# Finding better Sparse Neural Networks is hard



- Sparse neural networks are neural networks with a lot of zero-valued weight parameters.
- They can use up less memory and computation than dense ones, making them increasingly sought after as modern models grow to extreme sizes.
- $\min f(x)$
- Various neural network pruning techniques were developed to find sparse neural networks x with good performance or minimal loss f.
  - Still, it remains a constant challenge to push for higher sparsity with minimal performance degradation.

# **Degraded performance by loss sharpness**



- Sharpness of loss has been suggested as a major cause for diminishing the trainability of sparse networks in high sparsity (Lee et al., 2021).
- More generally, its strong correlations with degradation in generalization performance has been well studied (Keskar et al., 2017).



- Various techniques were proposed to explicitly minimize sharpness during training (Foret et al., 2021).
- This has been shown to be effective in improving the generalization performance and robustness of neural networks.

# **Enforcing flatness to improve sparsification**



Despite studies connecting degraded performance by sparsification to loss sharpness, very little work has attempted to develop a general framework to explicitly reduce sharpness while sparsifying neural networks.



To fill this gap, we propose SAFE, a general constrained-optimization-based framework to simultaneously enforce flatness and sparsification while optimizing the neural network.

# **SAFE:** Sparsification via ADMM with Flatness Enforcement

#### Sharpness-aware sparsity-constrained optimization problem

We first formulate this as an optimization problem, where the goal is to find a sparse solution  $x^*$  with at most d non-zero elements that minimizes the objective min max  $f(x+\epsilon)$  $\|x\|_0 \leq d \|\epsilon\|_2 \leq \rho$ function in the whole  $\epsilon$ -neighborhood, *i.e.*, seek flat minima.

#### Augmented Lagrangian based approach

To solve this, we form the augmented Lagrangian dual problem of the following:

$$\max_{u}, \min_{x,z} \left[ \mathcal{L}(x, z, u) := \max_{\|\epsilon\|_2 \le \rho} f(x + \epsilon) + I_{\|\cdot\|_0 \le d}(z) - \frac{\lambda}{2} \|u\|_2^2 + \frac{\lambda}{2} \|x - z + u\|_2^2 \right],$$

where we separate the sparsity-constraint satisfaction using variable z so that it can be handled more

#### **Alternating Direction Method of Multipliers**

$$\begin{aligned} x_{k+1} &= \arg\min_{x} \max_{\|\epsilon\|_{2} \le \rho} f(x+\epsilon) + \frac{\lambda}{2} \|x-z_{k}+u_{k}\|_{2}^{2} \quad (\text{x-min}) \\ z_{k+1} &= \arg\min_{z} I_{\|\cdot\|_{0} \le d}(z) + \frac{\lambda}{2} \|x-z+u\|_{2}^{2} \quad (\text{z-min}) \\ u_{k+1} &= u_{k} + x_{k+1} - z_{k+1} \end{aligned}$$

We apply dual ascent and minimize x and zin an alternating fashion, which gives us this ADMM iterate.

#### x-minimization: iterative minimization while enforcing flatness

We solve this iteratively using *Sharpness-aware minimization (SAM)*, where we approximately solve for  $\epsilon$  through first-order Taylor approximation:

$$\epsilon^{\star}(x) \approx \underset{\|\epsilon\|_{2} \leq \rho}{\operatorname{argmax}} f(x) + \epsilon^{\top} \nabla f(x) = \rho \frac{\nabla f(x)}{\|\nabla f(x)\|_{2}}$$

Applying this back to the objective and applying gradient descent gives us the following iteration for x-minimization

$$(\text{x-min}) \quad x_k^{(t+1)} = x_k^{(t)} - \eta^{(t)} \left[ \nabla f \left( x_k^{(t)} + \rho \frac{\nabla f(x_k^{(t)})}{\|\nabla f(x_k^{(t)})\|_2} \right) + \lambda (x_k^{(t)} - z_k + u_k) \right]$$

# ICML 2025 Spotlight Paper

# SAFE: Finding Sparse and Flat Minima to Improve Pruning

Dongyeop Lee, Kwanhee Lee, Jinseok Chung, and Namhoon Lee Pohang University of Science and Technology (POSTECH)

: We propose SAFE and SAFE+, constrained optimization algorithms to enforce flatness simultaneously with sparsity : strong/robust performance across various sparsity/settings



(Left) SAFE successfully finds sparse models at flat minima (Center) SAFE achieves strong image classification performance (Right) SAFE and SAFE+ achieves strong LLM pruning performance



Contact: {dongyeop.lee\_, kwanhee.lee, jinseokchung, namhoon.lee}@postech.ac.kr

			LLal	LLaMa-3				
			7B		13B			
Sparsity	Method	Wikitext/C4		Wiki	text/C4	Wikitext/C4		
0%	Dense	5.47 / 7.26		4.88	/ 6.72	6.23	/ 9.53	
50%	Magnitude SparseGPT Wanda ALPS SAFE SAFE	$16.03 \\ 6.99_{\pm 0.03} \\ 6.92_{\pm 0.01} \\ 6.87_{\pm 0.01} \\ \underline{6.78}_{\pm 0.01} \\ 6.56_{\pm 0.01}$	/ 21.33 / 9.20 $_{\pm 0.03}$ / 9.23 $_{\pm 0.00}$ / 8.98 $_{\pm 0.00}$ / 8.93 $_{\pm 0.00}$	$6.82 \\ 6.06_{\pm 0.03} \\ 5.98_{\pm 0.01} \\ 5.96_{\pm 0.02} \\ \underline{5.76}_{\pm 0.01} \\ 5.67_{\pm 0.01} $	/ 9.37 / $8.20_{\pm 0.01}$ / $8.28_{\pm 0.01}$ / $8.09_{\pm 0.04}$ / $7.85_{\pm 0.02}$	134.20 9.36±0.11 9.71±0.03 <u>9.05</u> ±0.12 9.59±0.06 <b>8.62</b> ±0.06	/ 273.3 / 13.96±0.02 / 14.88±0.04 / <u>13.40</u> ±0.06 / 14.60±0.04 / <b>13.26</b> ±0.06	
60%	Magnitude SparseGPT Wanda ALPS SAFE SAFE	$1864 \\ 10.19_{\pm 0.08} \\ 10.75_{\pm 0.07} \\ 9.55_{\pm 0.00} \\ \underline{9.20}_{\pm 0.04} \\ 8.30_{\pm 0.06}$	/ 2043 / 12.86 $_{\pm 0.05}$ / 13.87 $_{\pm 0.01}$ / <u>11.24<math>_{\pm 0.03}</math></u> / 11.51 $_{\pm 0.04}$ / <b>10.59</b> $_{\pm 0.00}$	$11.81$ $8.31_{\pm 0.09}$ $8.43_{\pm 0.07}$ $7.54_{\pm 0.03}$ $\overline{7.18}_{\pm 0.03}$ $6.78_{\pm 0.04}$	/ 14.62 / 10.85 $_{\pm 0.09}$ / 11.55 $_{\pm 0.01}$ / 9.87 $_{\pm 0.05}$ / 9.59 $_{\pm 0.03}$ / 9.02 $_{\pm 0.15}$	$\begin{array}{r} 5335\\ 15.46_{\pm 0.40}\\ 22.06_{\pm 0.19}\\ \underline{14.03}_{\pm 0.35}\\ 15.90_{\pm 0.25}\\ 12.18_{\pm 0.22}\end{array}$	/ 7438 / 21.25 $_{\pm 0.18}$ / 32.28 $_{\pm 0.37}$ / <u>18.72<math>_{\pm 0.15}</math></u> / 22.26 $_{\pm 0.16}$ / <b>17.30</b> $_{\pm 0.02}$	
4:8	Magnitude SparseGPT Wanda ALPS SAFE SAFE SAFE <sup>+</sup>	$\begin{array}{c} 15.91 \\ 8.42 {\scriptstyle \pm 0.05} \\ 8.64 {\scriptstyle \pm 0.03} \\ \underline{8.11} {\scriptstyle \pm 0.09} \\ 8.21 {\scriptstyle \pm 0.01} \\ \textbf{7.59} {\scriptstyle \pm 0.03} \end{array}$	/ $31.61$ / $10.73_{\pm 0.03}$ / $11.35_{\pm 0.01}$ / $10.21_{\pm 0.04}$ / $10.61_{\pm 0.04}$ / <b>9.88</b> $_{\pm 0.01}$	$\begin{array}{c} 7.32 \\ 7.02 {\pm} 0.06 \\ 7.01 {\pm} 0.02 \\ 6.81 {\pm} 0.07 \\ 6.60 {\pm} 0.02 \\ \textbf{6.37} {\pm} 0.03 \end{array}$	/ 9.96 / 9.33 $_{\pm 0.04}$ / 9.70 $_{\pm 0.03}$ / 9.33 $_{\pm 0.04}$ / <u>8.95<math>_{\pm 0.02}</math> / <b>8.61</b><math>_{\pm 0.01}</math></u>	$\begin{array}{c} 212.5\\ 12.16_{\pm 0.20}\\ 13.84_{\pm 0.04}\\ \underline{11.38}_{\pm 0.17}\\ 12.15_{\pm 0.14}\\ \textbf{10.51}_{\pm 0.13} \end{array}$	/ 336.3 / 17.36 $_{\pm 0.06}$ / 21.14 $_{\pm 0.06}$ / <u>16.10<math>_{\pm 0.10}</math></u> / 17.90 $_{\pm 0.15}$ / <b>15.67</b> $_{\pm 0.02}$	
2:4	Magnitude SparseGPT Wanda ALPS SAFE SAFE SAFE <sup>+</sup>	$\begin{array}{c} 37.77\\ 11.00 {\scriptstyle \pm 0.20}\\ 12.17 {\scriptstyle \pm 0.02}\\ \underline{9.99} {\scriptstyle \pm 0.19}\\ 10.53 {\scriptstyle \pm 0.13}\\ \textbf{8.96} {\scriptstyle \pm 0.07} \end{array}$	/ 74.70 / 13.54 $_{\pm 0.03}$ / 15.60 $_{\pm 0.11}$ / <u>12.04<math>_{\pm 0.04}</math></u> / 13.20 $_{\pm 0.07}$ / <b>11.34</b> $_{\pm 0.03}$	$\begin{array}{c} 8.88\\ 8.78_{\pm 0.09}\\ 9.01_{\pm 0.04}\\ 8.16_{\pm 0.17}\\ \hline \underline{7.64}_{\pm 0.05}\\ \hline \textbf{7.20}_{\pm 0.04}\end{array}$	/ 11.72 / 11.26 $\pm$ 0.11 / 12.40 $\pm$ 0.01 / 10.35 $\pm$ 0.18 / <u>10.10</u> $\pm$ 0.01 / <b>9.52</b> $\pm$ 0.01	$\begin{array}{c} 792.8\\ 15.87_{\pm 0.32}\\ 23.03_{\pm 0.38}\\ \underline{14.53}_{\pm 0.33}\\ 17.49_{\pm 0.27}\\ \textbf{13.39}_{\pm 0.23} \end{array}$	/ 2245 / 22.45 $_{\pm 0.12}$ / 34.91 $_{\pm 0.31}$ / <u>19.74<math>_{\pm 0.18}</math></u> / 24.45 $_{\pm 0.13}$ / <b>19.03<math>_{\pm 0.01}</math></b>	



#### z-minimization: Euclidean projection onto sparsity constraint

- $z_{k+1} = \arg\min_{z} I_{\|\cdot\|_0 \le d}(z) + \frac{\lambda}{2} \|x_{k+1} z + u_k\|_{1}$  $= \operatorname{proj}_{\|\cdot\|_0 \le d}(x_{k+1} + u_k).$
- z-minimization corresponds to projecting  $x_{k+1} + u_k$ onto the sparsity constraint in terms of Euclidean distance.
- This leads to the classic hard thresholding operator, where we zero out except d elements with the largest magnitude.

## **SAFE+:** Improving projection through generalized distance

#### Generalized quadratic distance to improve projection

$$\begin{aligned} z_{k+1} &= \operatorname{proj}_{\|\cdot\|_0 \le d}^{\mathbf{P}} (x_{k+1} + u_k) \\ &:= \operatorname*{arg\,min}_{\|z\|_0 \le d} \frac{1}{2} \|z - (x_{k+1} + u_k)\|_{\mathbf{P}}^2 \\ &= \operatorname*{arg\,min}_{\|z\|_0 \le d} \frac{1}{2} (z - (x_{k+1} + u_k))^{\top} \mathbf{P} (z - (x_{k+1} + u_k)). \end{aligned}$$

- However, this magnitude-based projection often yields subpar performance in practice.
- To improve this, we introduce a generalized distance  $\frac{1}{2} \| \cdot \|_{\mathbf{P}}^2$  with diagonal positive definite matrix P

#### Distance vs. saliencies

Criteria	Р
Magnitude	Ι
OBD	$\operatorname{diag}(H)$
SNIP	$\operatorname{diag}(\nabla f \nabla f^{\top})$
Wanda	$\operatorname{diag}(\mathbf{A}^{\top}\mathbf{A})$

- This generalized projection framework allows us to employ various saliency scores within the projection step
- Here we use this primarily for LLM pruning, though it is generally applicable to other domains

# **SAFE** demonstrates strong empirical performance

#### -1. Successful sparsity and flatness enforcement

- Weights are concentrated near zero  $\rightarrow$  **sparse**
- Wide minima + smaller largest Hessian eigenvalue  $\rightarrow$  **flat**

#### $\leftarrow 2$ . Strong image classification performance

- SAFE retains strong performance at high sparsities compared to baselines.
- Epoch: 200 (300 for ResNet-20) / Recompute batch statistics after final hard projection step

#### ←3. Strong LLM pruning performance

- Compared to methods specifically designed for LLMs, SAFE performs competitively and SAFE<sup>+</sup> outperforms all baselines across all models and sparsities.
- <u>Blockwise SAFE</u>: Following common practice, we sequentially apply SAFE to each transformer block to minimize the reconstruction error as  $\min_{\|\mathbf{x}\|_{0} \leq d} \|\mathsf{Block}(\mathsf{inputs}; \mathbf{x}) - \mathsf{Block}(\mathsf{inputs}; x_{\mathsf{original}})\|$ .
- SAFE<sup>+</sup>: Projection based on diagonal Hessian of the Layerwise reconstruction error is employed, which corresponds to using the Wanda score.
- Epoch: 30 / Dataset: 128 random samples from the first shard of the C4 dataset / seq len: 2048
- (Sidenote) As model scales, the computation of SAFE scales quadratically with the model width. In contrast, methods such as SparseGPT and ALPS scales cubically (see Appendix E).

#### 4. Strong robustness to various data noise

SAFE is robust to noisy label training (left), as well as common image corruptions and adversarial attacks (right) on ResNet-20/CIFAR-10.

Noise ratio						Common co	prruption (avg.)	Adversarial		
Sparsity	Method	25%	50%	75%	Sparsity	Method	intensity=3	intensity=5	$l_{\infty}$ -PGD	l <sub>2</sub> -PGD
70%	ADMM	$77.00_{\pm 0.91}$	$59.18_{\pm 0.55}$	$32.62_{\pm 0.89}$	000/	ADMM	<b>70.06</b> <sub>+0.03</sub>	$52.01_{\pm 0.38}$	<b>49.81</b> <sub>+1.02</sub>	<b>49.71</b> <sub>+1.06</sub>
	SAFE	$90.58_{\pm 0.30}$	86.51 $_{\pm 0.16}$	$67.01_{\pm 0.54}$	90%	SAFE	<b>73.98</b> $_{\pm 0.09}$	$55.11_{\pm 0.27}$	<b>56.43</b> <sub>±1.03</sub>	<b>56.36</b> $_{\pm 1.11}$
80%	ADMM	<b>76.18</b> $_{\pm 0.56}$	$62.67_{\pm 0.38}$	$32.86_{\pm 1.12}$	050/	ADMM	<b>68.87</b> <sub>+0.25</sub>	$50.56_{\pm 0.07}$	<b>49.84</b> <sub>+1.78</sub>	<b>49.68</b> +1.79
	SAFE	$91.25_{\pm 0.12}$	$86.55_{\pm 0.07}$	<b>66.49</b> $_{\pm 0.56}$	95%	SAFE	<b>72.92</b> $_{\pm 0.41}$	<b>54.86</b> <sub>±0.51</sub>	$51.40_{\pm 0.89}$	$51.36_{\pm 0.94}$
90%	ADMM	<b>79.40</b> $_{\pm 0.12}$	<b>66.64</b> $_{\pm 0.13}$	<b>36.84</b> <sub>±0.94</sub>	98%	ADMM	<b>65.46</b> <sub>±0.24</sub>	<b>48.65</b> <sub>±0.04</sub>	<b>43.33</b> <sub>±1.59</sub>	<b>43.42</b> <sub>±1.60</sub>
	SAFE	<b>90.68</b> $_{\pm 0.21}$	$86.49_{\pm 0.06}$	<b>64.72</b> ±0.61		SAFE	$68.20_{\pm 0.47}$	<b>49.96</b> <sub>±0.83</sub>	$43.34_{\pm 0.90}$	$43.41_{\pm 1.03}$
95%	ADMM	$77.71_{\pm 0.52}$	$67.10_{\pm 1.37}$	<b>39.68</b> <sub>±1.44</sub>	99%	ADMM	$59.21_{\pm 0.47}$	$43.81_{\pm 0.44}$	<b>30.29</b> <sub>±0.64</sub>	$30.32_{\pm 0.58}$
	SAFE	$89.86_{\pm 0.11}$	$85.18_{\pm 0.15}$	$64.25_{\pm 0.36}$		SAFE	$66.02_{\pm 0.56}$	<b>49.34</b> $_{\pm 1.03}$	<b>43.70</b> <sub>±1.28</sub>	<b>32.70</b> <sub>±1.28</sub>
					00 50/	ADMM	<b>55.72</b> $_{\pm 0.44}$	$41.55_{\pm 0.78}$	$23.25_{\pm 1.92}$	$23.25_{\pm 1.85}$
					99.0%	SAFE	$56.58_{\pm 0.36}$	$42.27_{\pm 0.63}$	<b>29.48</b> $_{\pm 0.68}$	<b>29.45</b> $_{\pm 0.74}$

### **SAFE** converges to stationary point within sparsity constraint

( $\delta$ -stationary point) We say a point  $\bar{x}$  is a  $\delta$ -stationary point of the sparsity-constrained optimization problem if  $\bar{x} \in \operatorname{arg\,min}_{a \in \mathcal{A}} \left\| a - \left( \bar{x} - \delta^{-1} \nabla f(\bar{x}) \right) \right\|$ ,

(Convergence of SAFE) Suppose that f is smooth and weakly convex. Assume further that  $\delta$  is chosen large enough so that  $\delta^{-1}\beta^2 - (\delta - \mu)/2 < 0$ . Let  $(\bar{x}, \bar{z}, \bar{u})$  be a limit point of SAFE algorithm. Then  $\bar{x}$  is a  $\delta$ -stationary point of the sparsity-constrained optimization problem.