



A Formal Methods Approach Towards Deep Learning Interpretability

Kriten Kessel, Christopher Lazarus, Javier Sagastuy
Stanford University



Summary

Deep Neural Networks are very useful at classifying tasks but their intrinsic complexity makes it really hard to explain the reasoning behind a classification outcome. In recent work [2], statistical methods were introduced to help assess the influence of human intelligible concepts in classification outcomes. We aim to asses and extend such methods by leveraging formal methods for Neural Network verification.

- **Problem:** Why did the network classify image γ with label k ?
- **Solution:** Come up with classes TCAV etc bla bla? or more like we tried to test the TCAV method??
- **Preliminary Results:** nothing nothing nothing.

Talk about TCAV ?

Concept Activation Vectors:

asdasdasd Output:

- predictions for the unlabelled nodes \hat{y}_L

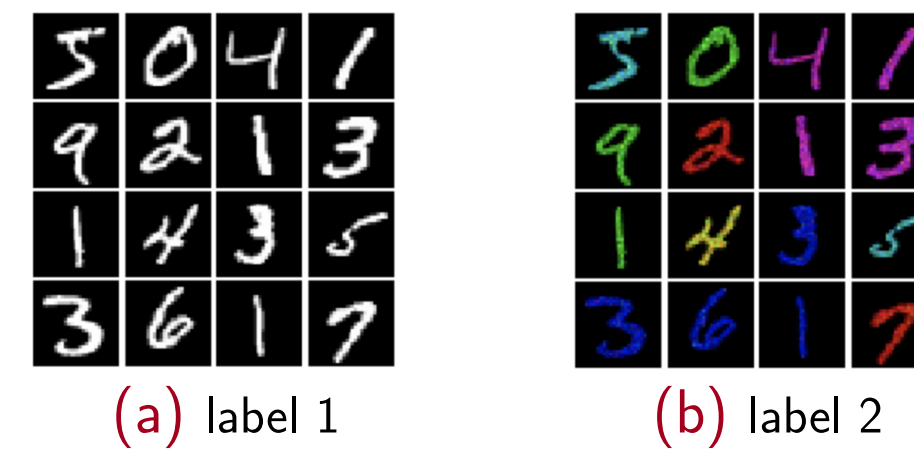
Neural Network Verification

GCN [1] for node classification: $\hat{y} = h(Ah(AXW_1)W_2)$

Approach: TCAV + Verification

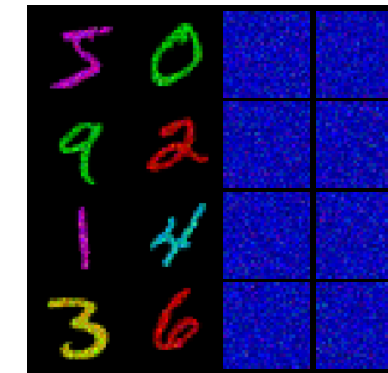
Data bla bla

- talk about data

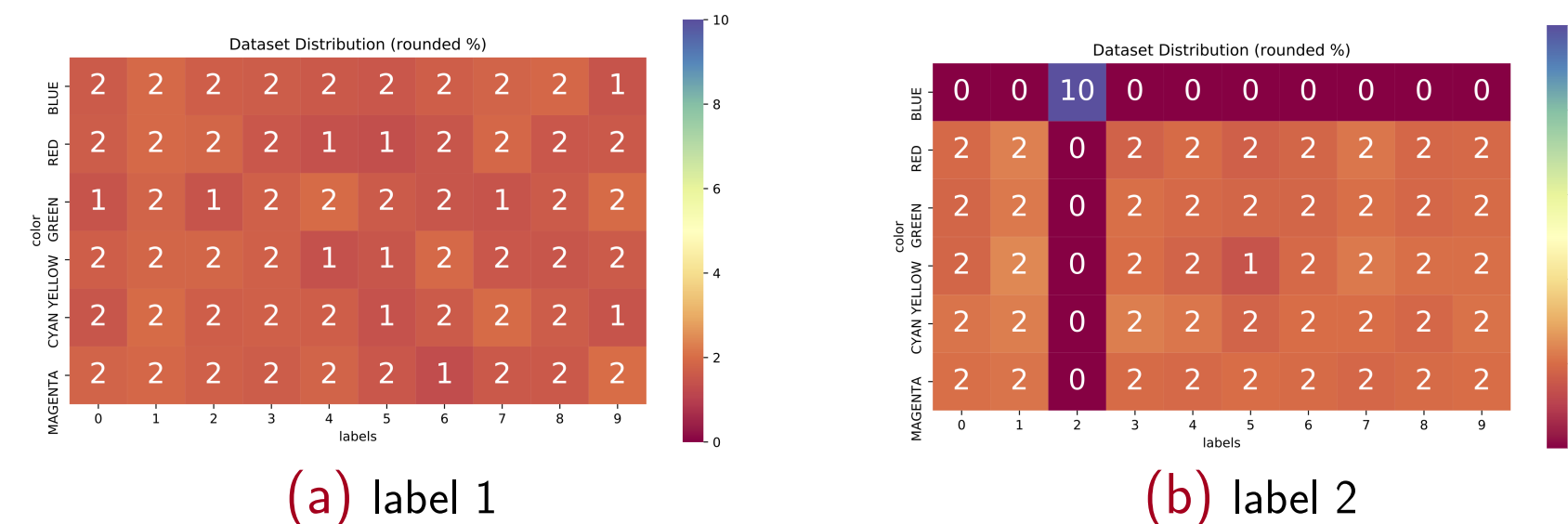


(a) label 1

(b) label 2



Maybe talk about classe sand support vector LALALALALAL



(a) label 1

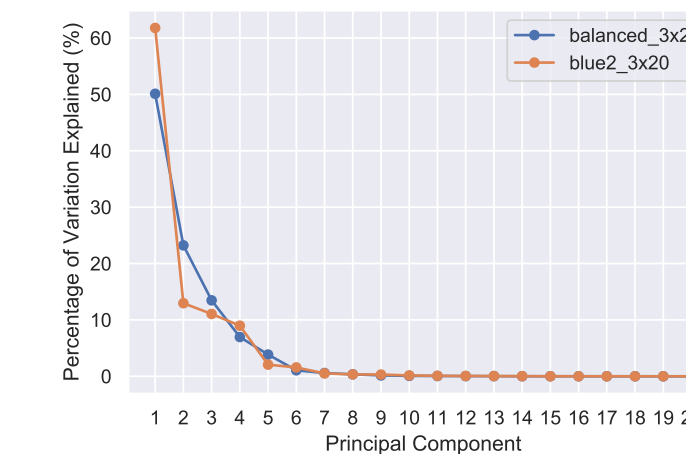
(b) label 2

$$\mathcal{L}_{LP} = -\mathbb{E}_{Z \sim q(Z|A, X)}[A_{ij} \log \tilde{A}_{ij} + (1 - A_{ij}) \log(1 - \tilde{A}_{i,j})] + \text{KL}(q(Z|A, X) || p(Z)).$$

2. **Edge hallucination** produces \hat{A} :
 - top K (K hyper-parameter)
 - sampling using gumbel softmax [4] trick (allows gradients to flow)
3. **Node classification**
 - $\hat{y} = \text{GCN}(\hat{A}, X)$

Results

Need to talk about significant CAVs Then talk about the avenues and boulevards.



A figure

model	layer	TCAV Score	significant
balanced.5x50	fc1	0.15 ± 0.10	yes
	fc4	0.14 ± 0.13	yes
blue2.5x50	fc1	1.00 ± 0.00	yes
	fc4	1.00 ± 0.00	yes
balanced.3x50	fc1	0.14 ± 0.13	yes
	fc2	0.08 ± 0.05	yes
blue2.3x50	fc1	0.80 ± 0.08	yes
	fc2	0.78 ± 0.11	yes
balanced.3x20	fc1	0.20 ± 0.06	yes
	fc2	0.14 ± 0.05	yes
blue2.3x20	fc1	1.00 ± 0.01	yes
	fc2	0.98 ± 0.01	yes

A table

A table

Node classification results

- Try to salvage somethibg
- Nothing worked :)

References

- [1] G. Katz, C. Barrett, D. L. Dill, K. Julian, and M. J. Kochenderfer. Reluplex: An efficient smt solver for verifying deep neural networks. In International Conference on Computer Aided Verification, pages 97–117. Springer, 2017
- [2] B. Kim, M. Wattenberg, J. Gilmer, C. J. Cai, J. Wexler, F. B. Viegas, and R. Sayres. Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (tcav). In ICML, 2018
- [3] Y. LeCun, C. Cortes, and C. Burges. Mnist handwritten digit database. ATT Labs [Online]. Available: <http://yann.lecun.com/exdb/mnist>, 2:18, 2010
- [4] C. Liu, T. Arnon, C. Lazarus, C. Barrett, and M. J. Kochenderfer. Algorithms for verifying deep neural networks. CoRR, abs/1903.06758, 2019