
Compact representations and pruning in residual networks

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We show that residual networks encode their input signals in the transient dynamics
2 of the neurons in each layer. These representations are similar for inputs from
3 the same class, and distinct for inputs from different classes. Based on the neural
4 transient dynamics, we provide a sufficient criterion to determine the depth of such
5 networks during training. This criterion is based on the convergence of the neural
6 dynamics in the last two successive layers of the residual block. This method
7 compresses the depth of the network and removes unnecessary deep layers.

8 1 Introduction

9 Residual networks (Resnets) [1] have been more successful in classification tasks in comparison
10 with many other standard methods. This success is attributed to the skip connections between layers
11 that facilitate the propagation of the gradient throughout the network, and in practice allow very
12 deep networks to undergo a successful training. Apart from mitigating the gradient problem in deep
13 networks, the skip connections introduce a dependency between variables in different layers that can
14 be seen as a system state. This novelty provides an opportunity for interesting theoretical analysis of
15 their functioning, and has been the underlying pillar for some interesting analysis of such networks
16 from a dynamical system point of view [2, 3, 4, 5, 6, 7, 8].

17 Some studies on Resnets have focused on tracking the features layer by layer [9, 10], and have chal-
18 lenged the idea that deeper layers in neural networks build up abstract features that are different than
19 those formed in lower layers. One supporting evidence for this challenge comes from lesion studies
20 on Resnets [11] and Highway networks [12] which show that after the network is trained, perturbing
21 the weights in the deep layers does not have a fundamental effect on the network performance, and
22 therefore, does not bring the performance to chance level. However, changing the weights which
23 are closer to initial layers, have a more damaging effect. Empirical studies in [9, 10] suggest an
24 alternative explanation for feature formation in deep layers; that is, successive layers estimate the
25 same features which, along the depth of the network, are more refined, and yield an estimate with
26 smaller standard deviation than earlier layers. Our approach in this paper is similar to the latter
27 studies, however, to understand the classification mechanism in Resnets, we focus on the role of the
28 intrinsic dynamics of the residuals over different layers of the network. Our study supports the idea in
29 [9, 10] by showing that features in different layers of a Resnet are formed by the transient dynamics
30 of residuals that may converge to their steady state values if they are stable. Based on this finding, we
31 suggest an algorithm that estimates the depth of the network adaptively during training.

32 2 Neural dynamics in Resnets

33 We consider a dense Resnet with N input dimensions, and arbitrary T layers with exactly N neurons
34 at each layer. In this network, the activity of neuron i at layer t is represented as $y_i(t)$, and the

35 activity of all neurons in the same layer is represented by the vector $\mathbf{y}(t)$. After the integration
 36 of the output from layer $t - 1$, the output of layer t is represented by $\mathbf{x}(t)$. The components of
 37 these residuals $\mathbf{y}(t)$ are calculated based on a linear function of $\mathbf{x}(t)$, i.e. $z_i(t) = \sum_{j=1}^N w_{ij}(t)x_j(t)$
 38 followed by a nonlinear function $f(z_i)$. Any hidden layer t represents a sample of the dynamical
 39 states \mathbf{x} after t steps. This implies that the network at different layers calculates samples of $\mathbf{x}(t)$.
 40 Input data is considered as the initial condition of the system, and is depicted by $\mathbf{x}(0)$. Interpreting
 41 the network as a dynamical system which evolves throughout the layers, the dynamics of neural
 42 activations are $\mathbf{x}(t + 1) = \mathbf{x}(t) + \mathbf{y}(t + 1)$, where $\mathbf{y}(t)$ is the output of the neurons, and in the
 43 rest of the paper, they are called "residuals". This equation implies a difference equation for the
 44 variable $\mathbf{x}(t)$, that is $\mathbf{x}(t + 1) - \mathbf{x}(t) = \mathbf{y}(t + 1)$. The left side of this equation resembles the forward
 45 Euler method of derivative of a continuous system, when the discretization step is equal to 1. This
 46 approximates a continuous system with dynamics that follow $\dot{x}_i(t) = y_i(t)$. The latter equation
 47 implies $x_i(t) = \int_0^t y_i(\tau)d\tau + x_i(0)$ where $x_i(0)$ stands for the input data that neuron i receives. In
 48 other words, $\mathbf{x}(t)$ sums up the input data as well as the activities of the neurons (residuals) over the
 49 layers. This signal feeds the next block in the network, or the classifier in the output layer.

50 To study the properties of the neural activities in each layer that shape the cumulative signal $\mathbf{x}(t)$, we
 51 considered a 784 dimensional network with the MNIST dataset as inputs. After training a 15-layer
 52 Resnet using the back-propagation algorithm, we studied the dynamics of the residuals (the signal
 53 of neural activities from the input layer up to the last layer). We observed rapid changes in the first
 54 initial layers of the network, and more steady behavior close to the final layers (figure 1A, B, see also
 55 [13] for more details on the layer-dependent dynamics of the residuals and their fixed points). The
 56 dynamics of the residuals do not change significantly after the 6th layer (figure 1A), and the standard
 57 deviation of the trajectories for different samples approaches zero thereafter (figure 1B). This implies
 that different trajectories for each sample converge to the same fixed point.

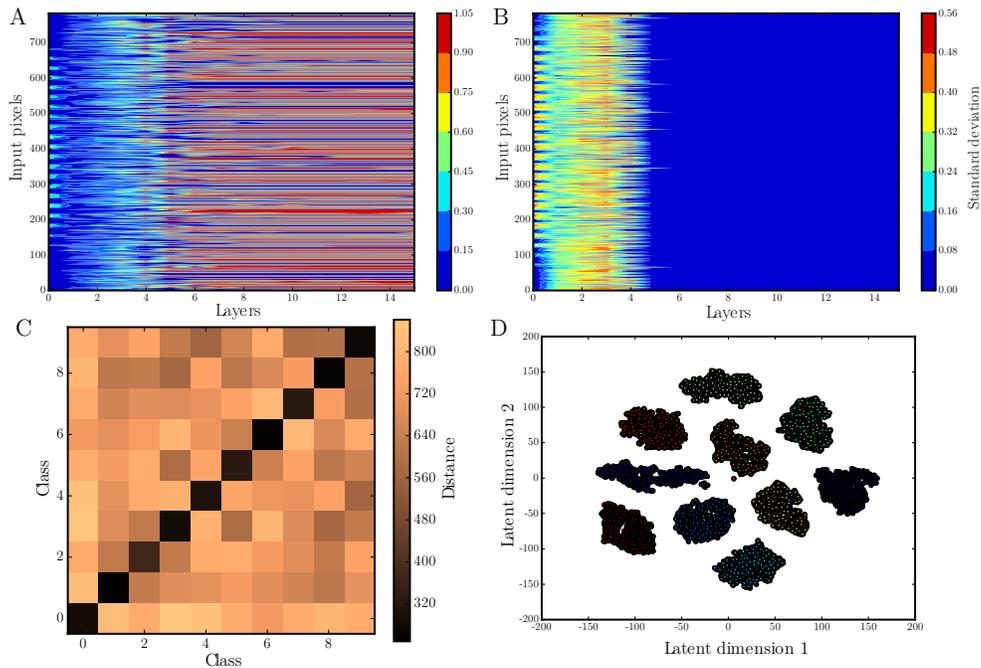


Figure 1: A: Mean of the residuals for 200 samples of each class from the MNIST dataset as a function of network depth (extracted from [13]). B: Standard deviation of the residuals (extracted from [13]). C: Distance matrix for 200 samples for each class, normalized by the number of samples. D: Dimensionality reduction of the cumulative signals for each class, at the final layer of the network, before the classifier.

58

59 In order to show the similarities and differences of the internal dynamics corresponding to each class,
 60 we compared the l_2 distance between neural trajectories of all neurons for 200 examples from each
 61 class. As illustrated in figure 1C, trajectories within each class have a smaller l_2 distance compared

62 with trajectories for samples from different classes. This shows that the network specifies different
 63 forms of neural trajectories for inputs from each class, that are distinct from each other. Moreover,
 64 in order to generalize the classification performance for inputs from the same class, this difference
 65 for those trajectories is less, meaning that inputs from the same class are more similar. To compare
 66 the cumulative signal $\mathbf{x}(t)$ among different classes, at the final hidden layer of the Resnet before the
 67 classifier, we employed the t-SNE algorithm [14] to visualize the properties of these high dimensional
 68 signals in a two-dimensional plot. As depicted in figure 1D, signals corresponding to different classes
 69 are mapped and clustered in different regions of the two-dimensional space. This reflects how Resnets
 70 encode inputs in different and distinct forms of their neural transient (layer-dependent) dynamics
 71 (even though residuals for different classes converge to identical values in this particular example).

72 3 Compressing network’s depth

73 In the MNIST network example, we observed that the neural trajectories (residuals) show less
 74 variabilities at deeper layers, some of which have already converged to their steady state values
 75 (figure 1A, B). Under this condition, the cumulative signal $\mathbf{x}(t)$ receives similar components in
 76 successive deep layers, and hence no new information about the variabilities in neural trajectories
 77 are encoded in this signal. In other words, almost always constant numbers (due to convergence) are
 78 added to this signal without providing new information about the input-induced neural transitions.
 79 Therefore, after achieving this state of neural transient dynamics, it seems viable to stop the forward
 80 propagation of information, and cut the extra layers of the network without losing much information
 81 in the cumulative signal. This provides the basis for an algorithm that adaptively sets the depth of the
 82 network considering the difference between neural activities of the last two successive layers of the
 83 network. The algorithm as such is the following:

```

84 while loss function is not minimum do
    for each epoch of the training data do
        residuals of the last hidden layer in the block  $\rightarrow r_1$ 
        residuals of the second last hidden layer in the block  $\rightarrow r_2$ 
        if  $l_1$  norm for  $r_1 - r_2 < threshold$  then
            | remove the layer corresponding to  $r_1$ 
        end
    end
end

```

85 which compares the l_1 norm of the difference between the neural activities of the last two layers
 86 [13]. If this difference is less than a threshold, the activities of the neurons could be considered
 87 approximately identical. If this condition is fulfilled, the last layer of the network is removed. To
 88 compress the network as much as possible, this algorithm was applied on the MNIST network with
 89 shared weights between layers [6], and resulted in a 5-layer network (for threshold = 0.01) without
 90 any significant changes in the classification accuracy. Note that the convergence criterion for the
 91 residuals is a sufficient condition to remove the last hidden layer, not a necessary condition. This
 92 means that even shallower networks might still be able to classify the inputs with the same accuracy,
 93 however, our algorithm does not provide the necessary conditions to achieve the same level of
 94 accuracy for those cases. This algorithm is more efficient than fully training shallower networks first,
 95 and then adding extra layers to compensate for a high-value loss function and retraining the network.

96 4 Conclusion

97 In this letter, we showed the importance of neural transient dynamics on input classification in Resnets.
 98 It was demonstrated here that the cumulative signal of the residuals has a distinct state for each
 99 input class, in a high dimensional state space of the neural network. Also, the neural trajectories
 100 across layers are similar for inputs of the same class while different for inputs from two distinguished
 101 classes. This form of input representation in the transient dynamics of neural trajectories in Resnets
 102 can potentially underlie more efficient methods of training in the future. Based on the convergence of
 103 the neural dynamics to a steady state, we proposed an algorithm that determines the sufficient number
 104 of layers for the network during training. This can be considered as network depth compression
 105 without losing the classification accuracy significantly.

References

- 106
107 [1] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep Residual Learning for Image Recognition.
108 *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016.
- 109 [2] Marco Ciccone, Marco Gallieri, Jonathan Masci, Christian Osendorfer, and Faustino Gomez. NAIS-NET:
110 Stable Deep Networks from Non-Autonomous Differential Equations. *arXiv*, 2018.
- 111 [3] Bo Chang, Lili Meng, Eldad Haber, Frederick Tung, and David Begert. Multi-level Residual Networks
112 From Dynamical Systems View. *ICLR*, pages 1–14, 2018.
- 113 [4] Eldad Haber and Lars Ruthotto. Stable Architectures for Deep Neural Networks. *arXiv*, 2017.
- 114 [5] Yiping Lu, Aoxiao Zhong, Quanzheng Li, Massachusetts General Hospital, and Bin Dong. Beyond Finite
115 Layer Neural Networks: Bridging Deep Architectures and Numerical Differential Equations. *arxiv*, pages
116 1–15, 2017.
- 117 [6] Qianli Liao and Tomaso Poggio. Bridging the Gaps Between Residual Learning, Recurrent Neural
118 Networks and Visual Cortex. *arXiv*, (047):1–16, 2016.
- 119 [7] Lars Ruthotto and Eldad Haber. Deep Neural Networks motivated by Partial Differential Equations. *arXiv*,
120 pages 1–7, 2018.
- 121 [8] Pratik Chaudhari, Adam Oberman, Stanley Osher, Stefano Soatto, and Guillaume Carlier. Deep Relaxation:
122 partial differential equations for optimizing deep neural networks. *Proceedings of the 34th International
123 Conference on Machine Learning, Sydney, Australia, PMLR*, 2017.
- 124 [9] Klaus Greff, Rupesh K. Srivastava, and Jürgen Schmidhuber. Highway and Residual Networks learn
125 Unrolled Iterative Estimation. *ICLR*, (2015):1–14, 2017.
- 126 [10] Brian Chu, Daylen Yang, and Ravi Tadinada. Visualizing Residual Networks. *arxiv*, 2017.
- 127 [11] Andreas Veit, Michael Wilber, Serge Belongie, and Cornell Tech. Residual Networks Behave Like
128 Ensembles of Relatively Shallow Networks. *NIPS*, pages 1–9, 2016.
- 129 [12] Rupesh Kumar Srivastava, Klaus Greff, and Jürgen Schmidhuber. Training Very Deep Networks. In *NIPS*,
130 pages 1–9, 2015.
- 131 [13] Anonymous. Residual Networks classify inputs based on their neural transient dynamics. *ICLR 2018
132 Submission*, (2018):1–11, 2018.
- 133 [14] Laurens Van Der Maaten and Geoffrey Hinton. Visualizing Data using t-SNE. *Journal of Machine
134 Learning Research*, 9:2579–2605, 2008.