

# Measuring mutual information within Neural networks

Andrius Grabauskas, ag939  
Robinson College  
**Thursday 17<sup>th</sup> January, 2019**

**Project Originator:** Andrius Grabauskas

**Project Supervisor:** Dr. Damon Wischik

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

**Overseers:** Dr. Robert Mullins      Prof. Pietro Lio'

## Abstract

## Introduction

### 1. Summary of results and structure of the paper

### 3. Numerical experiments and results

#### 3.1 Different Mutual Information estimators

##### 3.1.1 Justifying why binning might actually be an appropriate way to measure MI inside neural networks

- binning
- KDE
- KL estimator

- LNN order 1 and 2
- Local Gaussian Approximation - to be implemented

## 3.2 Different Data Sets

- Tishby's
- MNIST
- Fabricated with relevant and irrelevant inputs

## 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 (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

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 paper<sup>2</sup> was produced that tries to suggest that the theory is wrong, however this was dismissed by Tishby & Shwartz-Ziv<sup>3</sup>

## Starting Point

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

---

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

<sup>2</sup>[https://openreview.net/pdf?id=ry\\_WPG-A-](https://openreview.net/pdf?id=ry_WPG-A-)

<sup>3</sup>[https://openreview.net/forum?id=ry\\_WPG-A-&noteId=S1lBxcE1z](https://openreview.net/forum?id=ry_WPG-A-&noteId=S1lBxcE1z)

## 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 paper<sup>4</sup>)

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 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: MNIST<sup>5</sup> and NOT-MNIST<sup>6</sup>.

---

<sup>4</sup><https://arxiv.org/abs/1508.00536>

<sup>5</sup><http://yann.lecun.com/exdb/mnist/>

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

- 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.