

Andrius Grabauskas

Measuring mutual information in Neural Networks

Computer Science Tripos – Part II

Robinson College

Thursday 9th 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

Contents

| | | |
|----------|---|-----------|
| 1 | Introduction | 11 |
| 1.1 | The Information plane | 12 |
| 2 | Preparation | 15 |
| 2.1 | Entropy and Mutual Information | 15 |
| 2.2 | Neural Network | 16 |
| 2.3 | The Information plane | 17 |
| 2.3.1 | Setup | 17 |
| 2.3.2 | Visualization | 17 |
| 2.3.3 | Interpretation of the Information Plane | 18 |
| 2.3.4 | Compression In Neural Networks | 19 |
| 3 | Implementation | 25 |
| 3.1 | Mutual Information Estimation | 25 |
| 3.1.1 | Mutual Information Definition | 25 |
| 3.1.2 | Theoretically undefined | 25 |
| 3.2 | Calculating Mutual Information | 26 |
| 3.2.1 | Experimental setup | 26 |
| 3.2.2 | Discrete method | 27 |
| 3.2.3 | Kernel Density Estimation | 28 |
| 3.2.4 | Advanced methods | 29 |
| 3.3 | Implementation Optimizations | 30 |
| 3.3.1 | Maximising Resource Utilization | 30 |
| 3.3.2 | Minimising the Workload | 30 |
| 3.4 | As-If-Random Experiment | 34 |
| 3.5 | Repository Structure | 35 |
| 4 | Evaluation | 37 |
| 4.0.1 | Deterministic networks | 37 |
| 4.0.2 | Why Randomness is hard to capture | 37 |
| 5 | Conclusion | 39 |
| | Bibliography | 39 |
| A | Project Proposal | 41 |

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 | activation of layer l for input x | 16 |
| 2.2 | Definition of correlated random variables X, T_e^t and, Y | 17 |
| 2.3 | information plane for a neural network with 5 layers, which was only trained for one epoch. | 18 |
| 2.4 | Information plane for a neural network with 4 layers, which was trained for approximately 10 000 epochs. The Neural Network was trained on the same dataset as used by Tishby, 2017, hence entropy of the input $H(X)$ is 12 and entropy of the label $H(Y)$ is 1. | 18 |
| 2.5 | random distribution of $\theta'(e)$, where e is an epoch | 21 |
| 3.1 | The general algorithm for calculating mutual information inside a neural network. | 26 |
| 3.2 | Pseudo code for computing mutual information refer to Figure 3.3 and Figure 3.4 for entropy computation. | 27 |
| 3.3 | Algorithm for computing entropy – Discrete method | 27 |
| 3.4 | Algorithm for computing conditional entropy | 28 |
| 3.5 | Algorithm for computing entropy – Kernel Density Estimation method. The same algorithm as used by Saxe. | 29 |
| 3.6 | Delta Skip Exact. The skip value is assumed to be global it specifies how many epochs to skip until we measure again | 32 |
| 3.7 | Backtrack Algorithm | 32 |
| 3.8 | Example how distance between two epochs is measured. | 33 |
| 3.9 | Delta Skip Approximate | 34 |

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`

Acknowledgements

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 summarized this Markov chain view and applies it to the entirety of the training process.

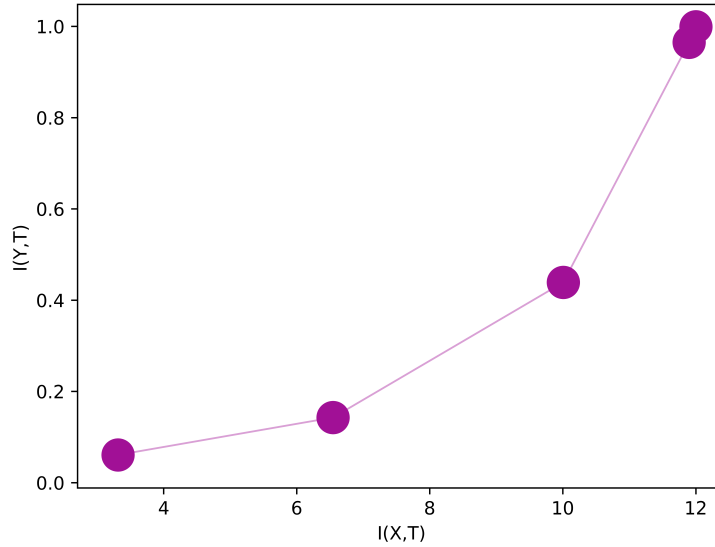


Figure 1.1: information plane for a neural network with 5 layers, which was only trained for one epoch.

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.

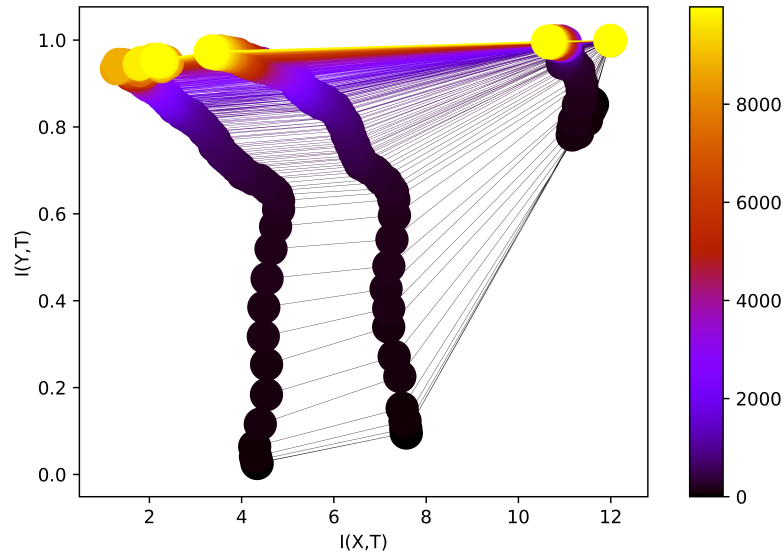


Figure 1.2: Information plane for a neural network with 4 layers, which was trained for approximately 10 000 epochs

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

2.1 Entropy and Mutual Information

Entropy – quantifies information content of a random variable. It's generally measured in bits and can be thought of as the expected information content when we sample a random variable once. Entropy for a discrete random variable is defined in Equation 2.1.

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

Consider a random variable X s.t $P(X = 1) = P(X = 0) = 0.5$, this tells us that $H(X) = 1$ and every time we sample X we are expected to gain 1 bit of information.

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

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. Entropy of Y given X is written as $H(Y|X)$. Conditional Entropy for a discrete random variable is defined in Equation 2.2

$$H(Y|X) = - \sum_{x \in X, y \in Y} p(x, y) \log \frac{p(x, y)}{p(x)} \quad (2.2)$$

Consider a random variable pair X and Y s.t

$$P(X = 0, Y = 0) = P(X = 0, Y = 1) = 0.25, P(X = 1, Y = 0) = 0.5, P(X = 1, Y = 1) = 0$$

then $H(Y|X) = 0.5$ and $H(X|Y) \approx 0.6887$, while $H(X) = 1$ and $H(Y) \approx 0.8112$

Mutual Information (MI) – is a measure between two random variables, it builds on top of entropy and measures how much information the variables have in common. It quantifies information gained about one variables when observing the other. Mutual information is generally computed either explicitly using the probabilities as in Equation 2.3 or using entropies as in Equation 2.4.

$$I(X, Y) = \sum_{y \in Y} \sum_{x \in X} p(x, y) \log \left(\frac{p(x, y)}{p(x)p(y)} \right) \quad (2.3)$$

$$I(X, Y) = H(X) - H(X|Y) \quad (2.4)$$

Consider the random variables X and Y as before in the conditional entropy section.

$$P(X = 0, Y = 0) = P(X = 0, Y = 1) = 0.25, P(X = 1, Y = 0) = 0.5$$

Here we can see that $I(X, Y) \approx 0.3113$

2.2 Neural Network

Neural Networks are a machine learning framework that is intended to solve the problem of prediction or classification of data.

Suppose we have a dataset:

$$(x_i, y_i) \text{ for } i = 1, \dots, N, x_i \in \mathbb{R}^d, y_i \in \mathbb{R}^{d'}$$

where x_i is the input to the Neural Network and y_i is the label or the value we want the neural network to predict. Suppose we are trying to classify if an image contains a cat or no, then x_i would be an encoding of an image and the label y_i could be 1 if the image contains a cat and 0 if it does not.

The network is structured in layers where every layer holds an intermediate representation of the final prediction output, lets call this intermediate representation an activation of that specific layer. For our purposes we will define a neural network to be a sequence of functions $f_\theta^1, f_\theta^2, \dots, f_\theta^N$ that are parameterized by weights θ s.t

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

We have one function for each layer transition of the neural network this allows us to extract the activation of a specific layer as in Figure 2.1, we will need to get access to intermediate layers in our experiments.

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

Figure 2.1: activation of layer l for input x

Stochastic Gradient Descent Each of the transition function are parameterized by the weights function θ . The θ function is directly controlled by Stochastic Gradient Descent (SGD) algorithm. SGD is an algorithm for training a neural network it gradually updates the weights θ according to some goal such as minimising the prediction error. SGD is an iterative process, it has a notion of epochs where one iteration of the algorithm advances epochs by one, this implies that θ function depends on which epoch we are currently at thus it must take the epoch number as an argument – $\theta(e)$.

2.3 The Information plane

2.3.1 Setup

The Information plane is a way of visualizing the Neural Networks training process through the information domain. By looking at mutual information between the input X the intermediate neural network layer activations T and the label Y we can see how the information flows through the network.

Mutual Information, as we defined it, is only applicable to Probability distributions, however we currently only have the dataset (x_i, y_i) for $i = 1 \dots N$, if we assume that every input x_i equally likely we can provide a routine that defines our probability distribution.

For convenience let's define F_θ^t to be the activation of layer t given input x , i.e

$$F_\theta^t(x) = f_\theta^t(f_\theta^{t-1}(\dots(f_\theta^1(x)))) \quad (2.5)$$

Consider Figure 2.2, the routine defines the random variables X, T_e^t, Y . Here T_e^t is the distribution of layer t for the epoch e , X and Y is the original data with assumption that it's uniformly distributed.

Using the probability distributions we now have values:

- $I(T_e^t, X)$ – Mutual Information between the input distribution and the layer activations
- $I(T_e^t, Y)$ – Mutual Information between the label distribution and the layer activations.

This allows us to generate the Information plane.

```

1 def rxy(e, t):
2     pick i ~ Uniform {1...N}
3     return (x_i, F_{\theta(e)}^t(x_i), y_i)
```

Figure 2.2: Definition of correlated random variables X, T_e^t and, Y

2.3.2 Visualization

The Information plane visualizes the whole training process, in order to generate the information plane we need $I(X, T_e^t)$ and $I(Y, T_e^t)$ for every epoch e and every layer t .

Consider for now that we only trained our neural network for one epoch, the Figure 2.3 shows an example of this.

In this case the network consists of 5 layers, in the figure every node corresponds to a distinct layer. The lines between the nodes help us distinguish epochs from each other and gives helps us to see the order of the layers, the upper-right-most node corresponds to the first layer in the neural network, the lower-left-most node corresponds to the fifth and last layer of the network.

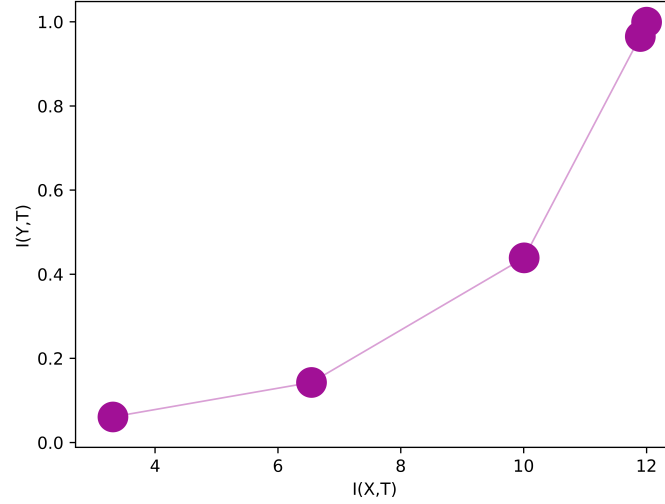


Figure 2.3: information plane for a neural network with 5 layers, which was only trained for one epoch.

Consider now Figure 2.4, it shows an information plane for a full training phase. The color signifies what epoch the data belongs to and let's us see how the network progressed over time. We can see that at the start $I(Y, T_1^4) \approx 0$ meaning the network has not preserved any information about the label distribution Y , but by the end of the training we see $I(Y, T_{10,000}^4) \approx 1$ which means we have preserved almost all the information about the label.

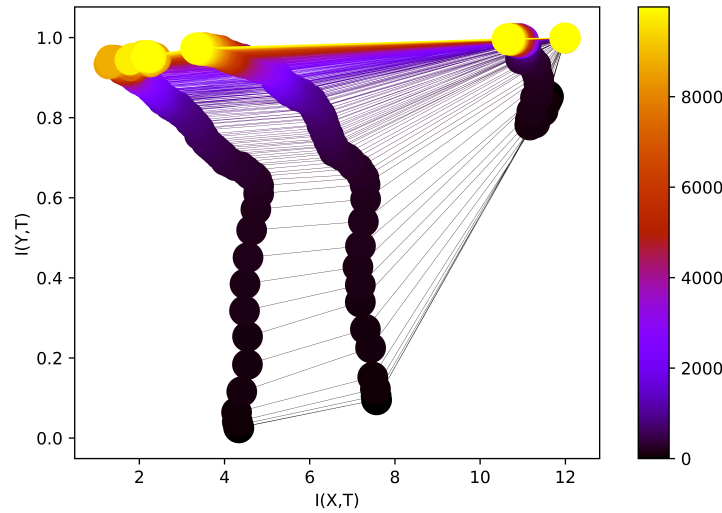


Figure 2.4: Information plane for a neural network with 4 layers, which was trained for approximately 10 000 epochs. The Neural Network was trained on the same dataset as used by Tishby, 2017, hence entropy of the input $H(X)$ is 12 and entropy of the label $H(Y)$ is 1.

2.3.3 Interpretation of the Information Plane

Let's once again consider Figure 2.4, we can see two phases in the figure Tishby has named them The Fitting Phase and The Compression Phase.

The Fitting Phase In Figure 2.4 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 $I(Y, T_e^t)$, 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 $I(X, T_e^t)$, the information about the input, as we advance through the epochs e , in our case we see very little change in $I(X, T_e^t)$.

During the fitting phase a neural network tries to memorize the data and make predictions based on the observations, this means that the network may learn useless features that only superficially correlate with the correct label.

The Compression Phase In Figure 2.4 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 $I(Y, T_e^t)$ is increasing with respect to epochs e .
- A slow decrease of $I(X, T_e^t)$ with respect to epoch e .

During the compression phase a neural network compresses representation of the input discarding more features that did not help with predicting correct labels. Discarding irrelevant features helps the neural network generalize and produce better predictions for new data.

2.3.4 Compression In Neural Networks

Mutual Information Let's introduce two properties of mutual information

Let's revisit Mutual Information and introduce two properties that will be useful.

Invertible Transformation Let u and v be invertible functions; then,

$$I(A, B) = I(u(A), v(B)) \quad (2.6)$$

Data Processing Inequality Let $A \rightarrow B \rightarrow C$ be a Markov chain; then,

$$I(a, b) \geq I(a, c) \quad (2.7)$$

In the neural network case this implies

$$H(X) \geq I(X, T_e^1) \geq I(X, T_e^2) \geq \dots \geq I(X, T_e^N), \text{ and} \quad (2.8)$$

$$I(X, Y) \geq I(T_e^1, Y) \geq I(T_e^2, Y) \geq \dots \geq I(T_e^N, Y) \quad (2.9)$$

i.e. through out the layers we can only lose or maintain information about the input distribution X and label distribution Y , equality is achieved iff transformation functions f_θ^t are invertible.

Viability of Compression

Let's consider the value $I(X, T_e^t)$, in order to generate the probability distribution X we have assumed that it's uniform, i.e

$$P(X = x_i) = 1/N \text{ for all } i = 1 \dots N \implies H(X) = \log_2(N)$$

we also know the probability distribution of T_e^t given X , i.e

$$P(T_e^t = t | X = x) = \begin{cases} 1, & \text{if } t = F_{\theta(e)}^t(x). \\ 0, & \text{otherwise.} \end{cases}$$

This implies if we have observed the value of X there is no uncertainty of the value of T_e^t , hence

$$H(T_e^t | X) = 0$$

this implies

$$I(X, T_e^t) = H(T_e^t) - H(T_e^t | X) = H(T_e^t)$$

Consider now $H(T_e^t)$,

we know from Equation 2.4

$$\begin{aligned} H(T_e^t) - H(T_e^t | X) &= H(X) - H(X | T_e^t) \\ \implies H(T_e^t) &= H(X) - H(X | T_e^t) \\ \implies H(T_e^t) + H(X | T_e^t) &= H(X) \\ \implies H(T_e^t) &\leq H(X) \end{aligned} \tag{2.10}$$

where equality is achieved iff $H(X | T_e^t) = 0$, i.e

$$P(X = x | T_e^t = t) = \begin{cases} 1, & \text{if } t = F_{\theta(e)}^t(x). \\ 0, & \text{otherwise.} \end{cases}$$

this is the case when if $F_{\theta(e)}^t$ is invertible and every x_i generates a unique layer activation t .

Recall the definition

$$F_{\theta}^t(x) = f_{\theta}^t(f_{\theta}^{t-1}(\dots(f_{\theta}^1(x))))$$

this implies

$$F_{\theta}^t \text{ is invertible } \iff f_{\theta}^i \text{ is invertible for } i = 1 \dots t$$

If equality is achieved in Equation 2.10 and F_{θ}^t is invertible that would imply that we have equalities in Equation 2.8 and Equation 2.9, this would imply that every layer contains the same information content as the input and no compression is can happen.

Let's consider a real world neural network, we have previously defined neural networks to be a sequence of parameterized functions $f_{\theta}^1, f_{\theta}^2, \dots, f_{\theta}^N$. In actual neural networks f_{θ}^t is a matrix that when applied to an activation vector of layer t produces the activation vector of layer $t + 1$.

The matrix is parameterized by weights θ which are controlled by the SGD algorithm, at the start of the training process weights θ are assigned random values, meaning we have completely random matrices, during every iteration process of SGD algorithm the weight are chosen at random and tweaked.

A matrix M is invertible iff its determinant is not zero $\det(M) \neq 0$, A matrix determinant can take values in the range $(-\infty, \infty)$, thus if we consider a random matrix M the probability that it takes any specific value is 0,

$$P(\det(M) = x) = 0, \text{ for } x \in \mathbb{R}$$

specifically $P(\det(M) \neq 0) = 1$, hence $\det(M) \neq 0$ and every random matrix M is invertible.

Given that every random matrix is invertible and that SGD generates random matrices we might conclude that no compression is possible, however this issue is resolvable if we consider weights θ to be a random distribution rather than an instance of a concrete value.

Let us define θ' to be a probability distribution defined by Figure 2.5

Let us define $\theta(e)$ to be an instance of $\theta'(e)$ for any epoch e

```

1  def random_θ'(epoch):
2      if epoch == 1:
3          pick θ ~ multinomialNormal(dimension = d)
4          return θ
5      old_θ = random_θ(epoch - 1)
6      Let θ = One update step of SGD applied to old_θ
7      return θ

```

Figure 2.5: random distribution of $\theta'(e)$, where e is an epoch

-
- consider single it of SGD the new θ is generated by a random process hence it's just a value of a from a random probability distribution.

- $i++$
-

- We have defined a NN to be a sequence of functions f, f, \dots, f . In real neural networks an f is a matrix application that takes in a activation vector of layer i and returns and activation of layer $i+1$
- NN is a collection of functions, where every function is a matrix application on a vector. \Rightarrow they are random

A matrix application is invertible iff A matrix has a determinant zero

- randomness
-

- data processing inequality, every function f , can only lose information about the input, they lose information via collisions

Then $I(T, X) = H(T) - H(T|X) = H(T)$

- Transformation invariant $I(X, Y) = I(f(X), g(Y))$, if f and g are invertible
- f are matrices that transform a vector generated by SGD, they are random so probability that they are invertible is 100
-

-
- justifying compression
 - information plane
-

Before developing a plan for how we are going to realize the project in code we needed to fully understand the ideas presented in the paper:

- We needed to identify the main ideas of the paper and understand why some parts of the paper are not agreed upon in the scientific community. Understand why his ideas are contentious and whether reproducing his experiments could bring more validity to his claims. This involved reading papers published by Tishby and academics who shown an opposing view to him.
- A main tool that the paper relies on is MIE (Mutual Information Estimation). Reading about MIE we quickly understood that MIE is a contentious part of the project as a result we had to do a decent amount of research regarding the subject. MIE is difficult because we are trying to estimate information between two continuous distributions using only a discrete sample set. This area has not seen much academic attention so the tools we ended up using could be greatly improved in the future.

Once we had a reasonable understanding of the ideas in the paper and which areas needed more attention we diverted our attention to figuring out the details of how the experiments were conducted figure out what hyper parameters Tishby decided are important and what assumptions he made whilst devising the experiments.

In addition we needed to find out what resources are available to us online, what programming frameworks we are going to use for the projects implementation, and to think about possible extensions to the project once the success criteria has been achieved.

- Online Resources: The two main papers by Tishby and by Saxe have made their code public online via Github, we made

Online Resources: The two main papers we were looking at has made their code available to the public via Github, the papers are Tishby's paper and the main opposing paper by Saxe.

- Programming frameworks: The original experiment implementation by Tishby has used the Tensorflow framework. We have decided to use the Keras framework as it produces code that is more concise and is easier to read/maintain. Furthermore rewriting the experiments in a different framework means that we cannot rely on the details of Tishby's and potentially avoid any mistakes that may exist in the original implementation.
- Thinking about how we could extend the project helped us understand the scope of the project and what areas were most important and/or interesting to us.

We came up with a couple of extensions before having written any code but the most interesting one only materialized after a good deal amount of work into the project (that is the AS-IF-Random experiment described below)

- Different Datasets : the most straight forward extension to the project just using different dataset to the one Tishby used. This is essentially just varying one of the parameters in the Neural Network. (Implemented)
- Quantized Neural Network : the idea behind this was to only allow single neurons to acquire values in a given range say 1...256. This would make the distribution within a DNN later discrete and hence it would make calculating mutual information straightforward. (Not Implemented)
- As-If-Random : one problem with Tishby's work is that he calculates mutual information for a single epoch at a time which by definition is zero (in his paper he tries to justify the result will explore this later) this extension tries to explore the weights of a neural network as random variables by calculating mutual information for multiple epochs at a time.

Chapter 3

Implementation

3.1 Mutual Information Estimation

As stated before we are trying to validate the fact that compression is happening within the hidden layers in a neural network as such having robust tools to measure mutual information is of utmost importance.

3.1.1 Mutual Information Definition

Mutual information is a measure of dependence between two variables X and Y , it quantifies the number of bits obtained about one when observing the other.

Suppose we have two random variables X and Y if we want to measure mutual information between the two of them we usually refer to one of the following equations:

Using the entropy of the distributions to infer the value

$$I(X, Y) = H(X) - H(X|Y) \quad (3.1)$$

$$I(X, Y) = H(X) + H(Y) - H(X, Y) \quad (3.2)$$

Or calculating the mutual information explicitly from probabilities, for discrete and continuous distributions respectively

$$I(X; Y) = \sum_{y \in \mathcal{Y}} \sum_{x \in \mathcal{X}} p(x, y) \log \left(\frac{p(x, y)}{p(x) p(y)} \right) \quad (3.3)$$

$$I(X; Y) = \int_{\mathcal{Y}} \int_{\mathcal{X}} p(x, y) \log \left(\frac{p(x, y)}{p(x) p(y)} \right) dx dy \quad (3.4)$$

3.1.2 Theoretically undefined

Calculating mutual information is mathematically defined for both continuous and discrete distributions by using one of the functions above. However since we are looking at neural networks we do not have the full distribution – only an empirical sample of the unknown joint distribution $P(X, Y)$. As such we cannot use conventional methods to calculate the mutual information and instead have to use tools to estimate it.

The field of estimating mutual information is relatively new as a result the tools that are available are quite primitive. There has been some papers published that use more advanced techniques to tackle the problem, we will talk about in subsection 3.2.4

3.2 Calculating Mutual Information

3.2.1 Experimental setup

```

1  Algorithm: Generate Information Plane Data
2  Input:
3   $T$  = number of layers
4   $X$  = training data
5   $Y$  = label data
6   $N$  = number of epochs
7  Output:
8   $I_x(\text{epoch}, \text{layer})$ , # mutual information with  $X$ 
9   $I_y(\text{epoch}, \text{layer})$  # mutual information with  $Y$ 
10 Algorithm:
11  $NN = \text{setup\_neural\_network}(T, X, Y)$ 
12 for  $e$  in  $0..N$ :
13      $NN.\text{run\_SGD\_once}()$ 
14     for  $t$  in  $0..T$ :
15         data = []
16         for  $x \in X$ :
17              $\hat{t} = NN.\text{predict}(x).\text{layer\_activation}(t)$ 
18             data.append( $\hat{t}$ )
19          $I_x(e, t) = \text{calculate\_mutual\_information}(\text{data}, X)$ 
20          $I_y(e, t) = \text{calculate\_mutual\_information}(\text{data}, Y)$ 

```

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

- X input is an empirical sample of all possible inputs.
- Y label data. Every $x \in X$ has a corresponding label $y \in Y$.
- $T_{e,i}$ data for a specific epoch e and specific layer i in the neural network. The dataset is generated by feeding every $x \in X$ through the neural network and recording the activations.

3.2.2 Discrete 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 3.5

$$I(X, Y) = H(X) - H(X|Y) = H(X) - \sum_{y \in Y} H(X|Y = y)p(Y = y) \quad (3.5)$$

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 3.2 through to Figure 3.4 outline the full algorithm.

```

1  Algorithm: Mutual Information
2  Input:
3   $X = x_1, x_2, \dots, x_n$ 
4   $Y = y_1, y_2, \dots, y_n$ 
5  Output:  $I(X : Y)$ 
6   $X$  = Bin close values of  $X$  together
7   $Y$  = Bin close values of  $Y$  together
8   $H(X)$  = Calculate entropy of  $X$ 
9   $H(X|Y)$  = Calculate conditional entropy of  $X$  given  $Y$ 
10  $I = H(X) - H(X|Y)$ 

```

Figure 3.2: Pseudo code for computing mutual information refer to Figure 3.3 and Figure 3.4 for entropy computation.

```

1  Algorithm: Entropy - Discrete Method
2  Input:  $X = x_1, x_2, \dots, x_n$ 
3  Output:  $H(X)$ 
4  for  $\hat{x} \in \text{Unique}(X)$ :
5      count = 0
6      for  $x \in X$ :
7          if  $\hat{x} = x$ :
8              count = count + 1
9       $P_x = \text{count} / \text{len}(X)$ 
10  $H(X) = - \sum_{x \in \text{Unique}(X)} P_x \log(P_x)$ 

```

Figure 3.3: Algorithm for computing entropy – Discrete method

```

1  Algorithm: Conditional Entropy
2  Input:
3   $X = x_1, x_2, \dots, x_n$ 
4   $Y = y_1, y_2, \dots, y_n$ 
5  Output:  $H(X|Y)$ 
6   $H(X|Y) = H(X)$ 
7  for  $\hat{y} \in \text{Unique}(Y)$ :
8       $\hat{X} = []$ 
9      for  $x, y \in \text{zip}(X, Y)$ :
10         if  $\hat{y} = y$ :
11              $\hat{X}.append(x)$ 
12          $H(X|Y) = H(X|Y) - H(\hat{X})$ 

```

Figure 3.4: 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

$$H(T|X) = 0 \quad (3.6)$$

$$I(T, X) = H(T) - H(T|X) = H(T) \quad (3.7)$$

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

3.2.3 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 Discrete 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 3.2, Figure 3.3, and Figure 3.4, however there is a change the way entropy is calculated so instead of Figure 3.3 we have Figure 3.5 below.

```

1  Algorithm: Entropy - KDE
2  Input:
3   $X = x_1, x_2, \dots, x_n$ 
4  var = noise variance 0.05 by default
5  Output:  $H(X)$ 
6  dists = compute distance matrix of  $X$ 
7  dists = dists / 2*var # divide every distance by a value
8  # an x is an observation of a hidden layer so, hence it's a vector
9  # which has an associated dimension
10 dim =  $x_1$ .dimension
11 normconst = (dim / 2)*log(2* $\pi$ *var)
12 lprobs = []
13 for row in dists:
14     lprobs.append(log(sum(exp(-row))) - log(n) - normconst)
15 H(X) = mean(lprobs) + (dim / 2)

```

Figure 3.5: 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 $t \in T$.

3.2.4 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 Figure 3.1, 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 Discrete 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 1.1, 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 3.6 and Figure 3.7. 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 3.7
12 else:
13     skip = skip*multiplier

```

Figure 3.6: 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 3.7: 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 3.8 how to compute it or to Figure 3.8 for a graphical representation.

$$D = \max_{t \in T} [\max(I_e(X, t) - I_{\hat{e}}(X, t), I_e(Y, t) - I_{\hat{e}}(Y, t))] \quad (3.8)$$

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

If we consider the Figure 3.8 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 3.8, and scale mutual information values by the entropy as in Equation 3.9.

$$D = \max_{s \in \{X, Y\}} \max_{t \in T} [(I_e(s, t) - I_{\hat{e}}(s, t)) / H(s)] \quad (3.9)$$

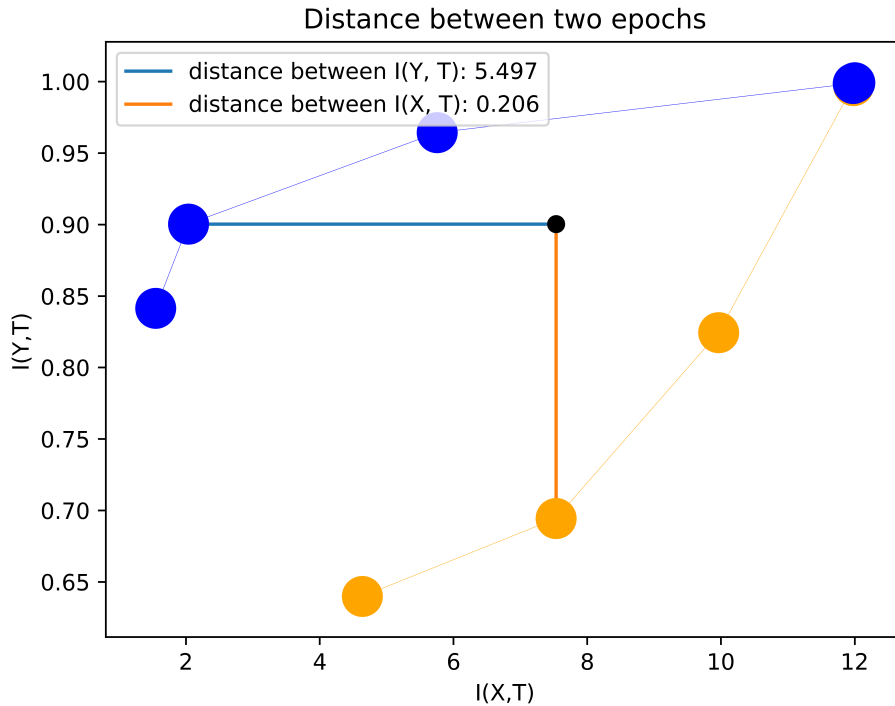


Figure 3.8: Example how distance between two epochs is measured.

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 3.9, 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 3.9: Delta Skip Approximate

3.4 As-If-Random Experiment

Tishby's paper relies heavily on the notion that weights of a neural network behave 'as if' they are random, however his and Saxe's experiments don't capture this idea as much as they could – refer to subsection 4.0.2.

An experiment that better captures the ideas of random is instead of

- Tishby and Saxe didn't capture randomness incredibly well
- We have a better 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 compression
- 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 blah not finished.

3.5 Repository Structure

Chapter 4

Evaluation

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

Appendix A

Project Proposal

Measuring mutual information within Neural networks

Andrius Grabauskas, ag939
Robinson College
Saturday 20th October, 2018

Project Originator: Andrius Grabauskas

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 – 2nd Feb**

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.

- **3rd Feb – 2nd Mar**

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>