

Measuring Mutual Information inside Deep Neural Networks

Andrius Grabauskas - ag939

Friday 8th February, 2019

Quick recap:

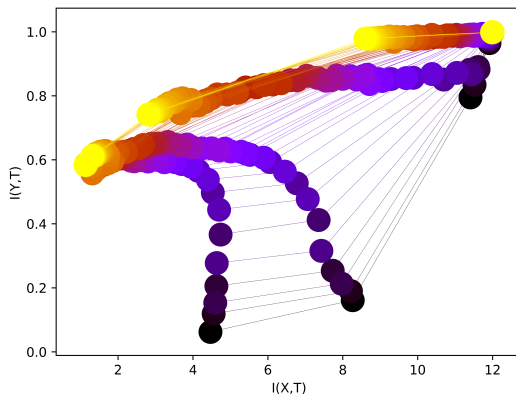
- We do not actually understand or know why DNN's learn.
- N. Tishby suggested that DNN's learn as a result of compression between the layers and the information bottleneck principle.¹
- N. Tishby's claims have been disputed by A. M. Saxe who claimed that his results are due to the parameters he chose.²
- I am here to:
 - ▶ Reproduce Tishby's results provided in his paper.
 - ▶ Explore new ways to verify Tishby's core thesis.
 - ▶ Provide stable code that could be extended and used by other in the future.

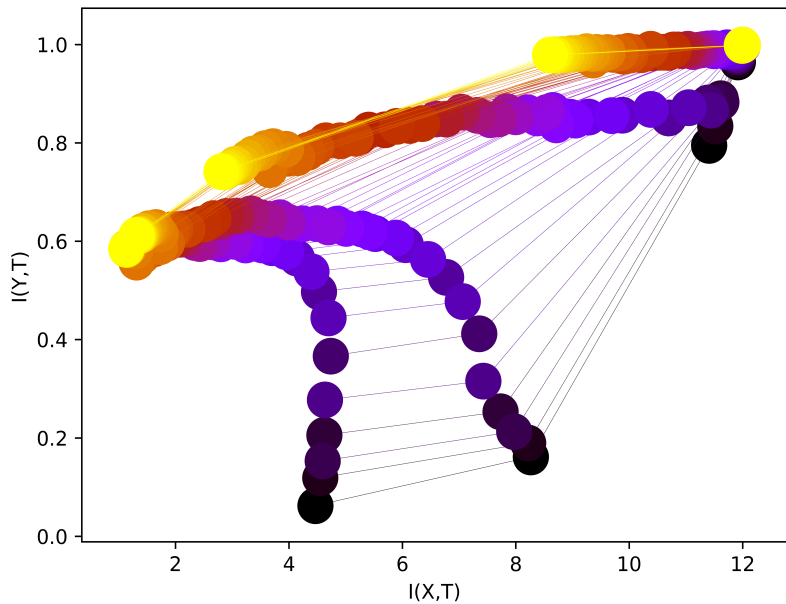
¹<https://arxiv.org/abs/1703.00810>

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

Success criteria

- Reproduce results in Tishby's paper
- Essentially reproduce this image and show that the results are stable with varying parameters.





Current progress

- Reimplemented Tishby's code and reproduced his results.
- Found evidence that Tishby's Mutual Information Estimator (MIE) might not be completely sound.
- Had lots of trouble with finding an alternative MIE as Neural Networks are have high dimensionality which doesn't play well with some MIE.
- Settled with Kernel Density Mutual (KDE) MIE, which agreed with Tishby's results. (Although it is by far not a perfect measure of information)

Remaining work, Issues

- Show that the results hold for a bigger dataset such as MNIST
- Clean up and comment the code such that it is easy to understand and extensible by others.
- At the core of Tishby's thesis is the notion that Weights of a NN are random
 - ▶ What does it actually mean for a weight of a NN to be random ?
 - ▶ given that at every point in time a weight is fixed and non changing.