t} in Figure 3.4, it can only obtain values in the range $1, \dots, bins$
- using binning to simulate randomness implies T can only take values in the range $1, \dots, bins$ – that is rather than being a continuous variable T is discrete.
- Which implies we can use the discrete entropy functions Equation 2.1 and Equation 2.2 to compute $H(T)$ $H(T|X)$. Which in turn lets use the MI function Equation 2.7.

- $j++j$
- By treating grouping values together we are simulating overlapping values generated by different θ values. This is a very crude method to simulate randomness.
- two values that are very close together but are separated by a boundary are treated as if they have nothing in common.

```

1 def rxy(e, l):
2     pick i ~ Uniform {1...N}
3     t = Fθ(e)l(xi)
4     t̂ = batchVector(t) # refer to Figure 3.5
5     return (xi, t̂, yi)

```

Figure 3.4

```

1 def batchVector(t):
2     # applies batchValue to every value of vector t
3     # Input: t - a value in ℝd
4     t̂ = [batchValue(a) for a in t]
5     return t̂
6
7 def batchValue(t):
8     # If we separate the range [low,high] into bins number of intervals
9     # batch returns the interval index that t would fall into.
10    # Input: t - a value in ℝ, s.t. low ≤ t ≤ high
11    # Globally defined values:
12    # bins - number of intervals to separate the range [low,high] into.
13    # low - lowest expected value of t
14    # high - highest expected value of t
15    s = (low - high)/bins
16    return ⌊(t - low)/s⌋

```

Figure 3.5

The Algorithm

```

1  def binning.estimateX(T, X):
2      # Estimates mutual information between the observed states T and input X
3      #  $\forall i \in \{0, \dots, N-1\}$ .  $T[i]$  is the observed state for  $X[i]$ .
4      # Since every  $x \in X$  is equally likely and unique  $H(T|X) = 0$ .
5      # hence,  $I(T, X) = H(T)$ 
6       $\hat{T} = \text{binning.batchData}(T)$ 
7      return binning.entropy( $\hat{T}$ )
8
9  def binning.estimateY(T, Y):
10     # Estimates mutual information between the observed states T and Y
11      $\hat{T} = \text{binning.batchData}(T)$ 
12      $H(\hat{T}) = \text{binning.entropy}(\hat{T})$ 
13      $H(\hat{T}|Y) = \text{binning.conditionalEntropy}(\hat{T}, Y)$ 
14     return  $H(\hat{T}) - H(\hat{T}|Y)$ 
15
16  def binning.entropy(A):
17     #  $A = [a_1, \dots, a_N]$ 
18     # Unique(A) - returns A with duplicate elements removed
19      $P = \text{getProbabilitiesOfData}(A)$  # refer to Figure 3.3
20      $H(A) = - \sum_{\alpha \in \text{Unique}(A)} P_{\alpha} \log_2(P_{\alpha})$ 
21     return H(A)
22
23  def binning.conditionalEntropy(A, B):
24     # Returns  $H(A|B)$ 
25     #  $A = [a_1, \dots, a_N]$ 
26     #  $B = [b_1, \dots, b_N]$ 
27      $P = \text{getProbabilitiesOfData}(B)$  # refer to Figure 3.3
28      $H(A|B) = 0$ 
29     for  $\beta$  in Unique(B)
30          $\hat{A} = []$ 
31         for a, b in zip(A, B) #  $\text{zip}(A, B) = [(a_1, b_1), \dots, (a_N, b_N)]$ 
32             if  $b == \beta$ 
33                  $\hat{A}.append(a)$ 
34                  $H(A|B) += P_{\beta} * \text{binning.entropy}(\hat{A})$ 
35     return  $H(A|B)$ 
36
37  def binning.batchData(T):
38      $\hat{T} = [\text{batchVector}(t) \text{ for } t \text{ in } T]$  # refer to Figure 3.5
39     return  $\hat{T}$ 

```

Figure 3.6: Implementation of binning MIE. Functions estimateX and estimateY are called by Algorithm from Figure 3.1.

binning.entropy Equation 2.1

binning.conditionalEntropy Equation 2.3 and Equation 2.4

3.2.5 MIE - Kernel Density Estimation

The Assumptions

Gaussian

- as well as the binning method – KDE as assumes that X is uniform.
- In this section let \log refer to the natural logarithm
- KDE assumes T to have Gaussian noise.

$$T_{e,l} = f_{\theta(e)}^l(X) + \epsilon, \text{ where } \epsilon \sim \mathcal{N}_{d_l}(0, \Sigma) \quad (3.20)$$

Σ here is the covariance matrix. In the implementation Σ is taken to be the identity matrix I multiplied by some σ . σ - is just a number in \mathbb{R}

- Consider $H(T|X)$ Under this assumption we get

$$P(T_{e,l} = f_{\theta(e)}^l(x) + t | X = x) = P(\mathcal{N}_{d_l}(0, \Sigma) = t) \quad (3.21)$$

- Equation 3.21 implies that

$$H(T_{e,l}|X) = H(\mathcal{N}_{d_l}(0, \Sigma)) \quad (3.22)$$

- The entropy of a Gaussian is analytically computable

$$H(\mathcal{N}_{d_l}(0, \Sigma)) = \frac{d_l}{2\log(2)} (\log(2\pi) + \log_e(\det(\Sigma)) + 1) \quad (3.23)$$

- Let us assume that $\Sigma = \sigma I$, where $\sigma \in \mathbb{R}$ is value representing noise and I is the identity matrix. This implies $\det(\Sigma) = \sigma$ and hence Equation 3.23 can be simplified to

$$H(\mathcal{N}_{d_l}(0, \Sigma)) = \frac{d_l}{2\log(2)} (\log(2\pi\sigma) + 1) \quad (3.24)$$

- Hence, from Equations 3.22 and 3.24 – Equation 3.25 holds

$$H(T_{e,l}|X) = \frac{d_l}{2\log(2)} (\log(2\pi\sigma) + 1) \quad (3.25)$$

3.2.6 Discrete vs Continuous Distributions

- KDE defines the random variable $T_{e,l}$ as a continuous distribution, while X and Y are still defined to be discrete.
- This means that we need to compute MI between between a continuous and a discrete distribution: $I(X, T_{e,l})$ and $I(Y, T_{e,l})$.
- Let V be a discrete random variable that takes values in $\{v_1, \dots, v_M\}$ and let C be a continuous distribution. Consider $I(V, C)$

$$\begin{aligned} I(V, C) &= H(C) - H(C|V) \\ &= H(C) - \sum_{i=1}^M P(v_i) H(C|V = v_i), \text{ Equation 2.3} \end{aligned} \quad (3.26)$$

- Equation 3.26 shows that in order to compute MI between a discrete and continuous random variables, we can instead compute a sum of entropies for different continuous random variables.
- For our specific MI values $I(X, T_{e,l})$ and $I(Y, T_{e,l})$ – we need to compute:

$$\begin{aligned} I(X, T_{e,l}) &= H(T_{e,l}) - H(T_{e,l}|X) \\ &= H(T_{e,l}) - \sum_{i=1}^N P(X = x_i) H(T_{e,l}|X = x_i) \\ &= H(T_{e,l}) - \frac{1}{N} \sum_{i=1}^N H(T_{e,l}|X = x_i) \end{aligned} \quad (3.27)$$

and

$$\begin{aligned} I(Y, T_{e,l}) &= H(T_{e,l}) - H(T_{e,l}|Y) \\ &= H(T_{e,l}) - \sum_{i=1}^C P(Y = y_i) H(T_{e,l}|Y = y_i) \end{aligned} \quad (3.28)$$

we cannot reduce Equation 3.28 further as the labels in the dataset might be imbalanced.

- In how to compute the upper and lower bounds for these Continuous entropies see Kolchinsky and Tracey[5] Section 4, and Kolchinsky and Tracey[6] Eq. 10. Or for the implementation see[2]

The Algorithm

```

1  def KDE.estimateX(T, X):
2      # Estimates mutual information between the observed states T and input X
3      #  $\forall i \in \{0, \dots, N-1\}$ .  $T[i]$  is the observed state for  $X[i]$ .
4      # Globally defined :  $\sigma$  - is the assumed noise of the observations T
5       $d = T[0].dimension$  # an observation  $t$  of  $T$  is a value in  $\mathbb{R}^d$ 
6      # KDE assumes  $H(T|X)$  is a Gaussian distribution
7      # we can compute entropy of a Gaussian analytically as follows
8       $H(T|X) = \frac{d}{2}(\log_e(2\pi\sigma) + 1)/\log_e(2)$ 
9       $H(T) = \text{KDE.entropy}(T)$ 
10     return  $H(T) - H(T|X)$ 
11
12  def KDE.estimateY(T, Y):
13      # Estimates mutual information between the observed states T and Y
14       $H(T) = \text{KDE.entropy}(T)$ 
15       $H(T|Y) = \text{KDE.conditionalEntropy}(T, Y)$ 
16     return  $H(T) - H(T|Y)$ 
17
18  def KDE.entropy(A):
19      # See Kolchinsky and Tracey[5] Section 4, and Kolchinsky
20      # and Tracey[6] Eq. 10
21      # Code Implementation[2]
22
23  def KDE.conditionalEntropy(A, B)
24      # Returns  $H(A|B)$ 
25      #  $A = [a_1, \dots, a_N]$ 
26      #  $B = [b_1, \dots, b_N]$ 
27       $P = \text{getProbabilitiesOfData}(B)$  # refer to Figure 3.3
28       $H(A|B) = 0$ 
29      for  $\beta$  in Unique(B)
30           $\hat{A} = []$ 
31          for a, b in zip(A, B) #  $\text{zip}(A, B) = [(a_1, b_1), \dots, (a_N, b_N)]$ 
32              if  $b == \beta$ 
33                   $\hat{A}.append(a)$ 
34                   $H(A|B) += P_\beta * \text{KDE.entropy}(\hat{A})$ 
35     return  $H(A|B)$ 

```

Figure 3.7: Implementation of KDE MIE. Functions estimateX and estimateY are called by Algorithm from Figure 3.1.

3.2.7 MIE - As If Random

The Assumptions

- X is uniform
- want a better way estimate randomness
- In the later phases of training the NNs weights changes very little.
- The hypothesis is that aggregating multiple layers together can give a better estimation for MI as for a specific $P(T|X = x) = 1/batch$
- only trusted in the later epochs as in the early period of training there are allot of changes

The Algorithm

3.2.8 AIR

- I, novel
- Tishby's method is not ideal
- consider the perfect method teh new rxtx
- hopefully a better way to measure mutual information, more computationally expensive. Only applicable for later epochs.
- AIR
- Tishby and Saxe didn't capture randomness incredibly well
- We have a bettter way
- How we've done it, only measure compression at the end of the training period when the network is stable and only "brownian randomness happens"(refer paper what randomness).
- the problem is every x in X is unique, if we sample multiple epochs instead of only one we have more compresssion
- we rely on the property of the network that there is very little change at the end of the training period
- Achieves compression even with ReLu yay!!!!!!!!!!

We believe that Tishby's experiments Saxe's experiments could be improved, with introducing a better notion of randomness. blah blah blah not finished.

3.2.9 Binning method

The method used by Tishby in his paper. The method estimates mutual information between the X or Y and the hidden layer T by assuming the observed empirical distribution of input samples is the true distribution. We use this assumption to compute entropies of T and varying subsets as outlined in Equation 3.29

$$I(X, Y) = H(X) - H(X|Y) = H(X) - \sum_{y \in Y} H(X|Y = y)p(Y = y) \quad (3.29)$$

However just calculating mutual information for a discrete distribution is not enough as this yields mutual information equal to 0, in order to sidestep this issue we need to simulate randomness Tishby achieves this by grouping multiple values together, which he calls binning. A more detailed explanation why this is done is given in subsection 4.0.1

Figure 3.8 through to Figure 3.10 outline the full algorithm.

1 **Algorithm: Mutual Information**

Figure 3.8: Pseudo code for computing mutual information refer to Figure 3.9 and Figure 3.10 for entropy computation.

1 **Algorithm: Entropy - Binning Method**

Figure 3.9: Algorithm for computing entropy – Binning method

1 **Algorithm: Conditional Entropy**

Figure 3.10: Algorithm for computing conditional entropy

When computing mutual information between a hidden layer T and the input set X we can abuse the fact that every element $x \in X$ is unique and uniquely identifies an element $t \in T$ hence

$$H(T|X) = 0 \quad (3.30)$$

$$I(T, X) = H(T) - H(T|X) = H(T) \quad (3.31)$$

This does not affect the result but increases performance of our algorithm.

3.2.10 Kernel Density Estimation

The KDE method used in Saxe’s paper but originally devised by Kolchinsky & Tracey (2017); Kolchinsky et al. (2017). As well as the Binning method KDE assumes that the observed empirical distribution in the hidden layer T is the true distribution.

However, instead of binning values together KDE assumes the distribution is a mixture of Gaussians. Using this fact we can get an upper bound when calculating entropy for a collection of data points.

The algorithm closely follows figures :Figure 3.8, Figure 3.9, and Figure 3.10, however there is a change the way entropy is calculated so instead of Figure 3.9 we have Figure 3.11 below.

1 **Algorithm: Entropy – KDE**

Figure 3.11: Algorithm for computing entropy – Kernel Density Estimation method. The same algorithm as used by Saxe.

As before when computing mutual information between a hidden layer T and the input set X we can use the fact that every element $x \in X$ is unique and hence uniquely identifies an element in $t \in T$.

3.2.11 Advanced methods

Since mutual information estimation is a contentious part of the project we wanted to experiment with more advanced techniques however we were not able to adopt the methods for this project, the methods that we have tried are outlined below.

Mutual Information by Gaussian approximation

”Estimating Mutual Information by Local Gaussian Approximation” (Shuyang Gao). A promising way to estimate mutual information. However the implementation proved to be too difficult and time consuming so we abandoned it. I contacted out to the author but he was not able to provide any concrete code.

Geometrical Adaptive Entropy Estimation

”Breaking the Bandwidth Barrier: Geometric Adaptive Entropy Estimation” (Weihao Gao). We were pointed to this paper by S. Gao the author of the previous method, as previously this looked like a promising method to measure entropy and mutual information. Furthermore, the code was available online unfortunately we were not able to adapt the code for multidimensional values as a result we were getting wrong and inconsistent results. We decided making this method work would require too much time and is out of scope of this project.

3.3 Implementation Optimizations

Producing data for the information plane requires a substantial amount of computational power and memory, in order for the computation to complete in a reasonable amount of time we have to utilize all the available resources and minimize the amount of work we are doing.

3.3.1 Maximising Resource Utilization

The obvious way to maximise the resource usage is to parallelize the workload. If we refer to ??, we can see two main ways we can do this.

The first way is to parallelize one of the two outer loops and run mutual information calculations in parallel, this is easy to do, however an issue occurs if we need a lot of memory since every instance of mutual information calculation manipulates the datasets creating copies, this is an issue for bigger datasets such as MNIST.

The second way is to parallelize the mutual information calculation itself, this is nice since we use minimal amounts of memory, however might be hard or impossible to implement as it depends on the method.

The second way is to parallelize the mutual information calculation itself, this method has a bonus that uses minimal amounts of memory improving performance for bigger datasets such as MNIST. However implementing this option might be hard as it heavily depends on the maths of the method, or impossible if the method has an inherently linear part to it. In our case KDE parallel performance is very good, however Tishby's Binning method has some bottlenecks that I wasn't able to remove.

3.3.2 Minimising the Workload

Even when using all the systems resources the calculations take a very long time to complete as such we need to find a way to speed it up. If we consider an information plane graph for ex. Figure 2.4b, we see that from epoch to epoch there is very little change that is occurring, skipping some epochs might be a good way to reduce workload while keeping the overall result unchanged. We have implemented a couple of ways to skip the epochs

Simple Skip

A very simple way to skip epochs calculates mutual information for every n^{th} epoch.

It's quite effective and is fast to implement and easy to parallelize, however it yields subpar results. At the beginning of the training period during the fitting phase the step sizes are too big and yields gaps in data as there are big changes between consecutive epochs. Toward the end of the training period during the compression phase the step size becomes too small and we are wasting computation as the changes between epochs is negligible. The next two methods address this problem, but have their own drawbacks.

Delta Skip – Exact

The Exact Delta Skip method introduces a distance metric which measures mutual information distance between two epochs. Using the distance metric the method tries to skip as many epochs

as possible while still guaranteeing that the distance is at most δ , and backtracks when necessary.

The Algorithm starts by measuring every consecutive epoch ($skip = 1$) as at the start of the training epochs are far apart – distance between them is more than δ . When the distance become smaller there is no need to measure every epoch so we exponentially increase $skip$ until the difference between consecutively measured epochs is larger than δ . At this point we run into the issue of backtracking since we made a guarantee that every measured epoch will be at most δ apart.

We can consider the backtracing problem as an array of unknown – let's call this array a section. Every unknown in a section corresponds to the saved state of the neural network, revealing the unknown is equivalent to computing mutual information for the epoch. In this context we can rephrase the problem of backtracking as finding a subset of the section such that the difference between every consecutive number is less than or equal δ .

Consider an example with $\delta = 2$. At the start all the values are unknown

?	?	?	?	?	?	?	?	?	?
---	---	---	---	---	---	---	---	---	---

We measure values at the start and the end of the section to get the range.

2	?	?	?	?	?	?	?	14
----------	---	---	---	---	---	---	---	-----------

We continue to split the section and measure the middle element until the difference between consecutive elements is less than or equal to δ or there are no more elements left.

2	?	?	?	10	?	?	?	14
---	---	---	---	-----------	---	---	---	----

2	?	6	?	10	?	12	?	14
---	---	----------	---	----	---	-----------	---	----

2	5	6	8	10	?	12	?	14
---	----------	---	----------	----	---	----	---	----

The algorithm stops at this stage as:

- The distances between 10, 12 and 14 are 2 and $\delta \leq 2$ holds.
- There is no cell between 2 and 5, even though the distance is $3 > \delta$.

Every unknown has to save the state of the neural network, that means that every unknown contains $|X| * |nodes|$ float numbers. To put it into perspective if we use the MNIST dataset which has 240,000 samples into a small network of 128 nodes every unknown would contain 0.25GB of data. As such this method is unsuited for large datasets in that case Delta Skip Approximate is more suited for the job. A way to remedy the problem is instead of saving all the activations, just save the neural network weights and compute the activations just before calculating Mutual Information for the epoch

Parallelizing The entire single threaded algorithm is described in Figure 3.12 and Figure 3.13. However there are two main ways how we can parallelize the algorithm.

- Parallelize Backtracking – every time we need to backtrack we can launch two threads to do the computation, it's quite easy to implement however we need to be careful as we don't want to change the global *skip* value.
- Compute Multiple sections – consider a section as before. Once we compute the latest epoch of the array we can launch backtracking and the next section computation in parallel as separate threads. We need to wait until the latest epoch is computed before computing the next section as we might need to update the *skip* value.

```

1  Algorithm: DeltaSkipExact
2  Input:
3  prev = mutual information result of the previous epoch
4  curr = mutual information result of the current epoch
5   $\delta$  = user specified maximum "distance" between epochs
6  multiplier = how much to multiply skip by
7  Output:
8  Algorithm:
9  dist = Distance(prev, curr)
10 if dist >  $\delta$ :
11     Backtrack(prev, curr) # Figure 3.13
12 else:
13     skip = skip*multiplier

```

Figure 3.12: Delta Skip Exact. The skip value is assumed to be global it specifies how many epochs to skip until we measure again

```

1  Algorithm: Backtrack
2  Input:
3  prev = mutual information result of the previous epoch
4  curr = mutual information result of the current epoch
5   $\delta$  = user specified maximum "distance" between epochs
6  Output:
7  if prev.epoch + 1 < curr.epoch:
8     mid_epoch = average(prev.epoch, curr.epoch)
9     mid = Calculate Mutual Information of mid_epoch
10 DeltaSkipExact(prev, mid,  $\delta$ , 1)
11 DeltaSkipExact(mid, curr,  $\delta$ , 1)

```

Figure 3.13: Backtrack Algorithm

Distance Metric Every epoch that we measure yields us with a vector of mutual information values, that is for every layer T we receive two values $I(X, T)$ and $I(Y, T)$. Given the information vectors for two epochs we need to find a reasonable way to measure distance between them.

The distance is used for the purposes of speeding up the computation and won't meaningfully impact the results. I've chosen to define the distance as the maximum shift between the two epochs refer to Equation 3.32 how to compute it or to Figure 3.14 for a graphical representation.

$$D = \max_{t \in T} [\max(I_e(X, t) - I_{\hat{e}}(X, t), I_e(Y, t) - I_{\hat{e}}(Y, t))] \quad (3.32)$$

The equation reduces the vector to a single value that is easy to compare and to track.

If we consider the Figure 3.14 we can see that the axis are different in scale this is due to input X and label Y having different entropy values $H(X) = 12$ and $H(Y) = 1$. As such we might wish to adjust Equation 3.32, and scale mutual information values by the entropy as in Equation 3.33.

$$D = \max_{s \in \{X, Y\}} \max_{t \in T} [(I_e(s, t) - I_{\hat{e}}(s, t)) / H(s)] \quad (3.33)$$

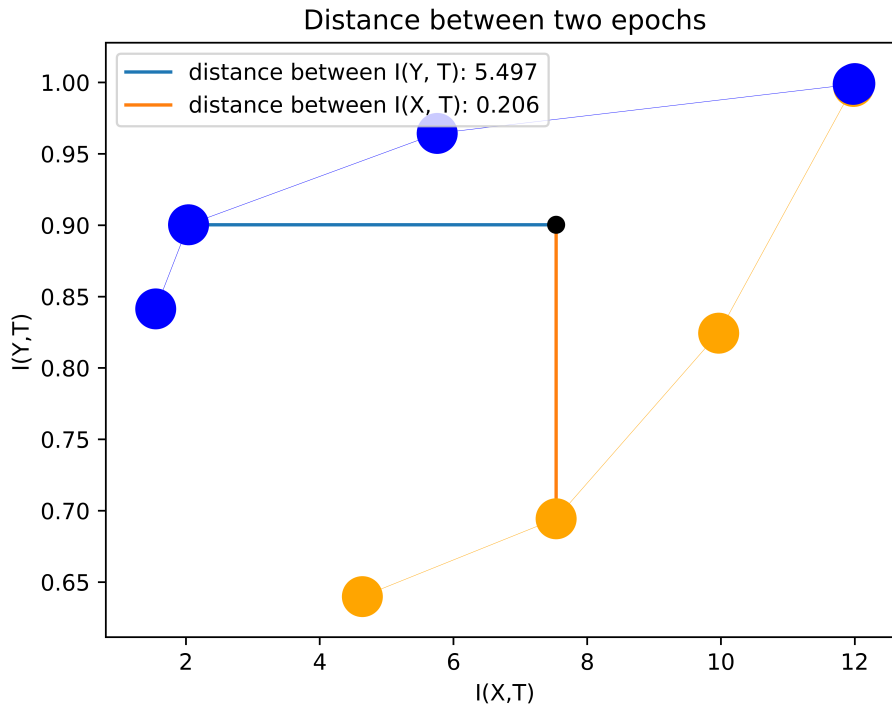


Figure 3.14: Example how distance between two epochs is measured.

Delta Skip – Approximate

The Exact method suffers the same problem as when we try to instantiate too many mutual information calculation instances, namely it runs out of memory if our dataset is too big. In order to solve this we use the approximate method which follows closely the Exact method.

The Algorithm – refer to Figure 3.15, as previously, in the Exact method it uses a distance metric and measures every n 'th epoch where n is adaptive and depends on how close the

epochs are. The critical difference between Exact and Approximate method happens when the distance between epochs is more than δ , the Exact method attempts to backtrack and fill in the gap whereas the Approximate doesn't backtrack and just continues to the next epoch, this is justified because the approximate method assumes that distance between epochs only shrinks and never increases.

This method does not suffer from memory issues but cannot be parallelized as in order to compute next epoch we need to know the distance between current epoch and the previous one. Delta Approximate method performs best when paired with highly parallel Mutual Information Estimator.

```
1  Algorithm: Delta Skip - Approximate
2  Input:
3  prev = mutual information result of the previous epoch
4  curr = mutual information result of the current epoch
5   $\delta$  = user specified estimated "distance" between epochs
6  Output:
7  Algorithm:
8  dist = Distance(prev, curr)
9  if dist >  $\delta$ :
10     skip = skip
11 else:
12     skip = skip*2
```

Figure 3.15: Delta Skip Approximate

3.4 Repository Structure

Chapter 4

Evaluation

Tishby's Experiment Tishby[?] talks about a subset of the hyper parameters mentioned above.

Tishby varies the following

- Even the implementation of the mathematical ideas can be considered a hyper parameter – as it is unlikely that implementations follow the mathematical model exactly. This is the reason why independent verification is important and partly the reason of this thesis.
- Notes on comparing MIEs need to run on the same NN as they vary, need to run in real time as saving the whole network
- $i++j$
- $i++j$
- Success Criteria
- we have successfully reproduced the results showed by Tishby and Saxe. However there are reasons to not trust either of them as they have flaws with them.
 - Tishby – used only a toy dataset
 - Saxe – Changed allot of parameters at once made the claim that no compression phase happens
- We need better tools for MIE
 - cannot judge subtleties if something has a compression phase our MIE are not trusted
 - we have seen KDE and Discrete show inconsistent results, when n 'th layers has less information about the input than the $n+1$ 'th layer.
- compression phase
 - judging if the networks have a compression phase is moot point as of yet as tools for measuring information are not good enough. Case and point Saxe argues that there cannot be compression in NN but every every experiments show that compression exists.

- Saxe states that there is no compression in NN, however their experiments disagree.
 - we don't have the tools to say if the compression phase is actually happening
 - Tishby says compression phase happens but he is using a toy dataset, which was shown to not have a compression phase by Saxe.
- performance
-

4.0.1 Deterministic networks

There's a very real argument to be made against compression in neural networks. Consider a generic neural network we can think of it as a function that is a series of matrix transformation, where a matrix corresponds to weights of a specific layer. However these matrices are all random (at least at the start of training) and hence probability of them being invertible is 100%.

Knowing that every single matrix is invertible allows us to conclude that that neural network as a whole is an invertible function, which means no information is lost and compression is impossible.

4.0.2 Why Randomness is hard to capture

Chapter 5

Conclusion

Bibliography

- [1] Dataset used by tishby. <https://github.com/ravidziv/IDNNs/tree/master/data>. Accessed : 14/05/2019, Refer to Tishby and Schwarz-Ziv, 2017, section 3.1.
- [2] Implementation of kde entropy function. <https://github.com/artemyk/ibsgd/blob/iclr2018/kde.py#L1>. Accessed : 14/05/2019,.
- [3] Mnist dataset. <http://yann.lecun.com/exdb/mnist/>. Accessed : 14/05/2019.
- [4] Joel Dapello Madhu Advani Artemy Kolchinsky Brendan D. Tracey David D. Cox Andrew M. Šaxe, Yamini Bansal. On the information bottleneck theory of deep learning, 2018.
- [5] Brendan D. Tracey Artemy Kolchinsky.
- [6] David H. Wolpert Artemy Kolchinsky, Brendan D. Tracey.
- [7] Noga Zaslavsky Naftali Tishby. Deep learning and the information bottleneck principle, 2015.
- [8] Naftali Tishby Ravid Schwarz-Ziv. Opening the black box of deep neural networks via information, 2017.

Appendix A

Project Proposal

Measuring mutual information within Neural networks

2359A

REDACTED College

Wednesday 15th May, 2019

Project Originator: 2359A

Project Supervisor: Dr. Damon Wischik

Director of Studies: Prof. Alan Mycroft

Overseers: Dr. Robert Mullins Prof. Pietro Lio'

Introduction and Description of the Work

The goal of this project is to confirm or deny the results produced by Shwartz-ziv & Tishby in their paper "Opening the black box of Deep Neural Networks via Information"¹

The paper tackles our understating of Deep Neural Networks (DNN's). As of yet there is no comprehensive theoretical understanding of how DNN's learn from data. The authors proposed to measure how information travels within the DNN's layers.

They found that training of neural networks can be split into to two distinct phases: memorization followed by the compression phase.

- memorization - each layer increases information about the input and the label
- compression - this is the generalization stage where each layer tries to forget details about the input while still increasing mutual information with the label thus improving performance of the DNN. This phase takes the wast majority of the training time.

They found that each layer in neural network tries to throw out unnecessary data from the input while preserving information about the output/label. As the network is trained each layer preserves more information about the label

¹<https://arxiv.org/abs/1703.00810>

The results they found were interesting but also contentious as they have not yet provided a formal proof, just experimental data as a result there are many peers that are cautious and sceptical of the theory even a paper² was produced that tries to suggest that the theory is wrong, however this was dismissed by Tishby & Shwartz-Ziv³

Starting Point

I have watched a talk that Prof. Tishby gave on this topic at Yandex, no other preparation was done.

Resources Required

The training DNN's and measuring mutual information will be computationally expensive so I will be using Azure cloud GPU service to acquire the required compute for this project. The GPU credits will be provided by Damon Wischik

For backups I intend to store my work on GitHub and my own personal machine. In case my laptop breaks I will get another one or use the MCS machines.

Substance and Structure of the Project

The aim of this project to reproduce the results provided by Prof. Tishby and his colleagues. The intention of my work is to help settle the debate surrounding the topic either strengthening the arguments in favour of the theory in case my results are inline with the aforementioned results or encourage discussion in case my results contradict the theory.

My work will require me to have a comprehensive understanding of Information theory, Information bottleneck and neural networks.

One of the more contentious parts of my project will be measuring mutual information between the input a layer in the DNN and the label. It will be computationally expensive to measure it in DNN since we will need to retrain the network in order to get a distribution rather than a single value. I will use Gaussian approximation to measure it (relevant paper⁴)

Will need to use Python to train the neural networks and GNUplot or alternative to plot the results.

Success Criteria

Reimplement the code that was used to generate the papers results. Confirm or deny the results produced in "Opening the black box of Deep Neural Networks via Information" paper on the same dataset as the paper. In order to do that I will need to: Train a neural network on the same dataset

²https://openreview.net/pdf?id=ry_WPG-A-

³https://openreview.net/forum?id=ry_WPG-A-¬eId=S1lBxcE1z

⁴<https://arxiv.org/abs/1508.00536>

that was used in the paper and measure mutual information between the layers. Analyse the results produced and address any discrepancies that may have occurred.

Extensions

Provided I achieve the success criteria there are two main ways to extend it.

- Use different datasets to test the theory. Using different datasets would confirm that the results are not data specific. Current datasets we are considering: MNIST⁵ and NOT-MNIST⁶.
- Explore different ways of measuring mutual information. One interesting way would be to explore a discrete neural network where every node would only be able assigned discrete values say 1...256. This would make the distribution within a DNN layer discrete and hence it would make calculating mutual information straightforward. However quantizing the neural network could possibly hurt the performance of the network.

⁵<http://yann.lecun.com/exdb/mnist/>

⁶<https://www.kaggle.com/quanbk/notmnist>

Schedule

- **20th Oct – 2nd Nov**

I expect to spend the first two weeks reading up on Information theory (primarily from Mackay's book⁷) and the information bottleneck method in order to understand the nuances of the paper.

- **3rd Nov – 30th Nov**

The following weeks I intend to spend reading up on DNN's doing some introductory courses, I will train the neural network on the same data as the paper but at this point will not yet try to measure the mutual information between the layers.

At this point I will also start examining the code⁸ provided and start to implement parts of it which don't deal with information measurement.

- **1st Dec – 28th Dec**

Will start reading up on mutual Information measurement with local Gaussian approximation.

Implementing mutual information measurement in code.

At this point I expect the computation to be too demanding for my machine and will need to use provided compute.

- **29th Dec – 1st Feb**

Having a working system to test data sets I will try to reproduce results from the paper on the same dataset. This will achieve my success criteria.

At this point my success criteria should be completed I will spend some time writing the skeleton of the thesis. Look for any discrepancies between my results and the ones provided in the paper.

- **2nd Feb – 15th Mar**

Assuming everything goes as planned I will start looking into implementing one of the extensions. Which are :

- Testing the theory on different datasets.
- Implementing a quantized neural network implementation.

or both, if time is in my favour.

- **16th Mar – 12th Apr**

Will use the remaining time to write up the dissertation.

⁷Information Theory, Inference, and Learning Algorithms by David J. C. MacKay

⁸<https://github.com/ravidziv/IDNNs>