

## Identifiable Steering via Sparse Autoencoding of Multi-Concept Shifts

Shruti Joshi, Andrea Dittadi, Sébastien Lachapelle, Dhanya Sridhar

**Motivation and Setup** 



## **Informal Theorem Statement**

SSAEs learn the concept shift vectors  $\delta_{V}^{c}$  and the linear encoder-decoder pair (r, q) up to permutation and scaling of the true solution, where columns of the decoder denote steering vectors for the concepts in  $c_{V}$ .

Unsupervised, but at least theoretically related to

Cosine Similarity  $(\tilde{\mathbf{z}}, \circ)$  (Higher is better)

true concepts via trivial transformations

BINARY(2, 2)

 $\tilde{\mathbf{z}}_{\mathrm{MD}}$ 

Causal Representation Learning  $\longleftrightarrow$  Mechanistic interpretability; propose a method that disentangles the latent concepts encoded in LLM activations

**Linear Representation Hypothesis** backbone treats activations as a linear mix of concepts

Provable disentanglement with weak supervision
— uses multi-concept contrast pairs, far cheaper
than single-concept perturbations



Exhibits strong **OOD** performance

Separates strongly **correlated** concepts

## Mean Correlation Coefficients (MCC) between the decoders of any two learned models

	SSAE	aff
LANG(1,1)	$0.995 \pm 0.001$	$0.985 \pm 0.004$
GENDER(1, 1)	$0.993 \pm 0.000$	$0.961\pm0.000$
BINARY(2,2)	$0.991 \pm 0.001$	$0.936 \pm 0.000$
$\operatorname{CORR}(2,1)$	$0.991 \pm 0.001$	$0.928 \pm 0.077$
CAT(135, 3)	$0.906 \pm 0.022$	$0.661\pm0.019$
TruthfulQA	$0.952\pm0.006$	$0.885\pm0.006$
	SSAE	aff
SYNTH(3, 2)	$0.999\pm0.0001$	$0.873\pm0.0561$
synth(4,3)	$0.999\pm0.0011$	$0.835 \pm 0.0097$
SYNTH(10, 7)	$0.993\pm0.0005$	$0.769\pm0.0103$

**Reproducible** across hyper-parameters

Learns steering from data with *multiple unknown concept* variations



LANG(1, 1)

 $\begin{array}{c|c} \mathbf{Steering} \ \mathbf{Method}(\circ) \\ \hline & & \\ \hline & & \\ \hline & & \\ \mathbf{\tilde{z}_{aff}} \\ \hline & & \\ \hline \end{array}$ 

 $\tilde{\mathbf{z}}_{PCA}$ 

GENDER(1, 1)

0.9

0.8

0.7

0.6

0.5

 $\mathbf{0.4}$ 

0.3

 $0.2^{-1}$ 

0.1

0.0



CORR(2, 1)

## SAMSUNG

TruthfulQA(1, 1)

CAT(135, 3)

Advanced Institute of Technology Al Lab Montreal



