

Beyond Sparsity: Improving Diversity in Sparse Autoencoders via Denoising Training

Xiang Pan¹ Yifei Wang² Qi Lei¹ ¹New York University ²Massachusetts Institute of Technology xiangpan@nyu.edu yifei_w@mit.edu ql518@nyu.edu

Problem: Feature Redundancy in SAEs



Deduplication Results



Methodology Overview

Standard SAE Training:

Input: LLM layer activations x
Encoder: z = TopK(ReLU(W₁x))
Decoder: x̂ = W₂z

- Current SAEs suffer from high feature redundancy
- Only 85% of features are unique in standard Top-K SAEs (1)
- Redundant features take up feature space capacity, preventing capture of important semantic directions
- High interpretation scores on redundant features can bias evaluation metrics

Example of a redundant feature: "A neuron about the [Beginning Of the Sentence]" — this feature reflects a preprocessing artifact, not a meaningful or interpretable concept.

Evaluation Metrics

1. Feature Coverage Score

 $1 \stackrel{d'}{-} (\mathbf{N} \mathbf{I}'^{\mathsf{T}} \mathbf{N} \mathbf{I}_{\mathsf{a}})_{\mathsf{m}}$

► Loss:
$$\mathcal{L} = \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 + \alpha \cdot \text{auxiliary}$$

Our Denoising Training:

- 1. Apply dropout as input noising mechanism
- 2. Reconstruct original input from the perturbed input

$$\mathcal{L} = \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 + lpha \cdot ext{auxiliary} + eta \|\mathbf{x} - \hat{\mathbf{z}}\|_2^2$$

Experimental Results





$$\frac{1}{d'} \sum_{i=1}^{I} \mathbb{1} \left(\max_{j} \frac{(\mathbf{v}_{2} \ \mathbf{v}_{2})_{ij}}{\|\mathbf{W}_{2}'[:,i]\|_{2} \|\mathbf{W}_{2}[:,j]\|_{2}} > \tau \right)$$

Measures fraction of feature bank directions covered by current SAE features.

2. Semantic Volume

Semantic Volume = log det($\mathbf{\tilde{V}}^T \mathbf{\tilde{V}} + \epsilon \mathbf{I}$) (2)

Quantifies diversity of feature explanation space using semantic embeddings.

3. Effective Rank

Effective Rank = exp
$$\left(-\sum_{i=1}^{r} \tilde{\sigma}_{i} \log \tilde{\sigma}_{i}\right)$$
 (3)

Measures how evenly the feature space is distributed across semantic dimensions.

Key Findings:

(1)

- Improved Feature Diversity: Higher feature coverage, semantic volume, and effective rank.
- Maintained Interpretation Quality: No degradation in interpretation scores.
- Better Reconstruction Quality

(4)