

CS6360: Advanced Topics in Machine Learning

Group Theory \cap Machine Learning

Pranav K Nayak

IIT Hyderabad

Paper Presentation 1

Homomorphism Autoencoder: Learning Group Structured Representations from Observed Transitions: Hamza Keurti, Hsiao-Ru Pan, Michael Besserve, Benjamin Grewe, Bernhard Schölkopf, *ICML 2023*

The paper attempts to model the effect of interventions as transformations in representation space.

They assert that this problem can be formulated as a problem of learning a homomorphism between the interventional structure of the world and the model's representations of it. This should allow it to be able to reverse-engineer the effects of potential interventions (transformations) through the knowledge of how its representations change.

The Learning Problem

- W is the latent space from which observations are generated, through the process g .
- O is the space of observations.
- Z is the space of representations, mapped to from O through the *inference process* h .

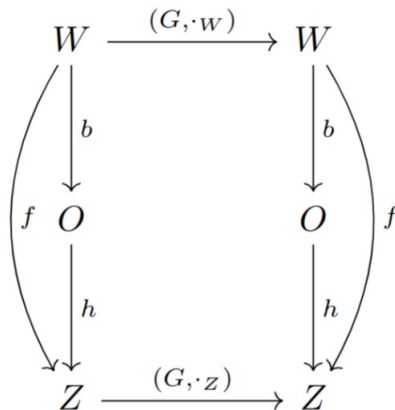


Figure: The proposed group structure of the learning problem

HAE and Two Losses

Definition (N -step Prediction Loss)

$$\mathcal{L}_{\text{pred}}^N(\rho, h) = \sum_t \sum_{j=1}^N \|h(o_{t+j}) - (\prod_{i=0}^{j-1} \rho(g_{t+i}))h(o_t)\|$$

Definition (N -step Reconstruction Loss)

$$\mathcal{L}_{\text{rec}}^N(\rho, h, d) = \sum_t \sum_{j=1}^N \|o_{t+j} - d(\prod_{i=0}^{j-1} \rho(g_{t+i}))h(o_t)\|$$

A weighted sum of both losses, $\mathcal{L}(\rho, h, d) = \mathcal{L}_{\text{rec}}^N(\rho, h, d) + \gamma \mathcal{L}_{\text{pred}}^N(\rho, h)$, is optimized for.

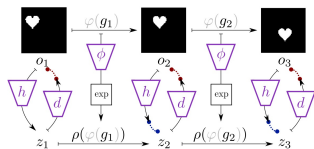


Figure: HAE's Architecture

Restrictions on the Representation Class

Theorem

If (ρ, h, d) are continuous and minimize the expectation of $\mathcal{L}_{pred}^2(\rho, h) + \gamma \mathcal{L}_{rec}^k(\rho, h, d)$, for $k \geq 0$, then ρ is a non-trivial group representation and (ρ, h) is a symmetry-based representation.

Informally, a symmetry-based representation is one that satisfies the following:

$$\rho(g_1, g_2, \dots, g_n)(z_1 \oplus \dots \oplus z_n) = \rho_1(g_1)(z_1) \oplus \dots \oplus \rho_n(g_n)(z_n) \text{ where } z_i = h(o_i).$$