

Neural Pruning via Growing Regularization

ICLR 2021 (Poster, Paper ID: 555)



Huan Wang



Can Qin



Yulun Zhang



Yun Fu

Northeastern University, Boston, MA, USA

Introduction: DNNs Effective, but not efficient (motivation)





How to do pruning: A typical way and its problem







Pruning criterion: Prune the weights whose absence leads to the *least* loss increase. **Solutions:**

(1) Trial and error? Too many weights (millions and billions params!), we cannot ablate them one by one.

(2) Use some analytical formula \Rightarrow Taylor expansion \Rightarrow approximate the Hessian (very hard for DNNs).

$$\Delta \mathcal{L} = \underbrace{\frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}}}_{\approx 0}^{\top} \Delta \boldsymbol{\theta} + \frac{1}{2} \Delta \boldsymbol{\theta}^{\top} \mathbf{H} \Delta \boldsymbol{\theta} + \mathcal{O}(||\Delta \boldsymbol{\theta}||^3)$$

Can we leverage the Hessian *without* knowing their specific values?

Leverage Hessian without knowing Hessian



□ How: Hessian can play its role by affecting weight magnitude (L1-norm).

Magnitude is known to be not accurate for pruning? Not necessarily. It is accurate as long as the gap is significant enough.





Parameter update in SGD:

0

$$w = w - \alpha(diff - \lambda w)$$

 α : learning rate, diff: gradient

when a weight stops update

 $\lambda w = diff$, where λ is the factor of weight decay.

(Ideally) The final position of a weight will be determined by both the task supervision and the regularization (task-agnostic prior), where λ can hold the balance.



$$\rightarrow$$
 $r_1 =$





 w_1

How r_1 and r_2 are related to the underlying Hessian?

Theoretical Analysis



$$\begin{aligned} \hat{\mathbf{w}}^* &= (\mathbf{H} + \delta\lambda \,\mathbf{I})^{-1}\mathbf{H}\mathbf{w}^*, \\ \text{(local quadratic approximation)} \end{aligned}$$
Case 1: Hessian is diagonal
$$\hat{w}_i^* &= \frac{h_{ii}}{h_{ii} + \delta\lambda} w_i^*, \Rightarrow r_i = \frac{\hat{w}_i^*}{w_i^*} = \frac{1}{\delta\lambda/h_{ii} + 1}, \\ \text{Case 2: Hessian is not diagonal (2d analysis)} \end{aligned}$$

$$\begin{cases} \hat{w}_1^* \\ \hat{w}_2^* \\ \end{pmatrix} &= \frac{1}{|\hat{\mathbf{H}}|} \begin{cases} (h_{11}h_{22} + h_{11}\delta\lambda - h_{12}^2)w_1^* + \delta\lambda h_{12}w_2^* \\ (h_{11}h_{22} + h_{22}\delta\lambda - h_{12}^2)w_2^* + \delta\lambda h_{12}w_1^* \\ \end{pmatrix} \approx \frac{1}{|\hat{\mathbf{H}}|} \begin{cases} (h_{11}h_{22} + h_{11}\delta\lambda - h_{12}^2)w_1^* \\ (h_{11}h_{22} + h_{22}\delta\lambda - h_{12}^2)w_2^* + \delta\lambda h_{12}w_1^* \\ \end{pmatrix} \approx \frac{1}{|\hat{\mathbf{H}}|} \begin{cases} (h_{11}h_{22} + h_{11}\delta\lambda - h_{12}^2)w_1^* \\ (h_{11}h_{22} + h_{22}\delta\lambda - h_{12}^2)w_2^* \\ \end{pmatrix}, \\ \Rightarrow r_1 &= \frac{1}{|\hat{\mathbf{H}}|} (h_{11}h_{22} + h_{11}\delta\lambda - h_{12}^2), r_2 &= \frac{1}{|\hat{\mathbf{H}}|} (h_{11}h_{22} + h_{22}\delta\lambda - h_{12}^2). \\ \end{aligned}$$
Conclusion: $h_{11} > h_{22} \Rightarrow r_1 > r_2$

A weight lying on a *sharper* local minimum will be pushed *less* towards zero.

Empirical Validation of the Theoretical Analysis



- The weight magnitude gap will be larger and larger.
- Eventually, the simple L1-norm will suffice to make a faithful criterion.



Figure 1: Row 1: Illustration of weight separation as L_2 penalty grows. Row 2: Normalized filter L_1 -norm over iterations for ResNet50 layer2.3.conv1 (please see the Appendix for VGG19 plots).

Regularization can also help in pruning schedule



$$\begin{aligned} \mathcal{E}(\mathbf{w}, \mathcal{D}) &= \mathcal{L}(\mathbf{w}, \mathcal{D}) + \frac{1}{2}\lambda \|\mathbf{w}\|_{2}^{2}, \\ \lambda_{j} &= \lambda_{j} + \delta\lambda, \text{ for all } j. \end{aligned}$$
 GReg-2

$$\begin{aligned} \mathcal{E}(\mathbf{w}, \mathcal{D}) &= \mathcal{L}(\mathbf{w}, \mathcal{D}) + \frac{1}{2}\lambda \|\mathbf{w}\|_2^2, \\ \lambda_j &= \lambda_j + \delta\lambda, j \in \{j \mid M[j] = 0\} \end{aligned}$$
 GReg-1

- Same pruning criterion as L1-norm pruning [1]: sort all the filters by their L1-norm, select those with the least L1-norms to prune (i.e., mask = 0).
- Two pruning schedules: (1) one-shot: remove the unimportant weights immediately; (2) pushing them towards zero first by the proposed growing regularization, then remove them.

[1] Li, Hao, et al. "Pruning Filters for Efficient ConvNets." ICLR (2017).



Table 1: Comparison between pruning schedules: one-shot pruning vs. our proposed GReg-1. Each setting is randomly run for 3 times, mean and std accuracies reported.

ResNet56 + CIFAR10: Baseline accuracy 93.36%, #Params: 0.85M, FLOPs: 0.25G									
Pruning ratio r (%)	50	70	90	92.5	95				
Sparsity (%) / Speedup	49.82/1.99×	$70.57/3.59 \times$	90.39/11.41×	93.43/14.76×	95.19/19.31×				
Acc. (%, L_1 +one-shot)	$92.97_{\pm 0.15}$	$91.88_{\pm 0.09}$	$87.34_{\pm 0.21}$	$87.31_{\pm 0.28}$	$82.79_{\pm 0.22}$				
Acc. (%, GReg-1, ours)	$93.06_{\pm 0.09}$	$92.23_{\pm 0.21}$	$89.49_{\pm 0.23}$	$88.39_{\pm 0.15}$	$85.97_{\pm 0.16}$				
Acc. gain (%)	0.09	0.35	2.15	1.08	3.18				
VGG19 + CIFAR100: Baseline accuracy 74.02%, #Params: 20.08M, FLOPs: 0.80G									
VGG19 + CII	FAR100: Baselin	ne accuracy 74.0)2%, #Params: 20).08M, FLOPs: 0	.80G				
VGG19 + CII Pruning ratio r (%)	FAR100: Baselin 50	ne accuracy 74.0	02%, #Params: 20 70	0.08M, FLOPs: 0 80	.80G 90				
VGG19 + CII Pruning ratio r (%) Sparsity (%) / Speedup	50 74.87/3.60×	ne accuracy 74.0 60 84.00/5.41×	02%, #Params: 20 70 90.98/8.84×	0.08M, FLOPs: 0 80 95.95/17.30×	.80G 90 98.96/44.22×				
VGG19 + CIIPruning ratio r (%)Sparsity (%) / SpeedupAcc. (%, L_1 +one-shot)	FAR100: Baselin 50 74.87/3.60× 71.49 _{±0.14}	$ \frac{100}{60} $ 84.00/5.41× 70.27+0.12	$\begin{array}{r} \hline 02\%, \text{\#Params: } 20\\\hline 70\\90.98/8.84\times\\\hline 66.05_{\pm 0.04}\end{array}$	0.08M, FLOPs: 0 80 95.95/17.30× 61.59 _{±0.03}	.80G 90 98.96/44.22× 51.36 _{±0.11}				
VGG19 + CIIPruning ratio r (%)Sparsity (%) / SpeedupAcc. (%, L_1 +one-shot)Acc. (%, GReg-1, ours)	FAR100: Baselin 50 $74.87/3.60 \times$ $71.49_{\pm 0.14}$ $71.50_{\pm 0.12}$	ne accuracy 74.0 60 $84.00/5.41 \times$ $70.27_{\pm 0.12}$ $70.33_{\pm 0.12}$	$\begin{array}{r} 02\%, \# \text{Params: } 20\\\hline 70\\90.98/8.84\times\\\hline 66.05_{\pm 0.04}\\\mathbf{67.35_{\pm 0.15}}\end{array}$	$\begin{array}{r} 0.08 \text{M}, \text{FLOPs: 0} \\ 80 \\ 95.95/17.30 \times \\ 61.59_{\pm 0.03} \\ \textbf{63.55}_{\pm 0.29} \end{array}$					

Given exactly the same weights to prune, different pruning schedules can lead to starkly different performances.

- GReg-1 > One-shot.
- Larger pruning ratio, the advantage of GReg-1 is more pronounced.



Method	Network	Base top-1 (%)	Pruned top-1 (%)	Top-1 drop	Speedup
L_1 (pruned-B) Li et al. (2017)		73.23	72.17	1.06	1.32×
Taylor-FO Molchanov et al. (2019)	DecNet24	73.31	72.83	0.48	$1.29 \times$
GReg-1 (ours)	Residents+	73.31	73.54	-0.23	$1.32 \times$
GReg-2 (ours)		73.31	73.61	-0.30	1.32 ×
ProvableFP Liebenwein et al. (2020)	PocNot50	76.13	75.21	0.92	1.43×
GReg-1 (ours)	Residence	76.13	76.27	-0.14	1.49 ×
AOFP Ding et al. (2019b)	D. N. (50	75.34	75.63	-0.29	1.49 × ⁻
GReg-2 (ours)*	Resiver50	75.40	76.13	-0.73	1.49 ×
IncReg Wang et al. (2019b)		75.60	72.47	3.13	$2.00 \times$
SFP He et al. (2018a)		76.15	74.61	1.54	$1.72 \times$
HRank Lin et al. (2020a)		76.15	74.98	1.17	$1.78 \times$
Taylor-FO Molchanov et al. (2019)		76.18	74.50	1.68	$1.82 \times$
Factorized Li et al. (2019)	ResNet50	76.15	74.55	1.60	2.33 ×
DCP Zhuang et al. (2018)		76.01	74.95	1.06	$2.25 \times$
CCP-AC Peng et al. (2019)		76.15	75.32	0.83	$2.18 \times$
GReg-1 (ours)		76.13	75.16	0.97	$2.31 \times$
GReg-2 (ours)		76.13	75.36	0.77	$2.31 \times$
C-SGD-50 Ding et al. (2019a)		75.34	74.54	$- \bar{0}.\bar{80}^{-} -$	$2.26 \times$
AOFP Ding et al. (2019b)	ResNet50	75.34	75.11	0.23	2.31 ×
GReg-2 (ours)*		75.40	75.22	0.18	2.31 ×
LFPC He et al. (2020)		76.15	74.46	1.69	$2.55 \times$
GReg-1 (ours)	ResNet50	76.13	74.85	1.28	2.56 imes
GReg-2 (ours)		76.13	74.93	1.20	2.56 imes
IncReg Wang et al. (2019b)		75.60	71.07	4.53	$3.00 \times$
Taylor-FO Molchanov et al. (2019)	RecNet50	76.18	71.69	4.49	$3.05 \times$
GReg-1 (ours)	Resilietato	76.13	73.75	2.38	3.06 ×
GReg-2 (ours)		76.13	73.90	2.23	3.06 ×

Table 3: Acceleration comparison on ImageNet. FLOPs: ResNet34: 3.66G, ResNet50: 4.09G.

* Since the base models of C-SGD and AOFP have a much lower accuracy than ours, for fair comparison, we train our own base models with similar accuracy.



Method	Base top-1 (%)	Pruned top-1 (%)	Top-1 drop	Sparsity (%)
GSM Ding et al. (2019c)	75.72	74.30	1.42	80.00
Variational Dropout Molchanov et al. (2017a)	76.69	75.28	1.41	80.00
DPF Lin et al. (2020b)	75.95	74.55	1.40	82.60
WoodFisher Singh & Alistarh (2020)	75.98	75.20	0.78	82.70
GReg-1 (ours)	76.13	75.45	0.68	82.70
GReg-2 (ours)	76.13	75.27	0.86	82.70

Table 4: Compression comparison on ImageNet with ResNet50. #Parameters: 25.56M.

The proposed methods can work *seamlessly* from filter pruning to unstructured pruning.

- Filter pruning: weight group is a filter.
- Unstructured pruning: weight group is a single weight element.

Conclusion



- we present two algorithms that exploit regularization in a new fashion that the penalty factor is uniformly raised to a large amount.
- □ The two algorithms show that:
- GReg-1: Pruning schedule is another important axis in pruning, which may deserve more research attention.
- GReg-2: Without any Hessian approximation, we can still tap into its power for pruning with the help of growing L2 regularization.
- □ Empirically, both algorithms achieve very promising results compared to many recent methods.



Our code and trained models are released at: https://github.com/MingSun-Tse/Regularization-Pruning

Great thanks for your attention!

Acknowledgements: This work is supported by the National Science Foundation Award ECCS-1916839 and the U.S. Army Research Office Award W911NF-17-1-0367.