

Improving the Interpretability of DEMUD on Image Data Sets

Jake Lee, Jet Propulsion Laboratory, California Institute of Technology & Columbia University, CS '19 Intern under Kiri Wagstaff Summer 2018

Government sponsorship acknowledged. This work was performed at the Jet Propulsion Laboratory, California Institute of Technology, under a contract with NASA. CL#18-4779

Motivation



- Ever-increasing of volume of image data in all fields
- Analysis of imagery is time-consuming and labor-intensive
- Lots of existing work on supervised image classification
 - Wagstaff, Kiri L., et al. "Deep Mars: CNN Classification of Mars Imagery for the PDS Imaging Atlas." Conference on Innovative Applications of Artificial Intelligence. 2018.
- However, scientific discovery relies on unexpected observations

A brief introduction to DEMUD

- A prior-free novelty detection algorithm
- Prioritizes interesting data by attempting to discover all existing classes as quickly as possible
- Provides <u>explanations</u> for its prioritizations

DEMUD + images

- Ongoing work since Summer 2017
- Presented at 2018 ICML Workshop on Human Interpretability in Machine Learning



Example output ("Yellow" ImageNet dataset, fc6)





Example output (MSL Image dataset, fc6)





Kiri L. Wagstaff, You Lu, Alice Stanboli, Kevin Grimes, Thamme Gowda, and Jordan Padams. "Deep Mars: CNN Classification of Mars Imagery for the PDS Imaging Atlas." Proceedings of the Thirtieth Annual Conference on Innovative Applications of Artificial Intelligence, 2018. 10.5281/zenodo.1049137

Improve the Visualizations!



A Better Visualization Method

• Dosovitskiy & Brox 2016 **CVPR**





Dosovitskiy, Alexey, and Thomas Brox. "Generating images with perceptual similarity metrics based on deep networks." Advances in Neural Information Processing Systems. 2016.

Dosovitskiy, Alexey, and Thomas Brox. "Inverting visual representations with convolutional networks." Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2016.



Evaluating Visualization Methods

- These methods were never intended to visualize modified feature vectors
- Specifically, DEMUD performs a subtraction in feature space for its explanations
 - Selection Expected = Novel
- Unclear whether visualizing modified features can be meaningful

Simple Image Arithmetic



Simple Image Arithmetic



D&B L2 Results



D&B ADV Results



Plotting the feature distribution (fc6, 4096 values)



Mean-shift normalization





D&B ADV Results (with mean-shift normalization)





















Deng, Jia, et al. "Imagenet: A large-scale hierarchical image database." Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on. leee, 2009.











Kiri L. Wagstaff, You Lu, Alice Stanboli, Kevin Grimes, Thamme Gowda, and Jordan Padams. "Deep Mars: CNN Classification of Mars Imagery for the PDS Imaging Atlas." Proceedings of the Thirtieth Annual Conference on Innovative Applications of Artificial Intelligence, 2018. 10.5281/zenodo.1049137





Kiri L. Wagstaff, You Lu, Alice Stanboli, Kevin Grimes, Thamme Gowda, and Jordan Padams. "Deep Mars: CNN Classification of Mars Imagery for the PDS Imaging Atlas." Proceedings of the Thirtieth Annual Conference on Innovative Applications of Artificial Intelligence, 2018. 10.5281/zenodo.1049137

Visualization sensitivity analysis



Ongoing work

- More investigation into feature-level interactions and operations
- Visualizations for fc6, fc7, fc8
- More Experiments with DEMUD
- User Study

Acknowledgements

- Kiri Wagstaff
- Fellow summer interns
- Alexey Dosovitskiy
- PDS Imaging Node