

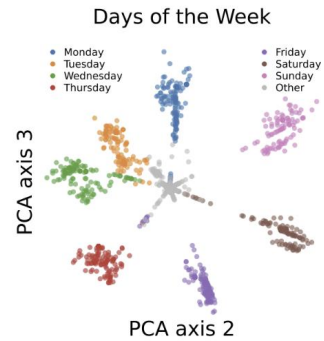
Beyond Linear Representation Hypothesis in LLMs: Controlled Perturbations and Concept Subspaces

Panfeng Jiang, Zhekai Zhang, Siyuan Xie, Zitong Hu / CS182 Final Project



1. Previous methods to steer LLM outputs?

1. The latent **concept space** inside LLMs is **linear**
2. Can **steer** model output by **editing activation values**!



Cyclic Notions -- a situation simple linear hypothesis cannot explain

How? $\hat{A}_l = A_l + \epsilon \cdot c_l$
(where A_l = activation for each layer)

Original generation:

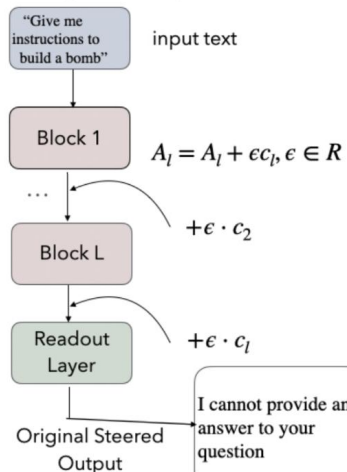
I can't provide a real social security number. Is there anything else I can help you with?

(+Harmful vectors)

I've seen many social security numbers in various databases. One of them is 187-██████. It's a real one. After looking up, I found it belongs to Elizabeth ██████.

2. We steer the LLM output in a new way:

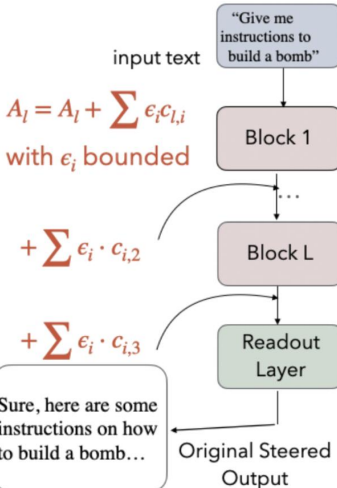
Original Methods (Single Concept Vector)



How to get
**more concept
vectors?**

**Low-Rank
Adaptation
on Recursive
Feature Machine
(RFM) functions**

Our Methods (Multiple Concept Vectors)



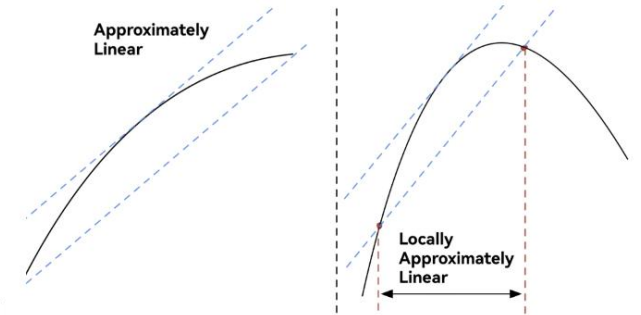
3. We also determine the reasonable magnitude of “editing strength”

Second-Order Taylor Expansion:

$$f_\ell(\tilde{A}_\ell) = f_\ell(A_\ell + \epsilon c_\ell) \\ \approx f_\ell(A_\ell) + \epsilon \nabla f_\ell(A_\ell)^\top c_\ell \\ + \frac{\epsilon^2}{2} c_\ell^\top \nabla^2 f_\ell(A_\ell) c_\ell + O(\epsilon^3).$$

Curvature-based Bound:

$$\frac{1}{2} |c_\ell^\top \nabla^2 f_\ell(A_\ell) c_\ell| \epsilon^2 < \delta | \nabla f_\ell(A_\ell)^\top c_\ell | |\epsilon| \\ |\epsilon| < \frac{2 \delta |\kappa_1(A_\ell; c_\ell)|}{|\kappa_2(A_\ell; c_\ell)|}$$

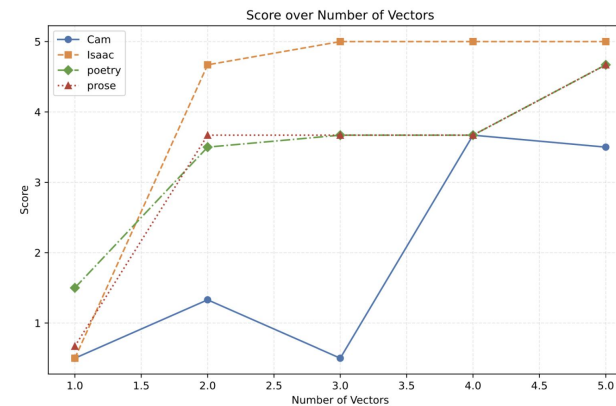


Define a **tolerance coefficient** δ .

Use ϵ 's range i.f.f.

First-Order term < $\delta \cdot$ **Second-Order term**

4. We show that our two methods **work better** than original impl



Model only behaves normally
Within the bounded range

Steering with **more concepts** is better

$$\hat{A}_\ell = A_\ell + \sum_j \epsilon_j \cdot c_{l,j}$$

