



Energy-Regularized Sequential Model Editing on Hyperspheres

Qingyuan Liu^{1*}, Jia-Chen Gu^{2*}, Yunzhi Yao³, Hong Wang⁴, Nanyun Peng²

¹Columbia ²UCLA ³ZJU ⁴USTC

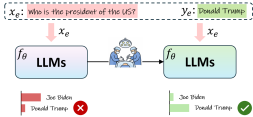


Paper

Code

Background & Motivation

Model Editing: Localized parameter updates modify target knowledge while **avoiding interference** with unrelated knowledge



Hyperspherical Uniformity

Hyperspherical Energy (HE) is used to quantify the uniformity of neuron weights on hyperspheres:

$$HE(X) = \sum_{i \neq j} (\|\hat{x}_i - \hat{x}_j\|^2 + \epsilon)^{-\alpha} = \sum_{i \neq j} (2(1 - \cos \theta_{ij}) + \epsilon)^{-\alpha}$$

Lower Energy

Higher Orthogonality

High-HE Weight

Low-HE Weight



Correlation between Uniformity and Editing



Energy increases

Weights become less uniform

Statistically significant negative correlation

Maintaining uniformity is critical in stabilizing editing

SPHERE

Sparse Projection for Hyperspherical Energy Regularized Editing

Project editing knowledge onto a **sparse** space which is **complementary to principal** hyperspherical directions



Enable new know incorporation



Mitigate ▲'s interfere with ●

● original weight
▲ update weight

Principal Space Estimation

$$v = \arg \max_{\|v\|=1} \left(\frac{1}{n} \|Wv\|^2 \right)$$

Maximize the variance of all neurons in W when projected onto v

Top-r: $U = [v_{d-r+1}, \dots, v_d] \in \mathbb{R}^{d \times r}$

Sparse Space Definition

$$P_\perp = I - \alpha U U^T \in \mathbb{R}^{d \times d}$$

Orthogonal

Suppression Strength:

how much contribution of U is removed

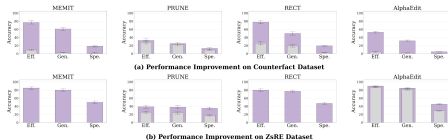
Sparse Space Projection

$$\hat{W} = W + \Delta W P_\perp$$

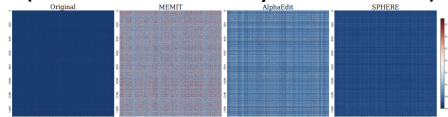
could be any editing method

Results & Analysis

Baseline editing methods can be improved with plug-and-play SPHERE by 43.27% (Eff.), 36.20% (Gen.), and 20.96% (Spe.) on avg.



SPHERE's angular distribution remains stable (elements: cosine similarity between neurons)



SPHERE-edited models can better preserve the general abilities

