
Fast Model Compression

Miguel Á. Carreira-Perpiñán and Arman Zharmagambetov
EECS, University of California, Merced

Motivation and summary In recent years, large neural networks have been successfully applied to various machine learning problems such as speech processing and computer vision. However, these neural nets contain a huge amount of parameters, which makes it difficult to deploy them in mobile phones or other devices with limited computation. This motivates the need for compressing a neural net while minimally hurting its performance. Many algorithms have been proposed that achieve significant compression based on pruning, quantization, low-rank decomposition and other techniques. However, most of these algorithms require access to the original training set. This imposes considerable resources in runtime and storage—for instance, modern image classification datasets such as ImageNet contain millions of high-resolution images. In some applications, for example in-device compression, it is desirable to do much faster (almost instant) compression.

In this work we focus on the framework of the “Learning-Compression” (LC) algorithm [1–3] because it can be applied to potentially any kind of compression type and combinations thereof. The basic idea is to replace the original loss function with an approximate, simpler loss that does not require access to the training set. The resulting optimization problem can be solved analytically or using the LC algorithm. We show that we can still achieve significant compression but much faster.

Fast model compression Assume we have a large, *reference* model with P parameters that has been trained on a loss L (e.g. cross-entropy on a given training set) to solve a task (e.g. classification). That is, $\bar{\mathbf{w}} = \arg \min_{\mathbf{w}} L(\mathbf{w})$. We define *compression* as finding a low-dimensional parameterization $\Delta(\Theta)$ of \mathbf{w} in terms of $Q < P$ parameters Θ . We seek a Θ such that its corresponding model has (locally) optimal loss. We define *model compression* as a constrained optimization problem: $\min_{\mathbf{w}, \Theta} \tilde{L}(\mathbf{w})$ s.t. $\mathbf{w} = \Delta(\Theta)$, where \tilde{L} is approximation of the original loss L which we define later. Note that the original LC algorithm uses the true loss L , therefore it requires access to dataset.

The *decompression mapping* $\Delta: \Theta \in \mathbb{R}^Q \rightarrow \mathbf{w} \in \mathbb{R}^P$ maps a low-dimensional parameterization to uncompressed model weights. The *compression mapping* $\Pi(\mathbf{w}) = \arg \min_{\Theta} \|\mathbf{w} - \Delta(\Theta)\|^2$ behaves as its “inverse” and appears in the C step of the LC algorithm. Our framework includes well-known types of compression (and combinations thereof), such as:

- *Pruning* defines $\mathbf{w} = \Delta(\theta) = \theta$ where \mathbf{w} is real and θ is constrained to have few nonzero values. The compression mapping involves some kind of thresholding.
- *Quantization* uses a discrete mapping Δ given by assigning each weight to one of K codebook values. The compression mapping is given by a form of rounding if we use a fixed codebook, such as binarization: $\{-1, +1\}$ (or by k -means if we use an adaptive codebook).
- *Low-rank compression* defines $\Delta(\mathbf{U}, \mathbf{V}) = \mathbf{U}\mathbf{V}^T$, where the weights matrix \mathbf{W} is constrained to be decomposed into low rank matrices \mathbf{U} and \mathbf{V} . The compression mapping is given by the singular value decomposition (SVD) of \mathbf{W} .

Approximation to the loss We approximate the original loss L using Taylor’s theorem. Assume a given loss on a training set, e.g. the cross-entropy loss $L(\mathbf{w}) = -\sum_n y_n \log f(\mathbf{x}_n; \mathbf{w})$. Then the loss function can be approximated as: $\tilde{L}(\mathbf{w}) = L_0 + \mathbf{g}^T(\mathbf{w} - \bar{\mathbf{w}}) + \frac{1}{2}(\mathbf{w} - \bar{\mathbf{w}})^T \mathbf{H}(\mathbf{w} - \bar{\mathbf{w}})$, where $L_0 = L(\bar{\mathbf{w}})$ is the loss value of the reference model, \mathbf{g} its gradient, \mathbf{H} its Hessian and $\bar{\mathbf{w}}$ are the weights of the reference model ($\bar{\mathbf{w}}$ need not be an exact minimizer, so its gradient need not be zero). This approximation is very good near $\bar{\mathbf{w}}$ but degrades progressively as we go away from it. Therefore we have to expect that the resulting solution will not be as accurate as the original LC algorithm. But this is the price we need to pay in order to get a fast compression. \mathbf{H} can be the full Hessian, diagonal, block diagonal, sparse or even zero if we use the first order approximation; here

we focus on the diagonal approximation. Then the above loss approximation takes the following form (by neglecting constant term L_0): $\tilde{L}(\mathbf{w}) = \sum_{i=1}^P [g_i(w_i - \bar{w}_i) + \frac{1}{2}h_i(w_i - \bar{w}_i)^2]$, where g_i and h_i are the elements of the gradient vector and diagonal elements of the Hessian, respectively.

Fast “Learning-Compression” (LC) algorithm This follows from using a penalty method (for simplicity, we describe the quadratic penalty) and alternating optimization, with the goal of separating the machine learning part (loss L) from the compression part (Δ). This results in an algorithm that alternates two generic steps while slowly driving the penalty parameter $\mu \rightarrow \infty$:

- **L (learning) step:** $\min_{\mathbf{w}} \sum_{i=1}^P [g_i(w_i - \bar{w}_i) + \frac{1}{2}h_i(w_i - \bar{w}_i)^2] + \frac{\mu}{2}\|\mathbf{w} - \Delta(\Theta)\|^2$. It is a separable quadratic optimization whose solution is $w_i = (h_i\bar{w}_i + \mu\Delta_i(\theta) - g_i)/(h_i + \mu)$. Note that in the original LC algorithm this step requires access to the training set and the neural net, which is computationally very costly (and requires SGD optimization in a GPU). Now the L step is data-independent and vastly faster.
- **C (compression) step:** $\min_{\Theta} \|\mathbf{w} - \Delta(\Theta)\|^2 \Leftrightarrow \Theta = \Pi(\mathbf{w})$. This means finding the best (lossy) compression of \mathbf{w} in the ℓ_2 sense. This step is identical to the original LC algorithm. It is independent of the loss, training set and task. It can be solved by calling a compression mapping (e.g. thresholding, SVD, etc.) corresponding to the desired compression type.

Analytical solution of the optimization problem In some particular cases the exact solution of the constrained optimization problem can be obtained analytically (without using iterative algorithms). For example, in the case of pruning [3], the optimization problem takes the form: $\min_{\mathbf{w}} \sum_{i=1}^P [g_i(w_i - \bar{w}_i) + \frac{1}{2}h_i(w_i - \bar{w}_i)^2]$ s.t. $\|\mathbf{w}\|_0 \leq \kappa$. Its exact solution is given by picking the weights having the largest values of $\alpha_i = g_i\bar{w}_i - \frac{1}{2}h_i\bar{w}_i^2 - g_i^2/(2h_i)$ and setting them to $w_i = \bar{w}_i - g_i/h_i$. If $\mathbf{g} = \mathbf{0}$ this solution corresponds to the Optimal Brain Damage algorithm of [4].

Another example is a problem of weight binarization: $\min_{\mathbf{w}} \sum_{i=1}^P [g_i(w_i - \bar{w}_i) + \frac{1}{2}h_i(w_i - \bar{w}_i)^2]$ s.t. $w_1, \dots, w_P \in \{-1, +1\}$. The problem separates over the weights and can be solved for each w_i by enumeration (try -1 and $+1$ and pick the one which gives the lowest value of the loss).

Experiments The figure shows compression results (as a tradeoff curve of error vs compression level) for the VGG-13 neural nets [5] on CIFAR-10. Currently, Tensorflow (and some other deep learning frameworks) are not able to provide just the diagonal of the Hessian. Therefore, we estimate it using the Gauss-Newton approximation. As we can see the fast LC algorithm is able to achieve low test error as long as we don’t compress much. The original LC shows better results on all experiments but its runtime is about 6 hours, whereas the fast compression runs only about 2 minutes for quantization, 1 second for pruning.

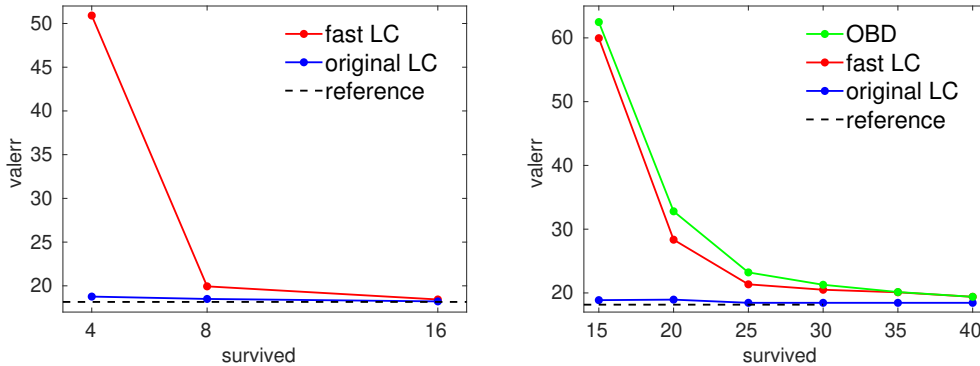


Figure 1: Error-compression curves for VGG-13 on CIFAR-10: quantization (left), pruning (right).

References

- [1] M. Á. Carreira-Perpiñán. Model compression as constrained optimization, with application to neural nets. Part I: General framework. arXiv:1707.01209 [cs.LG], July 5 2017.
- [2] M. Á. Carreira-Perpiñán and Y. Idelbayev. Model compression as constrained optimization, with application to neural nets. Part II: Quantization. arXiv:1707.04319 [cs.LG], July 13 2017.
- [3] M. Á. Carreira-Perpiñán and Y. Idelbayev. “learning-compression” algorithms for neural net pruning. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [4] Y. LeCun, J. S. Denker, and S. A. Solla. Optimal brain damage. In *Advances in Neural Information Processing Systems 2*, pages 598–605.
- [5] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *International Conference on Learning Representations*, 2015.