
Do deep neural networks possess concept space grid cells?

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Place and grid-cells are known to aid navigation in animals and humans. Together
2 with concept cells, they allow humans to form an internal representation of the
3 external world, namely the concept space. We investigate the presence of such a
4 space in deep neural networks by plotting the activation profile of its hidden layer
5 neurons. Although place cell and concept-cell like properties are found, grid-cell
6 like firing patterns are absent thereby indicating a lack of path integration or feature
7 transformation functionality in trained networks. Overall, we present a plausible
8 inadequacy in current deep learning practices that restrict deep networks from
9 performing analogical reasoning and memory retrieval tasks.

10 1 Introduction

11 Cells in the hippocampal region of the rat brain signal the animal’s location in space. This phenomenon
12 has provided support to the idea that the rat hippocampus operates like a cognitive map [1]. Specific
13 cells, termed as the “place cells” are known to selectively fire when the animal is in certain locations
14 in the environment [1]. The discovery of place cells was followed by the discovery of “grid” cells,
15 that selectively fire at multiple discrete and regularly spaced locations [2]. Extensive animal studies,
16 especially in rats, involving electrophysiological recording have shown the existence of grid cells
17 and place cells in localized brain structures, namely the medial entorhinal cortex and hippocampus
18 respectively. In humans, the role of the medial temporal lobe and particularly the hippocampus in
19 creation, consolidation and recall of declarative memories has been established [3]. Specific cells in
20 the medial temporal lobe (MTL), called the “concept” cells, seem to fire when a particular concept is
21 presented to the subject [4]. These findings established that humans are indeed capable of forming an
22 internal representation of the external world. Navigating this internally represented abstract concept
23 space might indeed form the basis of memory retrieval or context-driven reasoning.

24 Deep learning (DL) has brought about a new dawn in artificial intelligence by enabling models to
25 learn representations of data with multiple levels of abstraction. However, most of these models lack
26 the flexibility to perform relational reasoning. These issues have rekindled interest among researchers
27 to understand and replicate the brain’s learning process. The success of deep neural networks (DNN)
28 is generally attributed to their ability to identify optimal discriminative features in a given dataset.
29 This learnt feature space can be thought of as a “concept” space for the DNN. The ability to perform
30 relational reasoning depends on the network’s ability to navigate over this concept space. To solve
31 analogical problems, the network needs to apply the desired feature transformation to arrive at the
32 correct solution. For instance in the image space, the network needs to navigate the visual concept
33 space to understand that an *OR* operation between a “dog image” and a “ball image” could lead to an
34 image of “a dog playing with a ball”.

35 The ability to navigate the concept space depends on the properties of the constituent neurons. Akin
36 to memory retrieval or context-driven reasoning in humans, DNN would need neurons that show

37 place and grid-cell like activation properties. Artificial neurons that respond to specific classes of
38 images could be considered as showing place cell-like activity in the concept space. Consequently,
39 all final layer neurons of a DNN trained for a classification problem would be showing place-cell like
40 activity. However, the ability to navigate the concept space relies on the DNN’s ability to perform
41 “path integration”, or the ability to localize its position in the concept space given specific feature
42 transformations. Grid cells are known to be responsible for path integration in rodents and humans.
43 Therefore, the ability of the DNN to solve analogical problems relies on the presence of grid cell-like
44 properties among DNN neurons.

45 In this work¹, we aim to investigate a particular class of DNN, namely convolutional neural network
46 (CNN). We investigate the firing properties of neurons, specifically in the final and pre-final layers, to
47 understand their activation patterns while the network performs classification.

48 2 Methods and Results

49 We trained a CNN to identify hand-written digits 0 – 9 in the MNIST dataset. The dataset has 60,000
50 images for training and the network was tested on 10,000 images. The network architecture consists
51 of two 2D convolutional layers (filter size 5) followed by 2D Max-pooling layers. The output of the
52 convolutional layers was fed to a dropout layer and further input to two fully connected layers with 50
53 and 10 hidden units. Rectified Linear Unit (ReLU) was used as the activation function. The network
54 was implemented on PyTorch 1.1.0 and trained to optimize negative log-likelihood loss using the
55 Adam optimizer. The network was trained using a learning rate of 0.001 for 10 epochs (training batch
56 size = 128). At the end of 10 epochs, the network’s accuracy on the testset was **98.79 %**.

57 Once the network was trained to yield a satisfactory classification accuracy on the testset, we
58 investigated the activation properties of constituent neurons. We limited our analysis to the pre-
59 final and final layers (Layer 6 & 7 respectively) only. To visualize the concept space, we used
60 the t-distributed Stochastic Neighbourhood Embedding (tSNE) [5] to generate a low-dimensional
61 representation (2D representation here) of the feature space. We explored two concept spaces – one
62 generated from the tSNE plot of raw image data and the other generated from the activation of the
63 final layer neurons of the network. The first corresponds to the image concept space and the second
64 corresponds to the concept space learned by the network, also referred hereafter as the network
65 concept space. We generated the aforementioned spaces for the testset images (10,000 points). Figure
66 1 shows the two concept spaces.

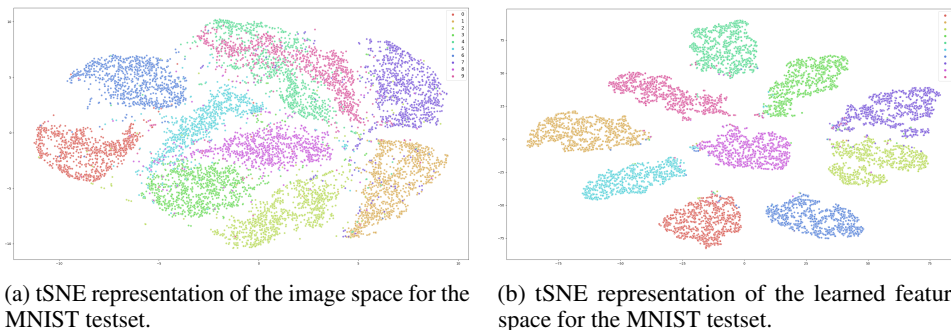
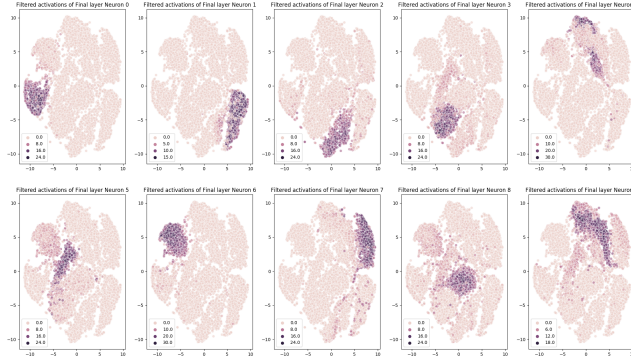


Figure 1: Concept spaces obtained from tSNE representations of higher dimensional data in 2 dimensions.

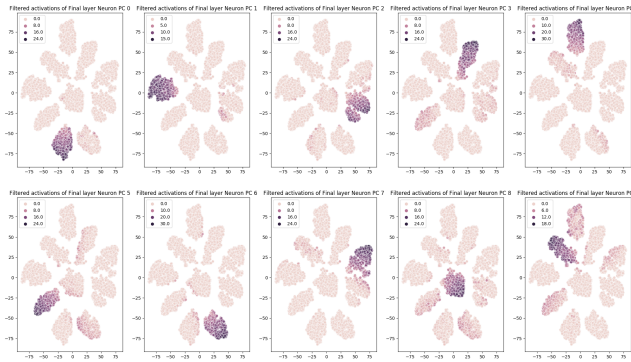
67 To observe the firing property of the neuron across the concept space, the tSNE plots were color-coded
68 using the respective neuron’s activation values. Figure 2 shows the firing pattern of Layer 7 neurons.
69 All the final layer neurons showed place-cell like activity, with preference to one cluster. Given the
70 training paradigm, this was expected. Further, we extended our analysis to observe the firing pattern
71 of the previous layer.

72 The pre-final layer had 50 neurons. However, none of the neurons had grid-cell like firing property.
73 We used Principal Component Analysis (PCA) and obtained 11 PCs that significantly explained the

¹Code and results available here



(a) Activation patterns of final layer neurons in image space



(b) Activation patterns of final layer neurons in network space

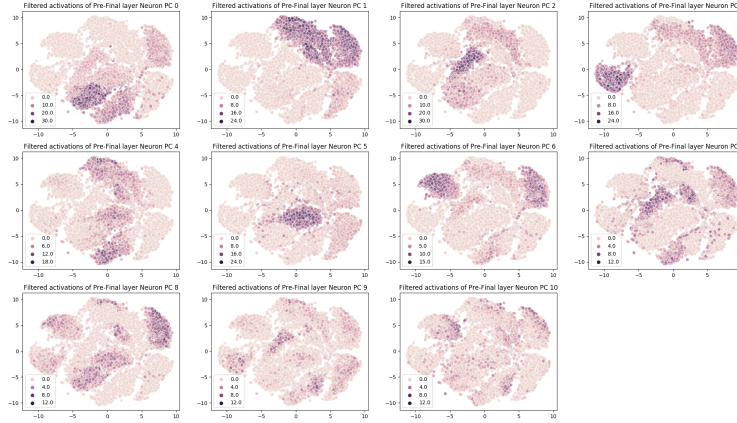
Figure 2: Activation patterns of the final layer neurons represented using color coding on the tSNE plots. All non-negative activities are clamped to 0 as they correspond to highlight regions with high neuronal activation, akin to firing rates.

74 firing properties of these neurons. The firing patterns of these 11 PC neurons are shown in Figure 3.
 75 Since the firing patterns of the PC neurons are obtained from a linear combination of those of the
 76 pre-final neurons, the absence of a grid-like firing pattern in the 50 neurons is reflected in the firing
 77 patterns shown in Figure 3. Instead these neurons show preferential firing to multiple clusters, thus
 78 indicating that these neurons respond to images of multiple digits.

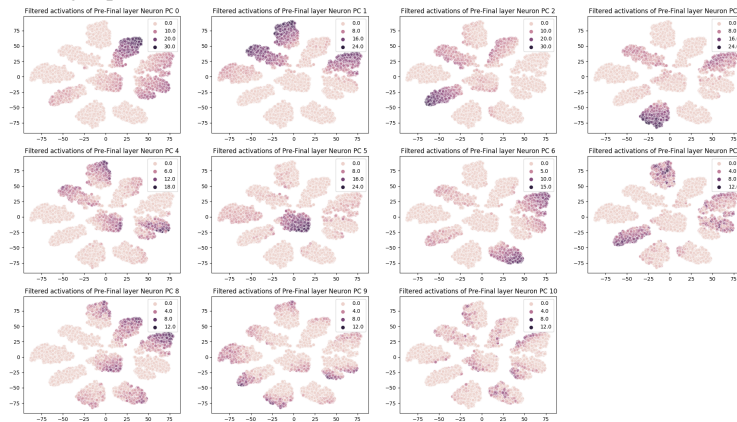
79 3 Discussions & Conclusion

80 In this work, we explored the firing patterns of final and pre-final layer neurons in a CNN trained
 81 to identify hand-written digits in the MNIST dataset. Although the final layer neurons showed
 82 place-cell like firing patterns in the learned concept space, the pre-final layer neurons failed to show
 83 any grid-cell like firing patterns. The distinction in firing patterns is clearer in the network space
 84 as compared to the image space, thus indicating that the network coded for certain regions in the
 85 concept space. However, the absence of grid-cell like firing indicates that the network failed to learn
 86 the transformations and manifolds in the concept space. This would imply that the network would be
 87 unable to navigate this space and thereby perform feature transformation operations in this space.
 88 Thus, the trained CNN would fail to perform analogical problems, similar to other deep models
 89 trained on more complicated datasets.

90 The network analyzed here was a fairly simple and shallow network as compared to practical deep
 91 networks. However, it serves as a toy example to illustrate the inadequacy of current training practices
 92 that restrict deep networks from performing analogical reasoning tasks. Since the images were
 93 shown to the network in isolation, the network was unable to learn the information of manifolds and



(a) Activation patterns of the principal components of pre-final layer neurons in image space



(b) Activation patterns of the principal components of pre-final layer neurons in network space

Figure 3: Activation patterns of the principal components of pre-final layer neurons represented using color coding on the tSNE plots.

94 relative transformation among the clusters. Adding examples that lie on manifolds in this space could
 95 alleviate this problem. More recently, Hill et al. proposed a framework to incorporate analogical
 96 reasoning in deep networks using a modified training strategy [6]. It would be interesting to observe
 97 the firing patterns of neurons in such a network, as done here, to inspect if the network navigated the
 98 concept space similar to the animal brain.

99 References

- 100 [1] John O’keefe and Lynn Nadel. *The hippocampus as a cognitive map*. Oxford: Clarendon Press, 1978.
- 101 [2] Edvard I Moser, Yasser Roudi, Menno P Witter, Clifford Kentros, Tobias Bonhoeffer, and May-Britt Moser.
 102 Grid cells and cortical representation. *Nature Reviews Neuroscience*, 15(7):466, 2014.
- 103 [3] Howard Eichenbaum. Hippocampus: cognitive processes and neural representations that underlie declarative
 104 memory. *Neuron*, 44(1):109–120, 2004.
- 105 [4] Rodrigo Quian Quiroga. Concept cells: the building blocks of declarative memory functions. *Nature*
 106 *Reviews Neuroscience*, 13(8):587, 2012.
- 107 [5] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning*
 108 *research*, 9(Nov):2579–2605, 2008.
- 109 [6] Felix Hill, Adam Santoro, David GT Barrett, Ari S Morcos, and Timothy Lillicrap. Learning to make
 110 analogies by contrasting abstract relational structure. *arXiv preprint arXiv:1902.00120*, 2019.