
A Flexible, Extensible Software Framework for Neural Net Compression

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We propose a software framework based on ideas of the Learning-Compression
2 algorithm [1, 2, 4], that allows one to compress any neural network by different
3 compression mechanisms (pruning, quantization, low-rank, etc.). By design, the
4 learning of the neural net (handled by SGD) is decoupled from the compression of
5 its parameters (handled by a signal compression function), so that the framework
6 can be easily extended to handle different combinations of neural net and compression
7 type. In addition, it has other advantages, such as easy integration with
8 deep learning frameworks, efficient training time, competitive practical performance
9 in the loss-compression tradeoff, and reasonable convergence guarantees.
10 Our toolkit is written in Python and Pytorch and we plan to make it available by
11 the workshop time, and eventually open it for contributions from the community.

12 With the great success of neural network in solving practical problems in various fields (vision, NLP,
13 etc.) there has been an emergence of research in neural network compression techniques that allows
14 to compress these large models in terms of memory, computation and/or power requirements. At
15 present many ad-hoc solutions have been proposed that typically solve *one specific type of compression*
16 (binarization and quantization [2, 5, 8, 11, 18, 23, 24], pruning [9, 10, 15, 16, 21], low-rank or
17 tensor factorization [6, 7, 12–14, 17, 19, 20, 22], etc.), as well as several submissions to the present
18 workshop.

19 Among the various research strands in neural net compression, in our view a fundamental problem
20 is that in practice one does not know what the best type of compression (or combination of compression
21 types) may be the best for a given network. In principle, it would be possible to try different
22 existing algorithms, assuming one can find an implementation for them. We seek a solution that
23 directly addresses this problem and can potentially allow non-expert end users to compress models
24 easily and effectively¹. It is based on a recently proposed compression framework, the LC algorithm
25 [1, 2, 4], that by design separates the “learning” part of the problem (which involves the dataset,
26 neural net model and loss function) from the “compression” part (how the network parameters will
27 be compressed). This has the advantage of modularity (a pillar of structured programming and software
28 development): we can change the compression type by simply calling a different compression
29 routine (e.g. k -means instead of the SVD) within the overall algorithm *in a principled way*, with
30 no other changes to the algorithm. We briefly describe the LC algorithm below and mention additional
31 advantages it provides—such as easy integration with deep learning frameworks, efficient
32 training time, competitive practical performance in the loss-compression tradeoff, and reasonable
33 convergence guarantees.

¹There have been some attempts to include compression as a black-box routines into deep learning frameworks (e.g. https://www.tensorflow.org/performance/post_training_quantization), but they are limited to a subset of simple methods.

34 In this paper, we describe our ongoing efforts in building a software implementation that can cap-
 35 italize on the modularity of the LC algorithm. At present this handles 1) (C step) various forms
 36 of quantization, pruning and low-rank compression, and we will soon add combinations of those
 37 and further compression types; and 2) (L step) various types of deep net models. Our framework is
 38 written in Python and Pytorch. We plan to make it available online as open source by the time of the
 39 workshop. We also hope that interested researchers and developers will eventually contribute their
 40 own routines for signal compression or for training of specific neural net architectures.

41 1 Model compression as a constrained optimization problem

42 Assume we have a previously trained model with weights $\bar{\mathbf{w}} = \arg \min_{\mathbf{w}} L(\mathbf{w})$. This is our *refer-*
 43 *ence* model, which represents the best loss we can achieve without compression. The “Learning-
 44 Compression” paper [1] defines *compression* as finding a low-dimensional parameterization $\Delta(\Theta)$
 45 of \mathbf{w} in terms of $Q < P$ parameters Θ . The goal is to find such Θ that its corresponding model has
 46 (locally) optimal loss. Therefore the *model compression as a constrained optimization* problem is
 47 defined as:

$$\min_{\mathbf{w}, \Theta} L(\mathbf{w}) \text{ s.t. } \mathbf{w} = \Delta(\Theta) \quad (1)$$

48 Compression and decompression are usually seen as algorithms, but here they are regarded as math-
 49 ematical mappings in parameter space. The *decompression mapping* $\Delta: \Theta \in \mathbb{R}^Q \rightarrow \mathbf{w} \in \mathbb{R}^P$ maps
 50 a low-dimensional parameterization to uncompressed model weights and the *compression mapping*
 51 $\Pi(\mathbf{w}) = \arg \min_{\Theta} \|\mathbf{w} - \Delta(\Theta)\|^2$ behaves as its “inverse”.

52 The problem in (1) is constrained, nonlinear, and potentially non-differentiable w.r.t. Θ (e.g. bi-
 53 naryzation). The LC-algorithm is obtained by converting this problem to an equivalent formulation
 54 using penalties and employing an alternating optimization. This results in an algorithm that alter-
 55 nates two generic steps while slowly driving the penalty parameter $\mu \rightarrow \infty$:

- 56 • **L (learning) step:** $\min_{\mathbf{w}} L(\mathbf{w}) + \frac{\mu}{2} \|\mathbf{w} - \Delta(\Theta)\|^2$. This is a regular training of the un-
 57 compressed model but with a quadratic regularization term. *This step is independent of the*
 58 *compression type.*
- 59 • **C (compression) step:** $\min_{\Theta} \|\mathbf{w} - \Delta(\Theta)\|^2 \Leftrightarrow \Theta = \Pi(\mathbf{w})$. This means finding the best
 60 (lossy) compression of \mathbf{w} (the current uncompressed model) in the ℓ_2 sense (orthogonal
 61 projection on the feasible set), and corresponds to our definition of the compression map-
 62 ping Π . *This step is independent of the loss, training set and task.*

63 The LC algorithm defines a continuous path $(\mathbf{w}(\mu), \Theta(\mu))$ which, under some mild assumptions,
 64 converges to a stationary point (typically a minimizer) of the constrained problem. The beginning
 65 of this path, for $\mu \rightarrow 0^+$, corresponds to training the reference model and then compressing it
 66 disregarding the loss (*direct compression*), a simple but suboptimal approach popular in practice.

67 **Optimizing the L and C steps** The L step can be solved by stochastic gradient descent (clipping
 68 the step sizes so they never exceed $\frac{1}{\mu}$). The C step can be solved by calling a compression routine
 69 corresponding to the desired compression type. For example, for quantization with an adaptive
 70 codebook we can run k -means with a codebook of size K weights [2]; for pruning, we can prune all
 71 but the top- κ weights (where κ depends on the sparsifying norm used) [3]. Fig. 1 gives a pseudo-
 72 code for the LC algorithm.

73 2 Overall software approach

74 In the pseudo-code there is only one compression mapping $\Delta(\Theta)$, however, practically we want to
 75 mix and match, e.g. compress first layer using quantization and all remaining ones using pruning.
 76 Therefore, there might be multiple constraints $\mathbf{w}_i = \Delta_i(\Theta_i)$ such that $\mathbf{w}_i \subset \mathbf{w}$ where \mathbf{w} is a set
 77 of all weights of the neural network. This means, during the C-step of the LC algorithm, instead of
 78 having one optimal compression $\Pi(\mathbf{w})$, we have multiple $\Pi_i(\mathbf{w}_i)$ which separates from each other
 79 and can be done in parallel.

80 **Data structures and functions defined by library.** Library defines several data structures that are
 81 necessary for internal housekeeping. It abstracts away from the architecture of the neural network,
 82 and sees its weights as a vector containing P values where P is cardinality of union of all \mathbf{w}_i , e.g.
 83 $P = |\mathbf{w}_1 \cup \mathbf{w}_2 \cup \dots \cup \mathbf{w}_k|$. This is achieved by a *weights view* data structure.

84 The compression mappings assumes a specific structure, e.g. low-rank compression expects a ma-
 85 trix, while pruning expects a vector. Therefore, we define a *compression view* data structure that
 86 allows to map \mathbf{w}_i to the format suitable to a *compression function*,

87 The *compression function* (Π_i) obtains \mathbf{w}_i in the suitable format and computes $\Pi_i(\mathbf{w}_i)$. Library
 88 comes equipped with a number of implemented compression functions: for quantization, pruning
 89 and low-rank. For example, if user wants no more than κ nonzero items, decompression mapping
 90 has form $\Delta(\Theta) = \Theta$ s.t. $\|\Theta\| \leq \kappa$, C-step will be $\Pi_i(\mathbf{w}_i) = \min_{\Theta_i} \|\mathbf{w}_i - \Theta_i\|^2$ s.t. $\|\Theta_i\| \leq \kappa$
 91 which solved by zeroing all but top κ values of \mathbf{w}_i in magnitude. Here we give a snippet from library
 92 that performs that:

```
93
pruned = np.zeros_like(weights)
indx = np.argsort(np.abs(weights), kind='mergesort')
remaining_indx = indx[-kappa:]
pruned[remaining_indx] = weights[remaining_indx]
```

94 **Input from the user.** For compression of a neural network, the software needs following
 95 things from the user: *L-step* implementation, *list of compression tasks* and a list of μ -values
 96 $(\mu_0, \mu_1, \dots, \mu_m)$. Note that the dataset, the loss function of the neural network is abstracted away
 97 from the library.

98 An *implementation of L-step* is a Python function which will be invoked with a new value of penalty
 99 (the term $\Delta(\Theta)$) and it must return an updated \mathbf{w} that minimizes the $L(\mathbf{w}) + \mu/2 \|\mathbf{w} - \Delta(\Theta)\|^2$.
 100 This step is similar to the reference network learning and requires minimal modifications to include
 101 penalty terms. We also note that it is independent of compression, and needs to be done once for a
 102 network, and can be re-used for other compression and combinations.

103 A *list of compression tasks* is list of tuples where each tuple contains: *compression view* of \mathbf{w}_i and
 104 *compression function* Π_i which maps subset of weights \mathbf{w} to a specific compression user wants to
 105 achieve, e.g. first layer to be quantized, second layer to be pruned.

106 **Library operations.** As soon as user invokes `lc.run` function with inputs described in previous
 107 paragraph, library initializes and/or parses necessary internal data structures and enters the loop over
 108 μ -values. For each μ value it invokes user supplied *L-step implementation* function and immediately
 109 afterwards, in parallel, traverses through *list of compression tasks*, obtains updated compression-
 110 view of \mathbf{w}_i and passes it to compression function Π_i to update Θ_i .

111 3 Conclusion

112 The fields of machine learning and signal compression have developed independently for a long
 113 time: machine learning solves the problem of training a deep net to minimize a desired loss on a
 114 dataset, while signal compression solves the problem of optimally compressing a given signal. Both
 115 are mature fields with well-understood algorithms and efficient implementations. Both converge
 116 in the problem of model compression, where we seek the smallest model (in the sense of memory,
 117 inference time or energy, etc.). The LC algorithm achieves this by seamlessly integrating the existing
 118 algorithms to train deep nets and to compress a signal. Here, we seek to develop an extensible
 119 software framework that can easily plug in existing deep net training techniques with existing signal
 120 compression techniques and their combinations. We intend to make the software open source by the
 121 time of the workshop and eventually to seek contributions from the community.

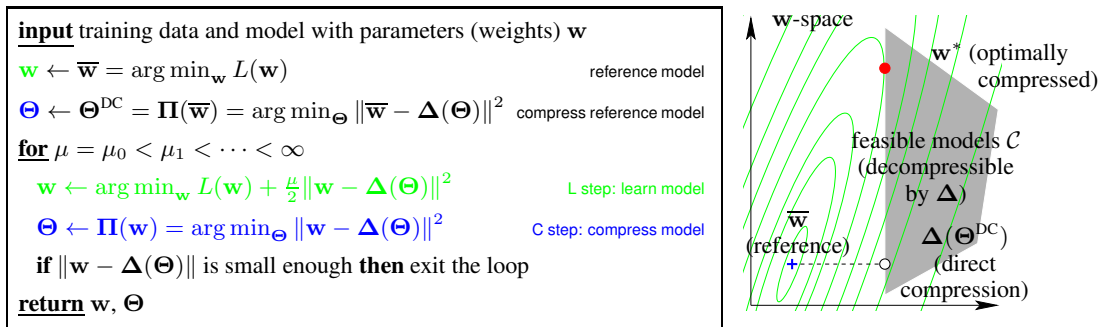


Figure 1: *Left*: A pseudo-code of LC algorithm. *Right*: an illustration of the idea of model compression by constrained optimization.

References

- 122
- 123 [1] M. Á. Carreira-Perpiñán. Model compression as constrained optimization, with application to
124 neural nets. Part I: General framework. arXiv:1707.01209 [cs.LG], July 5 2017.
- 125 [2] M. Á. Carreira-Perpiñán and Y. Idelbayev. Model compression as constrained optimization,
126 with application to neural nets. Part II: Quantization. arXiv:1707.04319 [cs.LG], July 13 2017.
- 127 [3] M. Á. Carreira-Perpiñán and Y. Idelbayev. Learning-compression algorithms for neural net-
128 work pruning. In *Proc. of the 2017 IEEE Computer Society Conf. Computer Vision and Pattern
129 Recognition (CVPR'17)*, Honolulu, HI, July 21–26 2017.
- 130 [4] M. Á. Carreira-Perpiñán and Y. Idelbayev. “learning-compression” algorithms for neural net
131 pruning. In *Proc. of the 2018 IEEE Computer Society Conf. Computer Vision and Pattern
132 Recognition (CVPR'18)*, Salt Lake City, UT, June 18–22 2018.
- 133 [5] M. Courbariaux, Y. Bengio, and J.-P. David. BinaryConnect: Training deep neural net-
134 works with binary weights during propagations. In C. Cortes, N. D. Lawrence, D. D. Lee,
135 M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems
136 (NIPS)*, volume 28, pages 3105–3113. MIT Press, Cambridge, MA, 2015.
- 137 [6] M. Denil, B. Shakibi, L. Dinh, M. Ranzato, and N. de Freitas. Predicting parameters in deep
138 learning. In C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q. Weinberger,
139 editors, *Advances in Neural Information Processing Systems (NIPS)*, volume 26, pages 2148–
140 2156. MIT Press, Cambridge, MA, 2013.
- 141 [7] E. L. Denton, W. Zaremba, J. Bruna, Y. LeCun, and R. Fergus. Exploiting linear struc-
142 ture within convolutional networks for efficient evaluation. In Z. Ghahramani, M. Welling,
143 C. Cortes, N. D. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information
144 Processing Systems (NIPS)*, volume 27, pages 1269–1277. MIT Press, Cambridge, MA, 2014.
- 145 [8] Y. Gong, L. Liu, M. Yang, and L. Bourdev. Compressing deep convolutional networks using
146 vector quantization. In *Proc. of the 3rd Int. Conf. Learning Representations (ICLR 2015)*, San
147 Diego, CA, May 7–9 2015.
- 148 [9] S. Han, J. Pool, J. Tran, and W. Dally. Learning both weights and connections for efficient
149 neural network. In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett,
150 editors, *Advances in Neural Information Processing Systems (NIPS)*, volume 28, pages 1135–
151 1143. MIT Press, Cambridge, MA, 2015.
- 152 [10] B. Hassibi and D. G. Stork. Second order derivatives for network pruning: Optimal brain
153 surgeon. In S. J. Hanson, J. D. Cowan, and C. L. Giles, editors, *Advances in Neural Information
154 Processing Systems (NIPS)*, volume 5, pages 164–171. Morgan Kaufmann, San Mateo, CA,
155 1993.
- 156 [11] K. Hwang and W. Sung. Fixed-point feedforward deep neural network design using weights
157 +1, 0, and −1. In *2014 IEEE Workshop on Signal Processing Systems (SiPS)*, pages 1–6,
158 Belfast, UK, Oct. 20–22 2014.
- 159 [12] Y. Ioannou, D. Robertson, R. Cipolla, and A. Criminisi. Deep roots: Improving cnn efficiency
160 with hierarchical filter groups. In *Proc. of the 2017 IEEE Computer Society Conf. Computer
161 Vision and Pattern Recognition (CVPR'17)*, pages 1231–1240, Honolulu, HI, July 21–26 2017.
- 162 [13] M. Jaderberg, A. Vedaldi, and A. Zisserman. Speeding up convolutional neural networks with
163 low rank expansions. In M. Valstar, A. French, and T. Pridmore, editors, *Proc. of the 25th
164 British Machine Vision Conference (BMVC 2014)*, Nottingham, UK, Sept. 1–5 2014.
- 165 [14] Y.-D. Kim, E. Park, S. Yoo, T. Choi, L. Yang, and D. Shin. Compression of deep convolutional
166 neural networks for fast and low power mobile applications. In *Proc. of the 4th Int. Conf.
167 Learning Representations (ICLR 2016)*, San Juan, Puerto Rico, May 2–4 2016.
- 168 [15] Y. LeCun, J. S. Denker, and S. A. Solla. Optimal brain damage. In D. S. Touretzky, editor, *Ad-
169 vances in Neural Information Processing Systems (NIPS)*, volume 2, pages 598–605. Morgan
170 Kaufmann, San Mateo, CA, 1990.

- 171 [16] B. Liu, M. Wan, H. Foroosh, M. Tappen, and M. Pinsky. Sparse convolutional neural networks.
172 In *Proc. of the 2015 IEEE Computer Society Conf. Computer Vision and Pattern Recognition*
173 (*CVPR'15*), pages 806–814, Boston, MA, June 7–12 2015.
- 174 [17] A. Novikov, D. Podoprikin, A. Osokin, and D. P. Vetrov. Tensorizing neural networks. In
175 C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in*
176 *Neural Information Processing Systems (NIPS)*, volume 28, pages 442–450. MIT Press, Cam-
177 bridge, MA, 2015.
- 178 [18] M. Rastegari, V. Ordonez, J. Redmon, and A. Farhadi. XNOR-net: ImageNet classification
179 using binary convolutional neural networks. In B. Leibe, J. Matas, N. Sebe, and M. Welling,
180 editors, *Proc. 14th European Conf. Computer Vision (ECCV'16)*, pages 525–542, Amsterdam,
181 The Netherlands, Oct. 11–14 2016.
- 182 [19] T. N. Sainath, B. Kingsbury, V. Sindhwani, E. Arisoy, and B. Ramabhadran. Low-rank ma-
183 trix factorization for deep neural network training with high-dimensional output targets. In
184 *Proc. of the IEEE Int. Conf. Acoustics, Speech and Sig. Proc. (ICASSP'13)*, pages 6655–6659,
185 Vancouver, Canada, Mar. 26–30 2013.
- 186 [20] C. Tai, T. Xiao, Y. Zhang, X. Wang, and W. E. Convolutional neural networks with low-rank
187 regularization. In *Proc. of the 4th Int. Conf. Learning Representations (ICLR 2016)*, San Juan,
188 Puerto Rico, May 2–4 2016.
- 189 [21] W. Wen, C. Wu, Y. Wang, Y. Chen, and H. Li. Learning structured sparsity in deep neural
190 networks. In D. D. Lee, M. Sugiyama, U. von Luxburg, I. Guyon, and R. Garnett, editors,
191 *Advances in Neural Information Processing Systems (NIPS)*, volume 29, pages 2074–2082.
192 MIT Press, Cambridge, MA, 2016.
- 193 [22] X. Zhang, J. Zou, K. He, and J. Sun. Accelerating very deep convolutional networks for
194 classification and detection. *IEEE Trans. Pattern Analysis and Machine Intelligence*, 38(10):
195 1943–1955, Oct. 2016.
- 196 [23] S. Zhou, Z. Ni, X. Zhou, H. Wen, Y. Wu, and Y. Zou. Dorefa-net: Training low bitwidth
197 convolutional neural networks with low bitwidth gradients. arXiv:1606.06160, July 17 2016.
- 198 [24] C. Zhu, S. Han, H. Mao, and W. J. Dally. Trained ternary quantization. In *Proc. of the 5th Int.*
199 *Conf. Learning Representations (ICLR 2017)*, Toulon, France, Apr. 24–26 2017.