Linear Representation Hypothesis & Geometry of LLMs Kiho Park, Yo Joong Choe, Victor Veitch University of Chicago

What Does "Linear" Even Mean?



Problem: It is not clear how these ideas relate to each other, nor which is the right notion of linear representation.



A concept W is defined by counterfactual outputs Y(W = 0), Y(W = 1). Concepts W and Z are causally separable if Y(w, z) is well-defined, that is, they can be manipulated freely and in isolation.

Formalizing Linear Representation Hypothesis

We first formalize the subspace notions of linear representations, then use softmax structure to connect them to measurement and intervention.





Intervention Representation $\mathbb{P}(W = 1 \mid Z, \lambda + c\bar{\lambda}_W) \text{ increasing in } c \in \mathbb{R}$ $\mathbb{P}(Z = 1 \mid W, \lambda + c\bar{\lambda}_W)$ constant in $c \in \mathbb{R}$

Causal Inner Product

Problems:

How do the unembedding and embedding representations relate? What is the right inner product for the representation space?



We introduce the causal inner product $\langle \cdot, \cdot \rangle_{C}$ such that $\langle \bar{\gamma}_W, \bar{\gamma}_Z \rangle_C = 0$ whenever W and Z are causally separable.

Theorem

This unifies the unembedding and embedding representations via $\langle \bar{\gamma}_W, \cdot \rangle_C = (\bar{\lambda}_W)^\top$. (This is the Riesz isomorphism.) We can estimate the causal inner product as $\langle \bar{\gamma}, \bar{\gamma}' \rangle_C = \bar{\gamma}^T \text{Cov}(\gamma)^{-1} \bar{\gamma}'$.

Experiments with LLaMA-2

We estimate the unembedding representations for various concepts by using the counterfactual pairs from a word analogy dataset.

$$\bar{\gamma}_W := \frac{\tilde{\gamma}_W}{\langle \tilde{\gamma}_W, \tilde{\gamma}_W \rangle_C} \text{ where } \tilde{\gamma}_W$$

Linear Representations Exist



Differences between counterfactual pairs are more parallel to $\bar{\gamma}_W$ than those between random pairs, supporting the linear representation hypothesis.

Unembedding $\gamma(y) \in \mathbb{R}^d$

Next word

''king'' Y(0,0)queen" Y(1,0) "roi" Y(0,1)"reine" Y(1,1)

 $verb \Rightarrow Ving (2)$ $verb \Rightarrow Ved (3)$ $Ving \Rightarrow 3pSg (4)$ $Ving \Rightarrow Ved (5)$ $3pSg \Rightarrow Ved (6)$ $verb \Rightarrow V + able (7)$ $verb \Rightarrow V + er$ $verb \Rightarrow V + ment$ (2) $adj \Rightarrow un + adj$ (1) adj⇒adj+ly(12 $male \Rightarrow female$



Heatmaps of $\langle \bar{\gamma}_W, \bar{\gamma}_{W'} \rangle$ show that the causal inner product (left) between the unembedding representations of causally separable concepts is close to zero. It clearly improves on the naive Euclidean inner product (right).

Unembedding Representation Yields Linear Probe $W = French \Rightarrow Spanish$ $W = male \Rightarrow female$ French Spanish

 $\bar{\gamma}$ French \Rightarrow Spanish (left) separates the embeddings of French and Spanish contexts, while $\bar{\gamma}_{male \Rightarrow female}$ (right) does not.

Embedding Representation Yields Steering Vector



Adding the embedding representation $\bar{\lambda}_W := \text{Cov}^{-1}(\gamma)\bar{\gamma}_W$

to context embeddings changes the target concept,

without changing other causally separable concepts.

Find us online! Email: parkkiho@uchicago.edu



Kiho Park

 $\gamma(y_i(1)) - \gamma(y_i(0))$

Causally Separable Concepts are Represented Orthogonally under the Causal Inner Product



10	$W = French \Rightarrow Spanish$					
F						
С						
0						
0		•				
		• •	•••			
-5						
		10	-5	0	5	10

 $\lambda = \lambda$ ("Long live the") $\lambda = \lambda$ ("Long live the") + $\alpha \overline{\lambda}_{male}$ female



YJ Choe



Victor Veitch



arXiv:2311.03658