

## Measuring mutual information within Neural networks

Andrius Grabauskas, ag939  
Robinson College  
**date**

**Project Originator:** Andrius Grabauskas

**Project Supervisor:** Dr. Damon Wischik      **Signature:**

**Director of Studies:** Prof. Alan Mycroft      **Signature:**

**Overseers:** Dr. Robert Mullins      **Signature:**

Prof. Pietro Lio'      **Signature:**

## 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"<sup>1</sup>

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

---

<sup>1</sup><https://arxiv.org/abs/1703.00810>

more information about the label

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.

## Starting Point

I've watched a talk given by Prof. Tishby gave on this topic at Yandex, no other preparation was done.

## Resources Required

I will use my own laptop to train and evaluate the neural networks as I expect them to be small given the difficulty of measuring mutual information and entropy in neural networks.

For backups I intend to store my work on GitHub and my own personal machine. In case my laptop breaks I'll 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 peers. 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 learn a great deal of Information theory and DNN theory. In order for my results to be valid.

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. Mutual information is known to be tricky to measure when it comes to Continuous random variables.

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

## Success Criteria

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 require:

- Train a neural network
- Measure the mutual information between input the layers and the label.
- analyze the results

## Extensions

The dataset used in the original paper is very non-standard so it would be very interesting to see how the theory holds up when we use it with different datasets.

Different datasets will most likely need different ways to approximate mutual information between layers.

## Schedule

- **15<sup>th</sup> Oct – 25<sup>th</sup> Nov**

Preliminary reading and examining code written provided by Prof. Tishby

- **26<sup>th</sup> Nov – 23<sup>rd</sup> Dec**

Reproduce the Papers results on the original dataset.

- **24<sup>th</sup> Dec – 6<sup>th</sup> Jan**

Examine result and produce graphs.

- **7<sup>th</sup> Jan – 17<sup>th</sup> Feb**

Test the theory on additional datasets

- **18<sup>th</sup> Feb – 17<sup>th</sup> Mar**

Final write up of the thesis