

VISUAL REPRESENTATION LEARNING DOES NOT GENERALIZE STRONGLY WITHIN THE SAME DOMAIN

Lukas Schott¹, Julius von Kügelgen^{2, 3, 4}, Frederik Träuble^{2, 4},
Peter Gehler⁴, Chris Russell⁴, Matthias Bethge^{1, 4}, Bernhard Schölkopf^{2, 4},
Francesco Locatello^{4, †}, Wieland Brendel^{1, †}

¹University of Tübingen

²Max Planck Institute for Intelligent Systems

³University of Cambridge

⁴Amazon Web Services

[†]Joint senior authors

lukas.schott@bethgelab.org

ABSTRACT

An important goal in machine learning is to uncover the underlying latent factors of variation as well as the true mechanism through which each factor acts in the world.

In this paper, we test whether 17 representation learning approaches correctly infer the generative factors of variation in simple datasets (DSprites, Shapes3D, MPI3D). In contrast to prior work that introduces novel factors of variation during test time, such as blur or other (un)structured noise, we here recombine, interpolate, or extrapolate only existing factors of variation from the training data set (e.g. small and medium sized squares during training and large squares during testing). Models that learn the correct mechanism should be able to generalise to this benchmark.

In total, we train and test 2000+ unsupervised, weakly-supervised and fully-supervised models and observe that all of them struggle to learn the underlying mechanism regardless of supervision signal and architectural bias. Moreover, the generalization capabilities of all tested models drop significantly as we move from artificial data sets towards more realistic real world data sets. Despite their inability to identify the correct mechanism, the models are surprisingly modular as their ability to infer other in-distribution factors remains fairly stable providing only a single factor is OOD.

These results point to an important yet understudied problem of learning mechanistic models of observations that can facilitate generalization.

