

Artificial Neural Networks: Structure and Behavior

Kendall Lowrey, Thanh Tam Nguyen

Abstract—Artificial Neural Networks have traditionally been treated as black boxes, both in their development and in their use. We present a method to discover the internal structure of neural networks by visualizing activation properties of the network with respect to input data: co-activations of multiple neurons. Our method combines statistical analysis techniques with a modified Sankey Diagram to show flow of data through neural networks unlike previous visualizations methods. Implications for this technique beyond behavioral and structural visualization include the optimization of an artificial neural network through parameter reduction and further understanding of their processing.

1 INTRODUCTION

Neural networks, both artificial and biological, are layers of individual nodes that are used to estimate or approximate functions operating on some input data. The function they approximate are trained or learned in similar ways: neurons that behave in the intended way are preserved while neurons that produce errors are adjusted. In biological networks, this training takes the form of forming stronger connections between neurons of different layers within a network such that when examined, the behavior or the network is directly related to its structure. Artificial networks, however, are trained on vast amounts of data, usually more data and faster than a human may comprehend. As a result, artificial neural networks (ANN) are often treated as a black box in their usage and study [2].

As ANNs are merely layers of nodes acting as matrix transformations, it may be intuitive to visualize them by laying out the nodes in vertical groupings with lines as their connecting arithmetic functions, as in figure 1. While this is indeed the 'structure' of the neural network, it says nothing to its behavior or function. Extensions to this idea may elucidate more complicated neural networks where layers have non-linear function behavior. Unfortunately, these representations still do little more than descriptive wording or pseudo-code to explain their behavior.

Historically, Hinton or Bond diagrams were used to abstractly represent neural networks and their neuron weights [7]. These methods have a higher learning curve and only work for small networks; modern neural networks are reaching a size that includes billions of parameters.

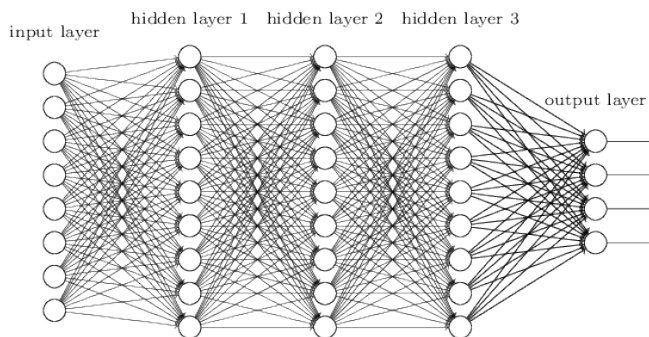


Fig. 1. Simple Neural Network Visualization

Recent developments in ANNs, specifically in using Convolutional Neural Networks [5], offer a different visualization method due to

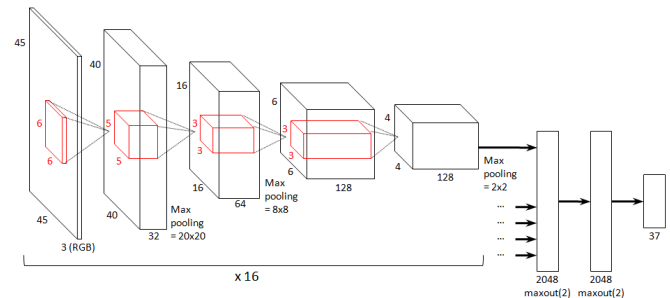


Fig. 2. Convolutional Neural Network

their inherent use of images as their data inputs. Successive layers of transformations in the neural network take the input image and modify it such that the final layer can make a safe decision as to what the input image contains. This requires the convolution of simple image primitive kernels over the input image. These kernels capture image features such as horizontal or vertical lines, points, and gradients, among others. While it is fascinating that these kernels are trained in the convolutional neural network from just training data alone, it is left to the user's imagination to integrate them among the entire network's behavior. Modern software tools may assist in this endeavor, but serves as validation of the networks trained behavior more than anything [3].

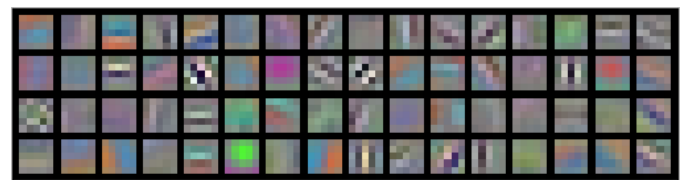


Fig. 3. Convolution primitives

Other than layers and their transforms, a neuron's connections and weights can be visualized in a number of ways. Other than by coloring the vertices's in a graph layout (which would quickly become crowded or unmanageable), a superimposed matrix can be used to show how a group of neurons are interconnected [1] [4]. This representation allows for a quick understanding for how the static properties of a network may be visualized, and is conveniently two dimensional.

In short, neural network visualizations frequently represent the definition of a particular network as nodes in a connected graph, specific layer's transformations on their inputs, or highlighting the connections of the graph and represent this data in a variety of ways.

2 CO-ACTIVATION VISUALIZATION OVERVIEW

The approach we take to glimpse the the internals of the black box are to consider the behavior of the network's components in the presence

- *GitHub*: klowrey E-mail: klowrey@cs.washington.edu
- *GitHub*: tam-nt E-Mail: nttam15@uw.edu
- *University of Washington, CSE 512, Spring 2015*

in their figure’s captions, we can briefly overview the results. This method of neural network structure produces visualizations that suggests the underlying structure of a neural network indeed affects its function. These results lets us make decisions and judgements as to a network’s ability to function well.

3.1 Conclusion and future work

We consider this a proof of concept into opening the black box that is a neural network and suggest more avenues of research. From what we understand, this kind of holistic network view has not been attempted before with success, and may prove useful in rendering the behavior of networks for researchers.

The most obvious next step would be to visualize the activations of the network after it has been structure with this method. As discussed in 8, certain areas of this visualization may activate when specific input conditions are met. Visualizing the effect of these input conditions would be an excellent addition to this method, but requires network specific support software to provide this functionality. Our robot control network and training software would have to be converted into Javascript, for example. For the MNIST network, we could see which neurons were activated in response to similar inputs, discovering which neurons are responsible for specific image features.

Another functional option for this research would be to visualize the network as it is being trained. We have discussed that an untrained network is indeed worth visualizing, but for very long training runs, we could use this method to give the researcher feedback beyond just a loss function metric. Visualizing the whole network during training could give insights into which parameters to tweak to sufficiently generalize the training network beyond what a singular loss value could provide.

REFERENCES

- [1] B. Alper, B. Bach, N. Henry Riche, T. Isenberg, and J.-D. Fekete. Weighted graph comparison techniques for brain connectivity analysis. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, pages 483–492. ACM, 2013.
- [2] J. M. Benítez, J. L. Castro, and I. Requena. Are artificial neural networks black boxes? *Neural Networks, IEEE Transactions on*, 8(5):1156–1164, 1997.
- [3] D. Bruckner. MI-o-scope: a diagnostic visualization system for deep machine learning pipelines. Technical report, DTIC Document, 2014.
- [4] S. Hinze and U. Schmoch. Opening the black box. In *Handbook of quantitative science and technology research*, pages 215–235. Springer, 2005.
- [5] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [6] L. Van der Maaten and G. Hinton. Visualizing data using t-sne. *Journal of Machine Learning Research*, 9(2579-2605):85, 2008.
- [7] J. Wejchert and G. Tesauero. Visualizing processes in neural networks. *IBM Journal of Research and Development*, 35(1.2):244–253, 1991.

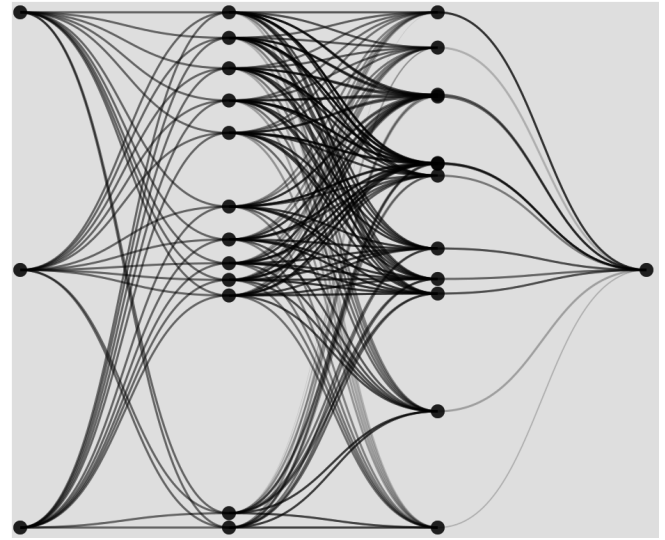


Fig. 5. The untrained XOR network. The network consists of two layers of 12 neurons each. As the network is untrained, we can see that there is very little visual clustering after the tSNE processing. Also, the links between the layers are almost uniformly dark. We take this to mean that the neurons are activating randomly due to the random weights they are initialized with: there is no structure to observe. This is excellent to know because we can begin network training knowing that the network has been initialized with a high degree of randomness – training will be more likely to be successful.

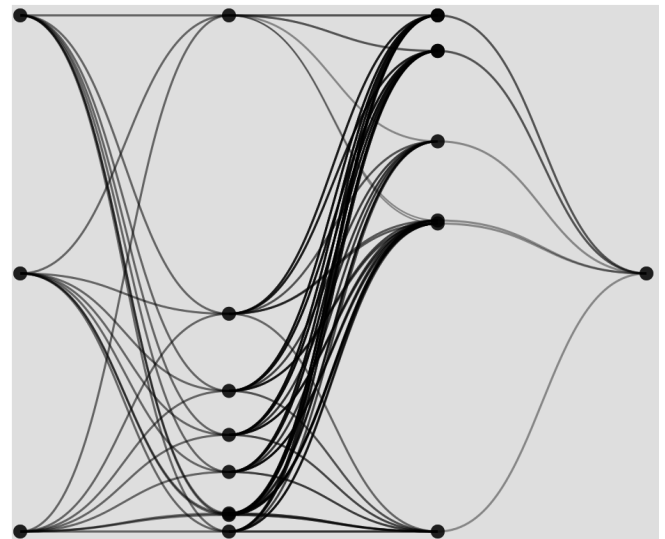


Fig. 6. The trained XOR network features distinct groupings both within a layer and between layers, suggesting that learning has taken place and thus affected the network’s structure. While learning can be verified numerically, we can use the visualizations to experiment further with the network beyond numerical results. For instance, we can observe neuron groupings and reduce the number of neurons potentially making the network more efficient.

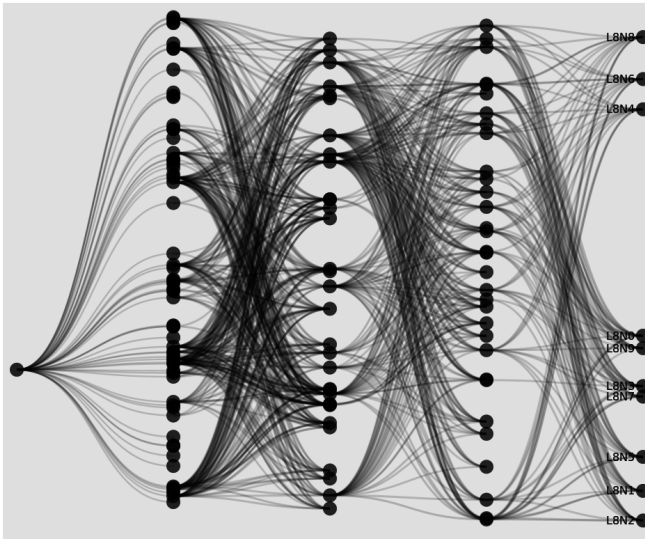


Fig. 7. The MNIST network of hand-written digits. There are many groupings throughout the network of 64, 32, 32 neurons, before a soft-max layer to the 10 different classes. We can see groupings in the output layers that reflects some of the mis-predictions other networks may face with the MNIST dataset. We also observe that through the link threshold there are some neurons that do not significantly contribute to the next layer; this could be cause to remove them from the network much like the neuron group consolidation idea demonstrated by the simplified XOR model in 10

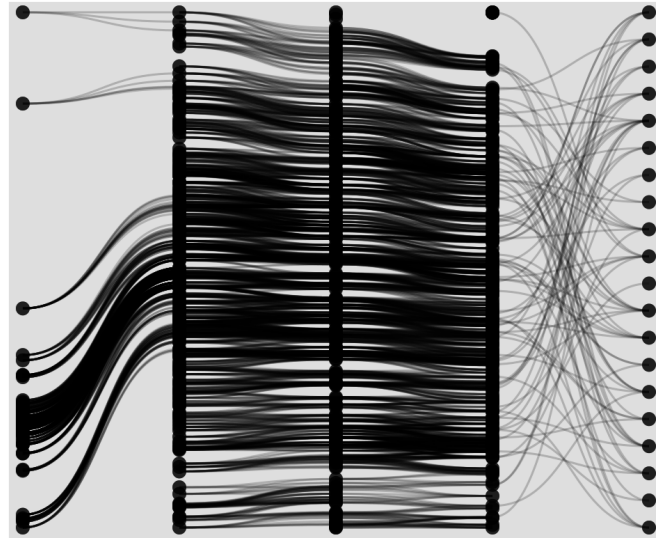


Fig. 9. The humanoid robot running control network. This network does not feature the same split as the walking robot control network, but the actions of the robot are more significantly varied – it does not hit a limit cycle in its brief dashing between locations. While we know training has happened (the links are not in random directions but focused), this particular visualization is not as useful as others.

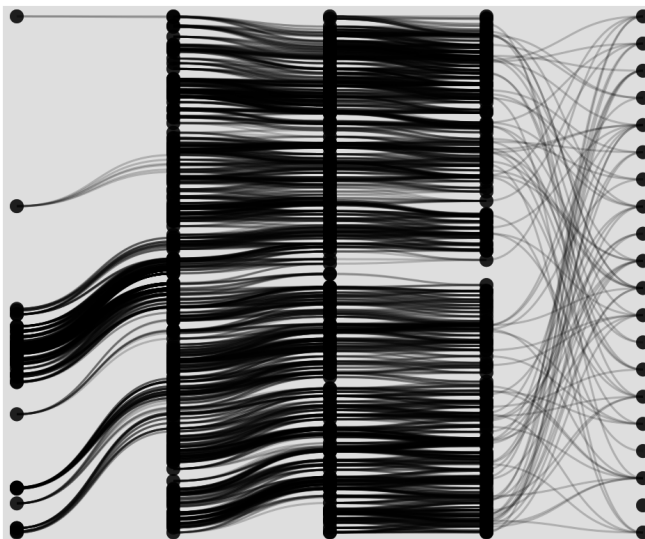


Fig. 8. The humanoid robot walking control network. An interesting feature of this network is the split at the third layer. As this is a bipedal robot that features a limit cycle walking gait (that was learned through trajectory optimization), we presume this split represents the symmetry encountered during walking, although more investigation would be required to confirm this.

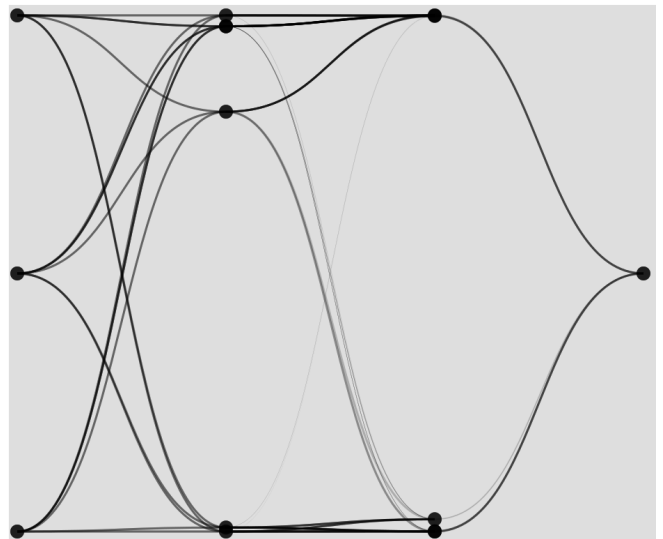


Fig. 10. We simplify the trained XOR network based on the visualized groupings and create a smaller network of 6 neurons and 5 neurons. The loss function for both networks trained on three-way XOR were the same (0.009) suggesting that the performance is comperable with less neurons. This network simplification is a potential application of this visualization technique.