

Lightweighted Sparse Autoencoder based on Explainable Contribution

JooHong Rhee and Hyunggon Park

Graduate Program in Smart Factory, Department of Electronic and Electrical Engineering
Ewha Womans University, Seoul, Republic of Korea

Introduction

Motivations & Challenges

- Performance improvement leads to heavy autoencoders
- Heavy autoencoders require high computing powers
→ Cannot be implemented into power-limited devices

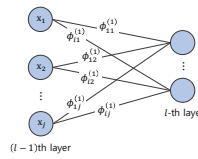
Objective: **Design lightweight autoencoder** while maintaining performance

Proposed Algorithm : SHAP-SAE

Link importance (LI) based on Shapley value

$$\phi_{ij}^{(l)} = \sum_{J \in \mathcal{A}(l)} \frac{|J|!(|l|-|J|-1)!}{|l|!} (v(J \cup \{i\}) - v(J))$$

- l : set of links in l -th layer
- J : subset excluding j -th link connected to i -th unit in l -th layer
- $v(J)$: value of subset J



- LI of link that connects j -th unit in $(l-1)$ th layer and i -th unit in l -th layer
- Measurement of LI based on their **contributions to output of layer**

Unit importance (UI)

$$\bar{v}_j^{(l-1)} = \frac{1}{n^{(l)}} \sum_{i=1}^{n^{(l)}} |\phi_{ij}^{(l)}|$$

- UI of j -th unit in $(l-1)$ th layer
- Average impact on output of l -th layer

- $n^{(l)}$: number units in l -th layer

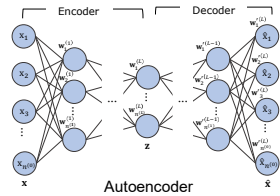
Step 1: Training Stage

1. Determine the set of optimal parameters for autoencoder

$$\{\theta^*, \theta'^*\} = \underset{\theta, \theta'}{\operatorname{argmin}} \mathcal{L}(x, \hat{x})$$

- Encoder: $\theta = \{W^{(l)}, b^{(l)} | 1 \leq l \leq L\}$
- Decoder: $\theta' = \{W'^{(l)}, b'^{(l)} | L+1 \leq l \leq 2L\}$

- Weight matrix : $\{W^{(l)}, W'^{(l)}\} \in \mathbb{R}^{n^{(l)} \times n^{(l-1)}}$
- Bias vector : $\{b^{(l)}, b'^{(l)}\} \in \mathbb{R}^{n^{(l)}}$
- $2L$: number of layers in autoencoder



Step 2: Sparsification Stage

1. Total LI $\phi_T^{(l)}$ as sum of individual LIs in l -th layer

2. Set of descending ordered Shapley values in l -th layer

$$\Phi^{(l)} = [\Phi^{(l)}(1), \Phi^{(l)}(2), \dots, \Phi^{(l)}(n^{(l-1)}n^{(l)})]$$

3. Support set $\Gamma^{(l)}$

$$\Gamma^{(l)} = \left\{ (i, j) \mid \sum_{k=1}^{k^*} \Phi^{(l)}(k) \geq m \cdot \phi_T^{(l)} \right\}$$

- m : importance level ($0 < m \leq 1$)
- k, k^* : integer

- Set of the pairs (i, j) of units i and j that have k^* largest contribution

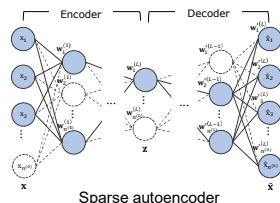
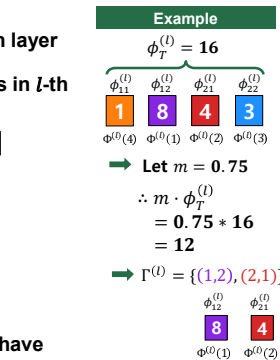
4. Prune autoencoder using mask function \mathcal{M}

- Mask function \mathcal{M} : $W^* = \mathcal{M}(W)$ and $W'^* = \mathcal{M}(W')$

$$\begin{cases} w_{ij}^{*(l)} = 0, & \text{if } (i, j) \in \Gamma^{(l)c} \text{ for } 1 \leq l \leq L \\ w_{ij}^{*(l)} = 0, & \text{if } (i, j) \in \Gamma^{(l)c} \text{ for } L+1 \leq l \leq 2L \end{cases}$$

- $\Gamma^{(l)c}$: complementary set of $\Gamma^{(l)}$
- W^* : pruned weight matrix of encoder
- W'^* : pruned weight matrix of decoder
- $w_{ij}^{*(l)}$: element of weight matrix W^*
- $w_{ij}^{*(l)}$: element of weight matrix W'^*

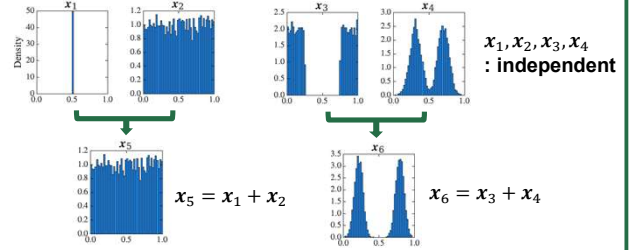
- The rest of weights remain unchanged



Experiment Results

SHAP-SAE with Synthetic dataset

Dataset



Explainability

- Impact of UIs on latent vector z

- $\bar{v}_5^{(0)} > \bar{v}_1^{(0)}$ or $\bar{v}_2^{(0)}$
- $\bar{v}_6^{(0)} > \bar{v}_3^{(0)}$ or $\bar{v}_4^{(0)}$

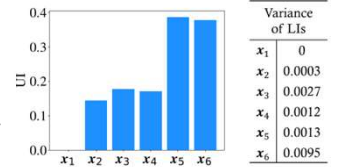
- ▶ x_5, x_6 : Include information of x_1, x_2, x_3, x_4

- Contribution of x_1, x_2, x_3, x_4 are marginal

- ▶ Redundant to x_5, x_6

- $\bar{v}_1^{(0)} = 0$

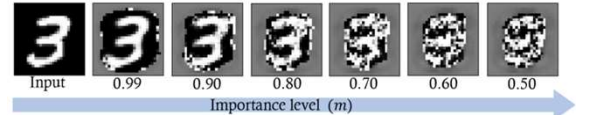
- ▶ x_1 : Set of constant values → Irrelevant to z



SHAP-SAE with MNIST dataset

Performance Analysis

- Gray pixels: pruned weights



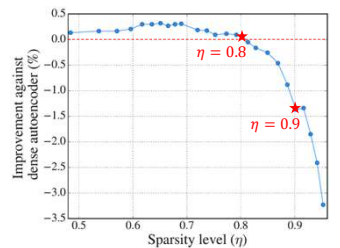
- As importance level m increases, mask function \mathcal{M} **removes less important weights**

- Outperforms dense autoencoder up to $\eta = 0.8$

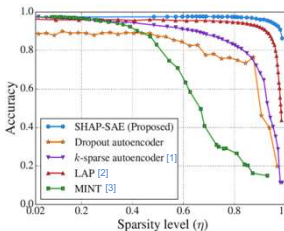
- Only 1.37% performance degradation with $\eta = 0.9$

- Sparsity level η : how many links are deactivated among all the links in autoencoder

$$\eta = \frac{\sum_{l=1}^{2L} |r^{(l)c}|}{\sum_{l=1}^{2L} |r^{(l)} \cup r^{(l)c}|}$$



Performance Comparisons



- **Outperforms other benchmarks** across all sparsity level

- **Remarkably robust** against high sparsity level

Conclusions

SHAP-SAE algorithm can

- Explicitly measure unit and link importance based on Shapley value
→ Activate only important units and links
- Allow sparse autoencoder to be explainable and robust against high sparsity level
- Be implemented into power-limited devices

[1] Makhzani, A. and Frey, B., "K-sparse autoencoders. *International Conference on Learning Representations*," 2014.

[2] Park, S., Lee, J., Mo, S., and Shin, J., "Lookahead: A farsighted alternative of magnitude-based pruning." *International Conference on Learning Representations*, 2020.

[3] Ganesh, M. R., Corso, J. J., and Sekhe, S. Y., "Mint: Deep network compression via mutual information-based neuron trimming." *2020 25th International Conference on Pattern Recognition (ICPR)*, pp. 8251-8258. IEEE, 2021.