
Rate Distortion For Model Compression: From Theory To Practice

Weihao Gao¹ Yu-Han Liu² Chong Wang³ Sewoong Oh⁴

Abstract

The enormous size of modern deep neural networks makes it challenging to deploy those models in memory and communication limited scenarios. Thus, compressing a trained model without a significant loss in performance has become an increasingly important task. Tremendous advances have been made recently, where the main technical building blocks are pruning, quantization, and low-rank factorization. In this paper, we propose principled approaches to improve upon the common heuristics used in those building blocks, by studying the fundamental limit for model compression via the rate distortion theory. We prove a lower bound for the rate distortion function for model compression and prove its achievability for linear models. Although this achievable compression scheme is intractable in practice, this analysis motivates a novel objective function for model compression, which can be used to improve classes of model compressor such as pruning or quantization. Theoretically, we prove that the proposed scheme is optimal for compressing one-hidden-layer ReLU neural networks. Empirically, we show that the proposed scheme improves upon the baseline in the compression-accuracy tradeoff.

1. Introduction

Deep neural networks have been successful, for example, in the application of computer vision (Krizhevsky et al., 2012), machine translation (Wu et al., 2016) and game playing (Silver et al., 2017). With increasing data and compu-

tational power, the number of weights in practical neural network model also grows rapidly. For example, in the application of image recognition, the LeNet-5 model (LeCun et al., 1998) only has 400K weights. After two decades, AlexNet (Krizhevsky et al., 2012) has more than 60M weights, and VGG-16 net (Simonyan & Zisserman, 2014) has more than 130M weights. Coates et al. (2013) even tried a neural network with 11B weights. The huge size of neural networks brings many challenges, including large storage, difficulty in training, and large energy consumption. In particular, deploying such extreme models to embedded mobile systems is not feasible.

Several approaches have been proposed to reduce the size of large neural networks while preserving the performance as much as possible. Most of those approaches fall into one of the two broad categories. The first category designs novel network structures with small number of parameters, such as SqueezeNet (Iandola et al., 2016) and MobileNet (Howard et al., 2017). The other category directly compresses a given large neural network using pruning, quantization, and matrix factorization, including (LeCun et al., 1990; Hassibi & Stork, 1993; Han et al., 2015b;a; Cheng et al., 2015). There are also advanced methods to train the neural network using Bayesian methods to help pruning or quantization at a later stage, such as (Ulrich et al., 2017; Louizos et al., 2017; Federici et al., 2017).

As more and more model compression algorithms are proposed and compression ratio becomes larger and larger, it motivates us to think about the fundamental question — How well can we do for model compression? The goal of model compression is to trade off the *number of bits* used to describe the model parameters, and the *distortion* between the compressed model and original model. We wonder *at least* how many bits is needed to achieve certain distortion? Despite many successful model compression algorithms, these theoretical questions still remain unclear.

In this paper, we fill in this gap by bringing tools from rate distortion theory to identify the fundamental limit on how much a model can be compressed. Specifically, we focus on compression of a pretrained model, rather than designing new structures or retraining models. Our approach builds upon rate-distortion theory introduced by Shannon (1959)

¹Department of Electrical and Computer Engineering, University of Illinois at Urbana-Champaign. Work done as an intern in Google. ²Google, Inc. ³Bytedance, Inc. ⁴Department of Computer Science, University of Washington. Correspondence to: Weihao Gao <wgao9@illinois.edu>, Yu-Han Liu <yuhanliu@google.com>, Chong Wang <mr.chongwang@gmail.com>, Sewoong Oh <sewoong@cs.washington.edu>.

and further developed by Berger (1971). The approach also connects to modeling neural networks as random variables in Mandt et al. (2017), which has many practical usages (Cao et al., 2018).

Our contribution for model compression is twofold: theoretical and practical. We first apply theoretical tools from rate distortion theory to provide a lower bound on the fundamental trade-off between *rate* (number of bits to describe the model) and *distortion* between compressed and original models, and prove the tightness of the lower bound for a linear model. This analysis seamlessly incorporate the structure of the neural network architecture into model compression via backpropagation. Motivated by the theory, we design an improved objective for compression algorithms and show that the improved objective gives optimal pruning and quantization algorithm for one-hidden-layer ReLU neural network, and has better performance in real neural networks as well.

The rest of the paper is organized as follows.

- In Section 2, we briefly review some previous work on model compression.
- In Section 3, we introduce the background of the rate distortion theory for data compression, and formally state the rate distortion theory for model compression.
- In Section 4, we give a lower bound of the rate distortion function, which quantifies the fundamental limit for model compression. We then prove that the lower bound is achievable for linear model.
- In Section 5, motivated by the achievable compressor for linear model, we proposed an improved objective for model compression, which takes consideration of the structure of the neural network. We then prove that the improved objective gives optimal compressor for one-hidden-layer ReLU neural network.
- In Section 6, we demonstrate the empirical performance of the proposed objective on fully-connected neural networks on MNIST dataset and convolutional networks on CIFAR dataset.

2. Related work on model compression

The study of model compression of neural networks appeared as long as neural network was invented. Here we mainly discuss the literature on directly compressing large models, which are more relevant to our work. They usually contain three types of methods — pruning, quantization and matrix factorization.

Pruning methods set unimportant weights to zero to reduce the number of parameters. Early works of model pruning

includes biased weight decay (Hanson & Pratt, 1989), optimal brain damage (LeCun et al., 1990) and optimal brain surgeon (Hassibi & Stork, 1993). Early methods utilize the Hessian matrix of the loss function to prune the weights, however, Hessian matrix is inefficient to compute for modern large neural networks with millions of parameters. More recently, Han et al. (2015b) proposed an iterative pruning and retraining algorithm that works for large neural networks.

Quantization, or weight sharing methods group the weights into clusters and use one value to represent the weights in the same group. This category includes fixed-point quantization by Vanhoucke et al. (2011), vector quantization by Gong et al. (2014), HashedNets by Chen et al. (2015), Hessian-weighted quantization by Choi et al. (2016) and Diameter-regularized Hessian-weighted quantization by Bu et al. (2019).

Matrix factorization assumes the weight matrix in each layer could be factored as a low rank matrix plus a sparse matrix. Hence, storing low rank and sparse matrices is cheaper than storing the whole matrix. This category includes Denton et al. (2014) and Cheng et al. (2015).

There are some recent advanced method beyond pruning, quantization and matrix factorization. Han et al. (2015a) assembles pruning, quantization and Huffman coding to achieve better compression rate. Bayesian methods (Ullrich et al., 2017; Louizos et al., 2017; Federici et al., 2017) are also used to retrain the model such that the model has more space to be compressed. He et al. (2018) uses reinforcement learning to design a compression algorithm.

Despite these aforementioned works for model compression, no one has studied the fundamental limit of model compression, as far as we know. More specifically, in this paper, we focus on the study of theory of model compression for pretrained neural network models and then derive practical compression algorithms given the proposed theory.

3. Rate distortion theory for model compression

In this section, we briefly introduce the rate distortion theory for data compression. Then we extend the theory to compression of model parameters.

3.1. Review of rate distortion theory for data compression

Rate distortion theory, firstly introduced by Shannon (1959) and further developed by Berger (1971), is an important concept in information theory which gives theoretical description of lossy data compression. It addressed the minimum average number of R bits, to transmit a random variable

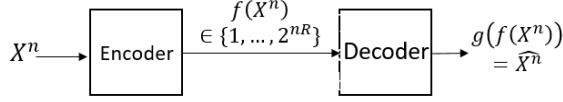


Figure 1. An illustration of encoder and decoder.

such that the receiver can reconstruct the random variable with distortion D .

Precisely, let $X^n = \{X_1, X_2 \dots X_n\} \in \mathcal{X}^n$ be i.i.d. random variables from distribution P_X . An encoder $f_n : \mathcal{X}^n \rightarrow \{1, 2, \dots, 2^{nR}\}$ maps the message X^n into codeword, and a decoder $g_n : \{1, 2, \dots, 2^{nR}\} \rightarrow \mathcal{X}^n$ reconstruct the message by an estimate \hat{X}^n from the codeword. See Figure 1 for an illustration.

A distortion function $d : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}^+$ quantifies the difference of the original and reconstructed message. Distortion between sequence X^n and \hat{X}^n is defined as the average distortion of X_i 's and \hat{X}_i 's. Commonly used distortion function includes Hamming distortion function $d(x, \hat{x}) = \mathbb{I}[x \neq \hat{x}]$ for $\mathcal{X} = \{0, 1\}$ and square distortion function $d(x, \hat{x}) = (x - \hat{x})^2$ for $\mathcal{X} = \mathbb{R}$.

Now we are ready to define the rate-distortion function for data compression.

Definition 1 A rate-distortion pair (R, D) is achievable if there exists a series of (probabilistic) encoder-decoder (f_n, g_n) such that the alphabet of codeword has size 2^{nR} and the expected distortion $\lim_{n \rightarrow \infty} \mathbb{E}[d(X^n, g_n(f_n(X^n)))] \leq D$.

Definition 2 Rate-distortion function $R(D)$ equals to the infimum of rate R such that rate-distortion pair (R, D) is achievable.

The main theorem of rate-distortion theory (Cover & Thomas (2012, Theorem 10.2.1)) states as follows,

Theorem 1 Rate distortion theorem for data compression.

$$R(D) = \min_{P_{\hat{X}|X}: \mathbb{E}[d(X, \hat{X})] \leq D} I(X; \hat{X}). \quad (1)$$

The rate distortion quantifies the fundamental limit of data compression, i.e., at least how many bits are needed to compress the data, given the quality of the reconstructed data. Here is an example for rate-distortion function.

Example 1 If $X \sim \mathcal{N}(0, \sigma^2)$, the rate distortion function is given by

$$R(D) = \begin{cases} \frac{1}{2} \log_2(\sigma^2/D) & \text{if } D \leq \sigma^2 \\ 0 & \text{if } D > \sigma^2 \end{cases}.$$

If the required distortion D is larger than the variance of the Gaussian variable σ^2 , we simply transmit $\hat{X} = 0$; otherwise, we will transmit \hat{X} such that $\hat{X} \sim \mathcal{N}(0, \sigma^2 - D)$, $X - \hat{X} \sim \mathcal{N}(0, D)$ where \hat{X} and $X - \hat{X}$ are independent.

3.2. Rate distortion theory for model compression

Now we extend the rate distortion theory for data compression to model compression. To apply the rate distortion theory to model compression, we view the weights in the model as a multi-dimensional random variable $W \in \mathbb{R}^m$ following distribution P_W . The randomness comes from multiple sources including different distributions of training data, randomness of training data and randomness of training algorithm. The compressor can also be random hence we describe the compressor by a conditional probability $P_{\hat{W}|W}$. Now we define the distortion and rate in model compression, analogously to the data compression scenario.

Distortion. Assume we have a neural network f_w that maps input $x \in \mathbb{R}^{d_x}$ to $f_w(x)$ in output space \mathcal{S} . For regressors, $f_w(x)$ is defined as the output of the neural network on \mathbb{R}^{d_y} . Analogous to the square distortion in data compression, We define the distortion to be the expected ℓ_2 distance between f_w and $f_{\hat{w}}$, i.e.

$$d(w, \hat{w}) \equiv \mathbb{E}_X [\|f_w(X) - f_{\hat{w}}(X)\|_2^2]. \quad (2)$$

For classifiers, $f_w(x)$ is defined as the output probability distribution over C classes on the simplex Δ^{C-1} . We define the distortion to be the expected distance between f_w and $f_{\hat{w}}$, i.e.

$$d(w, \hat{w}) \equiv \mathbb{E}_X [D(f_{\hat{w}}(X) \| f_w(X))]. \quad (3)$$

Here D could be any statistical distance, including KL divergence, Hellinger distance, total variation distance, etc. Such a definition of distortion captures the difference between the original model and the compressed model, averaged over data X , and measures the quality of a compression algorithm.

Rate. In data compression, the rate is defined as the description length of the bits necessary to communicate the compressed data \hat{X} . The compressor outputs \hat{X} from a finite code book \mathcal{X} . The description consists the code word which are the indices of \hat{x} in the code book, and the description of the code book.

In rate distortion theory, we ignore the code book length. Since we are transmitting a sequence of data X^n , the code word has to be transmitted for each X_i but the code book is only transmitted once. In asymptotic setting, the description length of code book can be ignored, and the rate is defined as the description length of the code word.

In model compression, we also define the rate as the code word length, by assuming that an underlying distribution

P_W of the parameters exists and infinitely many models whose parameters are i.i.d. from P_W will be compressed. In practice, we only compress the parameters once so there is no distribution of the parameters. Nevertheless, the rate distortion theory can also provide important intuitions for one-time compression, explained in Section 5.

Now we can define the rate distortion function for model compression. Analogously to Theorem 1, the rate distortion function for model compression is defined as follows,

Definition 3 *Rate distortion function for model compression.*

$$R(D) = \min_{P_{\hat{W}|W} : \mathbb{E}_{W, \hat{W}}[d(W, \hat{W})] \leq D} I(W; \hat{W}). \quad (4)$$

In the following sections we establish a lower bound of the rate-distortion function.

4. Lower bound and achievability for rate distortion function

In this section, we study the lower bound for rate distortion function in Definition 3. We provide a lower bound for the rate distortion function, and prove that this lower bound is achievable for linear regression models.

4.1. Lower bound for linear model

Assume that we are going to compress a linear regression model $f_w(x) = w^T x$. We assume that the mean of data $x \in \mathbb{R}^m$ is zero and the covariance matrix is diagonal, i.e., $\mathbb{E}_X[X_i^2] = \lambda_{x,i} > 0$ and $\mathbb{E}_X[X_i X_j] = 0$ for $i \neq j$. Furthermore, assume that the parameters $W \in \mathbb{R}^m$ are drawn from a Gaussian distribution $\mathcal{N}(0, \Sigma_W)$. The following theorem gives the lower bound of the rate distortion function for the linear regression model.

Theorem 2 *The rate-distortion function of the linear regression model $f_w(x) = w^T x$ is lower bounded by*

$$R(D) \geq \underline{R}(D) = \frac{1}{2} \log \det(\Sigma_W) - \sum_{i=1}^m \frac{1}{2} \log(D_i),$$

where

$$D_i = \begin{cases} \mu / \lambda_{x,i} & \text{if } \mu < \lambda_{x,i} \mathbb{E}_W[W_i^2], \\ \mathbb{E}_W[W_i^2] & \text{if } \mu \geq \lambda_{x,i} \mathbb{E}_W[W_i^2], \end{cases}$$

where μ is chosen that $\sum_{i=1}^m \lambda_{x,i} D_i = D$.

This lower bound gives rise to a “weighted water-filling” approach, which differs from the classical “water-filling” for rate distortion of colored Gaussian source in Cover &

Thomas (2012, Figure 13.7). The details and graphical explanation of the “weighted water-filling” can be found in Appendix A.

4.2. Achievability

We show that, the lower bound give in Theorem 2 is achievable. Precisely, we have the following theorem.

Theorem 3 *There exists a class of probabilistic compressors $P_{\hat{W}^*|W}^{(D)}$ such that $\mathbb{E}_{P_W \circ P_{\hat{W}^*|W}^{(D)}}[d(W, \hat{W}^*)] = D$ and $I(W; \hat{W}^*) = \underline{R}(D)$.*

The optimal compressor is Algorithm 1 in Appendix A. Intuitively, the optimal compressor does the following

- Find the optimal water levels D_i for “weighted water filling”, such that the expected distortion $D = \mathbb{E}_{W, \hat{W}}[d(W, \hat{W})] = \mathbb{E}_{W, \hat{W}}[\hat{W}^T \Sigma_X (W - \hat{W})]$ is minimized given certain rate.
- Add a noise Z_i which is independent of $\hat{W}_i = W_i + Z_i$ and has a variance proportional to the water level. That is possible since W is Gaussian.

We can check that the compressor makes all the inequalities become equality, hence achieve the lower bound. The full proof of the lower bound and achievability can be found in Appendix A.

5. Improved objective for model compression

In the previous sections, we study the rate-distortion theory for model compression. In rate-distortion theory, we assume that there exists a prior distribution P_W on the weights W , and prove the tightness of the lower bound in the asymptotic scenario. However, in practice, we only compress one particular pre-trained model, so there are no prior distribution of W . Nonetheless, we can still learn something important from the achivability of the lower bound, by extracting two “golden rules” from the optimal algorithm for linear regression.

5.1. Two golden rules

Recall that for linear regression model, to achieve the smallest rate given certain distortion (or, equivalently, achieve the smallest distortion given certain rate), the optimal compressor need to do the following: (1) find appropriate “water levels” such that the expected distortion $E_{W, \hat{W}}[d(W, \hat{W})] = \mathbb{E}_{W, \hat{W}, X}[(W^T X - \hat{W}^T X)^2] = \mathbb{E}_{W, \hat{W}}[(W - \hat{W})^T \Sigma_X (W - \hat{W})]$ is minimized. (2) make sure that \hat{W}_i is independent with $W_i - \hat{W}_i$, in other words, $\mathbb{E}_{W, \hat{W}}[\hat{W}^T \Sigma_X (W - \hat{W})] = 0$. Hence, we extract the following two “golden rules”:

1. $\mathbb{E}_{W, \hat{W}}[\hat{W}^T \Sigma_X (W - \hat{W})] = 0$
2. $\mathbb{E}_{W, \hat{W}}[(W - \hat{W})^T \Sigma_X (W - \hat{W})]$ should be minimized, given certain rate.

For practical model compression, we adopt these two “golden rules”, by making the following amendments. First, we discard the expectation over W and \hat{W} since there is only one model to be compressed. Second, the distortion can be written as $d(w, \hat{w}) = (w - \hat{w})^T \Sigma_X (w - \hat{w})$ only for linear models. For non-linear models, the distortion function is complicated, but can be approximated by a simpler formula. For non-linear regression models, we take first order Taylor expansion of the function $f_{\hat{w}}(x) \approx f_w(x) + (\hat{w} - w)^T \nabla_w f_w(x)$, and have

$$\begin{aligned} d(w, \hat{w}) &= \mathbb{E}_X [\|f_w(X) - f_{\hat{w}}(X)\|_2^2] \\ &\approx \mathbb{E}_X [(w - \hat{w})^T \nabla_w f_w(X) (\nabla_w f_w(X))^T (w - \hat{w})] \\ &= (w - \hat{w})^T I_w (w - \hat{w}) \end{aligned}$$

where the “weight importance matrix” defined as

$$I_w = \mathbb{E}_X [\nabla_w f_w(X) (\nabla_w f_w(X))^T], \quad (5)$$

quantifies the relative importance of each weight to the output. For linear regression models, weight importance matrix I_w equals to Σ_X .

For classification models, we will first approximate the KL divergence. Using the Taylor expansion $x \log(x/a) \approx (x - a) + (x - a)^2/(2a)$ for $x/a \approx 1$, the KL divergence $D_{KL}(P||Q)$ for can be approximated by $D_{KL}(P||Q) \approx \sum_i (P_i - Q_i) + (P_i - Q_i)^2/(2P_i) = \sum_i (P_i - Q_i)^2/(2P_i)$, or in vector form $D_{KL}(P||Q) \approx \frac{1}{2}(P - Q)^T \text{diag}[P^{-1}](P - Q)$. Therefore,

$$\begin{aligned} d(w, \hat{w}) &= \mathbb{E}_X [D_{KL}(f_{\hat{w}}(X)||f_w(X))] \\ &\approx \frac{1}{2} \mathbb{E}_X [(f_w(X) - f_{\hat{w}}(X))^T \text{diag}[f_w^{-1}(X)] \\ &\quad (f_w(X) - f_{\hat{w}}(X))] \\ &\approx \frac{1}{2} \mathbb{E}_X [(w - \hat{w})^T (\nabla_w f_w(X)) \text{diag}[f_w^{-1}(X)] \\ &\quad (\nabla_w f_w(X))^T (w - \hat{w})]. \end{aligned}$$

So the weight importance matrix is given by

$$I_w = \mathbb{E}_X [(\nabla_w f_w(X)) \text{diag}[f_w^{-1}(X)] (\nabla_w f_w(X))^T]. \quad (6)$$

This weight importance matrix is also valid for many other statistical distances, including reverse KL divergence, Hellinger distance and Jenson-Shannon distance.

Now we define the two “golden rules” for practical model compression algorithms,

1. $\hat{w}^T I_w (w - \hat{w}) = 0$,
2. $(w - \hat{w})^T I_w (w - \hat{w})$ is minimized given certain constraints.

In the following subsections we will show the optimality of the “golden rules” for a one-hidden-layer neural network, and apply the “golden rules” to derive new objective function for pruning and quantization.

5.2. Optimality for one-hidden-layer ReLU network

We show that if a compressor of a one-hidden-layer ReLU network satisfies the two “golden rules”, it will be the optimal compressor, with respect to mean-square-error. Precisely, consider the one-hidden layer ReLU neural network $f_w(x) = \text{ReLU}(w^T x)$, where the distribution of input $x \in \mathbb{R}^m$ is $\mathcal{N}(0, \Sigma_X)$. Furthermore, we assume that the covariance matrix $\Sigma_X = \text{diag}[\lambda_{x,1}, \dots, \lambda_{x,m}]$ is diagonal and $\lambda_{x,i} > 0$ for all i . We have the following theorem.

Theorem 4 *If compressed weight \hat{w}^* satisfies $\hat{w}^* I_w (\hat{w}^* - w) = 0$ and*

$$\hat{w}^* = \arg \min_{\hat{w} \in \hat{\mathcal{W}}} (w - \hat{w})^T I_w (w - \hat{w}),$$

where $\hat{\mathcal{W}}$ is some class of compressors, then

$$\hat{w}^* = \arg \min_{\hat{w} \in \hat{\mathcal{W}}} \mathbb{E}_X [(f_w(X) - f_{\hat{w}}(X))^2].$$

The proof uses the techniques of Hermite polynomials and Fourier analysis on Gaussian spaces, inspired by Ge et al. (2017). The full proof can be found in Appendix B. Generalizing this result to other activation functions and deeper neural networks are possible future directions.

Here $\hat{\mathcal{W}}$ denotes a class of compressors, with some constraints. For example, $\hat{\mathcal{W}}$ could be the class of pruning algorithms where no more than 50% weights are pruned, or $\hat{\mathcal{W}}$ could be the class of quantization algorithm where each weight is quantized to 4 bits. Theoretically, it is not guaranteed that the two “golden rules” can be satisfied simultaneously for every $\hat{\mathcal{W}}$, but in the following subsection we show that they can be satisfied simultaneously for two of the most commonly used class of compressors — pruning and quantization. Hence, minimizing the objective $(w - \hat{w})^T I_w (w - \hat{w})$ will be optimal for pruning and quantization.

5.3. Improved objective for pruning and quantization

Pruning and quantization are two most basic and useful building blocks of modern model compression algorithms, For example, DeepCompress (Han et al., 2015a) iteratively prune, retrain and quantize the neural network and achieve state-of-the-art performances on large neural networks.

In pruning algorithms, we choose a subset $S \in [m]$ and set $\hat{w}_i = 0$ for all $i \in S$ and $\hat{w}_i = w_i$ for $i \notin S$. The compression ratio is evaluated by the proportion of unpruned weights $r = (m - |S|)/m$. Since either \hat{w}_i or $w_i - \hat{w}_i$ is zero, so the first “golden rule” is automatically satisfied, so we have the following corollary.

Corollary 1 For any fixed r , let

$$\hat{w}_r^* = \arg \min_{S: \frac{d-|S|}{d}=r} (w - \hat{w})^T I_w (w - \hat{w}),$$

Then

$$\hat{w}_r^* = \arg \min_{S: \frac{d-|S|}{d}=r} \mathbb{E}_X [(f_w(X) - f_{\hat{w}}(X))^2].$$

In quantization algorithms, we cluster the weights into k centroids $\{c_1, \dots, c_k\}$. The algorithm optimize the centroids as long as the assignments of each weight $A_i \in [k]$. The final compressed weight is given by $\hat{w}_i = c_{A_i}$. Usually k -means algorithm are utilized to minimize the centroids and assignments alternatively. The compression ratio of quantization algorithm is given by

$$r = \frac{mb}{m \sum_{j=1}^k \frac{m_j}{m} [\log_2 \frac{m}{m_j}] + kb},$$

where m is the number of weights and b is the number of bits to represent one weight before quantization (usually 32). By using Huffman coding, the average number of bits for each weight is given by $\sum_{j=1}^k (m_j/m) [\log_2 (m/m_j)]$, where m_j is the number of weights assigned to the j -th cluster. The definition of compression ratio of pruning and quantization is consistent since both of them equals to the number of bits representing compressed model parameters divided by the number of bits representing original model parameters.

If we can find the optimal quantization algorithm with respect to $(w - \hat{w})^T I_w (w - \hat{w})$, then each centroids c_j should be optimal, i.e.

$$\begin{aligned} 0 &= \frac{\partial}{\partial c_j} (w - \hat{w})^T I_w (w - \hat{w}) \\ &= -2 \left(\sum_{i: A_i=j} e_i^T \right) I_w (w - \hat{w}) \end{aligned}$$

where e_i is the i -th standard basis. Therefore, we have

$$\begin{aligned} \hat{w} I_w (\hat{w} - w) &= \left(\sum_{j=1}^k c_j \left(\sum_{i: A_i=j} e_i \right) \right)^T I_w (w - \hat{w}) \\ &= \sum_{j=1}^k c_j \left(\left(\sum_{i: A_i=j} e_i^T \right) I_w (w - \hat{w}) \right) = 0. \end{aligned}$$

Hence the first “golden rule” is satisfied if the second “golden rule” is satisfied. So we have

Corollary 2 For any fixed number of centroids k , let

$$\hat{w}_k^* = \arg \min_{\{c_1, \dots, c_k\}, A \in [k]^m} (w - \hat{w})^T I_w (w - \hat{w}),$$

then

$$\hat{w}_k^* = \arg \min_{\{c_1, \dots, c_k\}, A \in [k]^m} \mathbb{E}_X [(f_w(X) - f_{\hat{w}}(X))^2].$$

As corollaries of Theorem 4, we proposed to use $(w - \hat{w})^T I_w (w - \hat{w})$ as the objective for pruning and quantization algorithms, which can achieve the minimum MSE for one-hidden-layer ReLU neural network.

6. Experiments

In the previous section, we proved that a pruning or quantization algorithm that minimizes the objective $(w - \hat{w})^T I_w (w - \hat{w})$ also minimizes the MSE loss for one-hidden-layer ReLU neural network. In this section, we show that this objective can also improve pruning and quantization algorithm for larger neural networks on real data.¹

We test the objectives on the following neural network and datasets.²

1. 3-layer fully connected neural network on MNIST.
2. Convolutional neural network with 5 convolutional layers and 3 fully connected layers on CIFAR 10 and CIFAR 100.

In Section 6.1, we use the weight importance matrix for classification in Eq. (6), which is derived by approximating the distortion of KL-divergence. This weight importance matrix does not depend on the training labels, so the induced pruning/quantization algorithms is called “unsupervised compression”. Furthermore, if the training labels are available, we treat the loss function $\mathcal{L}_w(X, Y) : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^+$ as the function to be compressed, and derive several pruning/quantization objectives. The induced pruning/quantization methods are called “supervised compression” and are studied in Section 6.2.

6.1. Unsupervised Compression Experiments

Recall that for classification problems, the weight importance matrix is defined as

$$I_w = \mathbb{E}_X [\nabla_w f_w(X) \text{diag}[f_w^{-1}(X)] (\nabla_w f_w(X))^T].$$

¹We leave combinations of pruning, model retraining and quantization like Han et al. (2015a) as future work.

²We load the pretrained models from <https://github.com/aaron-xichen/pytorch-playground>.

For computational simplicity, we drop the off-diagonal terms of I_w , and simplify the objective to $\sum_{i=1}^m \mathbb{E}_X \left[\frac{(\nabla_{w_i} f_w(X))^2}{f_w(X)} \right] (w_i - \hat{w}_i)^2$. To minimize the proposed objective, a pruning algorithm just prune the weights with smaller $\mathbb{E}_X \left[\frac{(\nabla_{w_i} f_w(X))^2}{f_w(X)} \right] w_i^2$ greedily. A quantization algorithm uses the weighted k -means algorithm (Choi et al., 2016) to find the optimal centroids and assignments. We compare the proposed objective with the baseline objective $\sum_{i=1}^m (w_i - \hat{w}_i)^2$, which were used as building blocks in DeepCompress (Han et al., 2015a). We compare the objectives in Table 6.1.

Name	Minimizing objective
Baseline	$\sum_{i=1}^m (w_i - \hat{w}_i)^2$
Proposed	$\sum_{i=1}^m \mathbb{E}_X \left[\frac{(\nabla_{w_i} f_w(X))^2}{f_w(X)} \right] (w_i - \hat{w}_i)^2$

Table 1. Comparison of unsupervised compression objectives.

For pruning experiment, we choose the same compression rate for every convolutional layer and fully-connected layer, and plot the test accuracy and test cross-entropy loss against compression rate. For quantization experiment, we choose the same number of clusters for every convolutional and fully-connected layer. Also we plot the test accuracy and test cross-entropy loss against compression rate. To reduce the variance of estimating the weight importance matrix I_w , we use the *temperature scaling* method introduced by Guo et al. (2017) to improve model calibration.

We show that results of pruning experiment in Figure 2, and the results of quantization experiment in Figure 3. We can see that the proposed objective gives better validation cross-entropy loss than the baseline, for every different compression ratios. The proposed objective also gives better validation accuracy in most scenarios. Occasionally the proposed objective can not improve the accuracy (top left of Figure 2). We conjecture that the reason is the ill-calibration of the original model. We relegate the results for CIFAR100 in Appendix C.

6.2. Supervised Compression Experiments

In the previous experiment, we only use the training data to compute the weight importance matrix. But if we can use the training label as well, we can further improve the performance of pruning and quantization algorithms. If the training label is available, we can view the cross-entropy loss function $\mathcal{L}(f_w(x), y) = \mathcal{L}_w(x, y)$ as a function from $\mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^+$, and define the distortion function as

$$d(w, \hat{w}) = \mathbb{E}_{X, Y} [(\mathcal{L}_w(X, Y) - \mathcal{L}_{\hat{w}}(X, Y))^2].$$

Taking first order approximation of the loss function gives the supervised weight importance matrix,

$$I_w = \mathbb{E} [\nabla_w \mathcal{L}_w(X, Y) (\nabla_w \mathcal{L}_w(X, Y))^T].$$

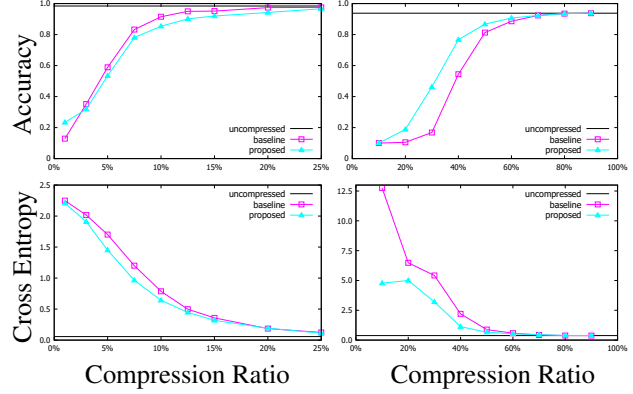


Figure 2. Result for unsupervised pruning experiment. Left: fully-connected NN on MNIST (Top: test accuracy, Bottom: test cross entropy). Right: ConvNN on CIFAR10 (Top: test accuracy, Bottom: test cross entropy).

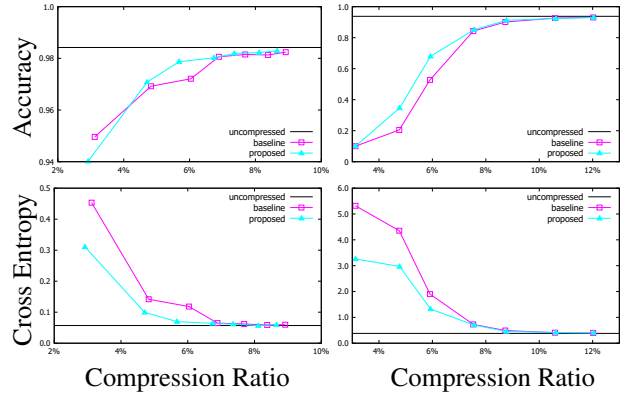


Figure 3. Result for unsupervised quantization experiment. Left: fully-connected NN on MNIST (Top: test accuracy, Bottom: test cross entropy). Right: ConvNN on CIFAR10 (Top: test accuracy, Bottom: test cross entropy).

We write \mathbb{E} instead of $\mathbb{E}_{X, Y}$ for simplicity. Similarly, we drop the off-diagonal terms for ease of computation, and simplify the objective to $\sum_{i=1}^m \mathbb{E} [(\nabla_{w_i} \mathcal{L}_w(X, Y))^2] (w_i - \hat{w}_i)^2$, which is called gradient-based objective. Note that for well-trained model, the expected value of gradient $\mathbb{E}[\nabla_w \mathcal{L}_w(X, Y)]$ is closed to zero, but the second moment of the gradient $\mathbb{E}[\nabla_w \mathcal{L}_w(X, Y) (\nabla_w \mathcal{L}_w(X, Y))^T]$ could be large. We compare this objective with the baseline objective $\sum_{i=1}^m (w_i - \hat{w}_i)^2$. We also compare with the hessian-based objective $\sum_{i=1}^m \mathbb{E}[\nabla_{w_i}^2 \mathcal{L}_w(X, Y)] (w_i - \hat{w}_i)^2$, which is used in (LeCun et al., 1990) and (Hassibi & Stork, 1993) for network pruning and (Choi et al., 2016) for network quantization. To estimate the diagonal entries of the Hessian matrix of the loss function with respect to the model parameters, we implemented Curvature Propagation (Martens et al., 2012) treating each layer and activation as a node. The running time is proportional to the running time of the usual gradient back-propagation by a factor that does not depend on the

size of the model. Manually optimizing the local Hessian calculation at each node reduces memory usage and allows us to use larger batch size and larger number of samples for more accurate estimates.

Furthermore, if we take second order approximation of the loss function, and drop the off-diagonal terms of the squared gradient matrix and squared hessian tensor, we have the following approximation

$$\begin{aligned} d(w, \hat{w}) &= \mathbb{E}[(\mathcal{L}_w(X, Y) - \mathcal{L}_{\hat{w}}(X, Y))^2] \\ &\approx \mathbb{E}[(\nabla_w \mathcal{L}_w(X, Y))^T (w - \hat{w}) \\ &\quad + \frac{1}{2} (w - \hat{w})^T \nabla_w^2 \mathcal{L}_w(X, Y) (w - \hat{w})]^2 \\ &\approx \sum_{i=1}^m \mathbb{E}[(\nabla_{w_i} \mathcal{L}_w(X, Y))^2] (w_i - \hat{w}_i)^2 \\ &\quad + \frac{1}{4} \sum_{i=1}^m \mathbb{E}[(\nabla_{w_i}^2 \mathcal{L}_w(X, Y))^2] (w_i - \hat{w}_i)^4, \end{aligned}$$

which is called gradient+hessian based objective. For pruning algorithm, we can prune the weights with smaller $\mathbb{E}[(\nabla_{w_i} \mathcal{L}_w(X, Y))^2] w_i^2 + \frac{1}{4} \mathbb{E}[(\nabla_{w_i}^2 \mathcal{L}_w(X, Y))^2] w_i^4$ greedily. For quantization algorithm, we use an alternative minimization algorithm in Appendix C to find the minimum. We conclude the different supervised objectives in Table 6.2.

Name	Minimizing objective
Baseline	$\sum_{i=1}^m (w_i - \hat{w}_i)^2$
Gradient	$\sum_{i=1}^m \mathbb{E}[(\nabla_{w_i} \mathcal{L}_w(X, Y))^2] (w_i - \hat{w}_i)^2$
Hessian	$\sum_{i=1}^m \mathbb{E}[\nabla_{w_i}^2 \mathcal{L}_w(X, Y)] (w_i - \hat{w}_i)^2$
Gradient + Hessian	$\sum_{i=1}^m \mathbb{E}[(\nabla_{w_i} \mathcal{L}_w(X, Y))^2] (w_i - \hat{w}_i)^2 + \frac{1}{4} \sum_{i=1}^m \mathbb{E}[(\nabla_{w_i}^2 \mathcal{L}_w(X, Y))^2] (w_i - \hat{w}_i)^4$

Table 2. Comparison of supervised compression objectives.

We show that results of pruning experiment in Figure 4, and quantization experiment in Figure 5. Generally, the gradient objective and hessian objective both give better performance than baseline objective, while gradient objective is slightly than hessian objective at some points. Gradient + hessian objective gives the best overall performance. We relegate the results for CIFAR100 in Appendix C.

Remark. Here we define the supervised distortion function as $d(w, \hat{w}) = \mathbb{E}_{X, Y} [(\mathcal{L}_w(X, Y) - \mathcal{L}_{\hat{w}}(X, Y))^2]$, analogously to the distortion of regression. However, since the goal of classification is to minimize the loss function, the following definition of distortion function $\tilde{d}(w, \hat{w}) = \mathbb{E}_{X, Y} [\mathcal{L}_{\hat{w}}(X, Y) - \mathcal{L}_w(X, Y)]$ is also valid and has been adopted in LeCun et al. (1990) and Choi et al. (2016). The main difference is — $d(w, \hat{w})$ focus on the quality of *compression algorithm*, i.e., how similar is the compressed model compared to uncompressed model, whereas $\tilde{d}(w, \hat{w})$ focus on the quality of *compressed model*, i.e. how good

is the compressed model. So $d(w, \hat{w})$ is a better criteria for the compression algorithm. Additionally, by taking second order approximation of $d(w, \hat{w})$, we have gradient+hessian objective, which shows better empirical performance than hessian objective, derived by taking second order approximation of $\tilde{d}(w, \hat{w})$.

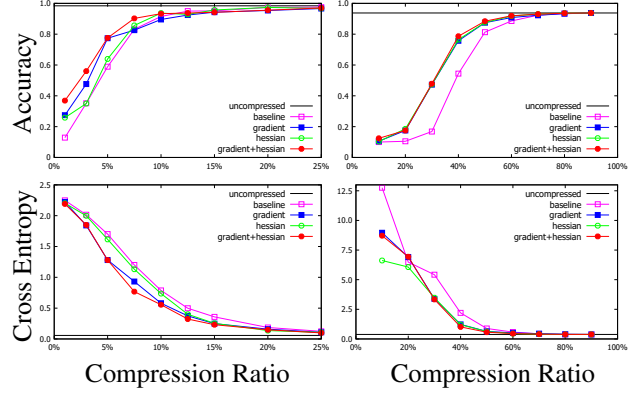


Figure 4. Result for supervised pruning experiment. Left: fully-connected NN on MNIST (Top: test accuracy, Bottom: test cross entropy). Right: ConvNN on CIFAR10 (Top: test accuracy, Bottom: test cross entropy).

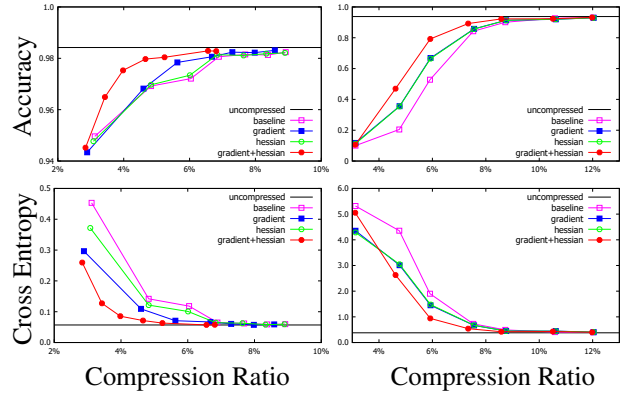


Figure 5. Result for supervised quantization experiment. Left: fully-connected NN on MNIST (Top: test accuracy, Bottom: test cross entropy). Right: ConvNN on CIFAR10 (Top: test accuracy, Bottom: test cross entropy).

7. Conclusion

In this paper, we investigate the fundamental limit of neural network model compression algorithms. We prove a lower bound for the rate distortion function for model compression, and prove its achievability for linear model. Motivated by the rate distortion function, we propose the weight importance matrix, and show that for one-hidden-layer ReLU network, pruning and quantization that minimizes the proposed objective is optimal. We also show the superiority of proposed objective in real neural networks.

Acknowledgement

The authors thank Denny Zhou for initial comments and helpful discussions. This work is partially supported by Google and NSF award 1815535.

References

- Berger, T. Rate distortion theory: A mathematical basis for data compression. 1971.
- Bu, Y., Gao, W., Zou, S., and Veeravalli, V. V. Information-theoretic understanding of population risk improvement with model compression. *arXiv preprint arXiv:1901.09421*, 2019.
- Cao, W., Wang, X., Ming, Z., and Gao, J. A review on neural networks with random weights. *Neurocomputing*, 275:278–287, 2018.
- Chen, W., Wilson, J., Tyree, S., Weinberger, K., and Chen, Y. Compressing neural networks with the hashing trick. In *International Conference on Machine Learning*, pp. 2285–2294, 2015.
- Cheng, Y., Yu, F. X., Feris, R. S., Kumar, S., Choudhary, A., and Chang, S.-F. An exploration of parameter redundancy in deep networks with circulant projections. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 2857–2865, 2015.
- Choi, Y., El-Khamy, M., and Lee, J. Towards the limit of network quantization. *arXiv preprint arXiv:1612.01543*, 2016.
- Coates, A., Huval, B., Wang, T., Wu, D., Catanzaro, B., and Andrew, N. Deep learning with cots hpc systems. In *International Conference on Machine Learning*, pp. 1337–1345, 2013.
- Cover, T. M. and Thomas, J. A. *Elements of information theory*. John Wiley & Sons, 2012.
- Denton, E. L., Zaremba, W., Bruna, J., LeCun, Y., and Fergus, R. Exploiting linear structure within convolutional networks for efficient evaluation. In *Advances in neural information processing systems*, pp. 1269–1277, 2014.
- Federici, M., Ullrich, K., and Welling, M. Improved bayesian compression. *arXiv preprint arXiv:1711.06494*, 2017.
- Ge, R., Lee, J. D., and Ma, T. Learning one-hidden-layer neural networks with landscape design. *arXiv preprint arXiv:1711.00501*, 2017.
- Gong, Y., Liu, L., Yang, M., and Bourdev, L. Compressing deep convolutional networks using vector quantization. *arXiv preprint arXiv:1412.6115*, 2014.
- Guo, C., Pleiss, G., Sun, Y., and Weinberger, K. Q. On calibration of modern neural networks. *arXiv preprint arXiv:1706.04599*, 2017.
- Han, S., Mao, H., and Dally, W. J. Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding. *arXiv preprint arXiv:1510.00149*, 2015a.
- Han, S., Pool, J., Tran, J., and Dally, W. Learning both weights and connections for efficient neural network. In *Advances in neural information processing systems*, pp. 1135–1143, 2015b.
- Hanson, S. J. and Pratt, L. Y. Comparing biases for minimal network construction with back-propagation. In *Advances in neural information processing systems*, pp. 177–185, 1989.
- Hassibi, B. and Stork, D. G. Second order derivatives for network pruning: Optimal brain surgeon. In *Advances in neural information processing systems*, pp. 164–171, 1993.
- He, Y., Lin, J., Liu, Z., Wang, H., Li, L.-J., and Han, S. Amc: Automl for model compression and acceleration on mobile devices. 2018.
- Howard, A. G., Zhu, M., Chen, B., Kalenichenko, D., Wang, W., Weyand, T., Andreetto, M., and Adam, H. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.
- Iandola, F. N., Han, S., Moskewicz, M. W., Ashraf, K., Dally, W. J., and Keutzer, K. Squeezenet: Alexnet-level accuracy with 50x fewer parameters and 0.5 mb model size. *arXiv preprint arXiv:1602.07360*, 2016.
- Jiao, J., Gao, W., and Han, Y. The nearest neighbor information estimator is adaptively near minimax rate-optimal. *arXiv preprint arXiv:1711.08824*, 2017.
- Krizhevsky, A., Sutskever, I., and Hinton, G. E. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pp. 1097–1105, 2012.
- LeCun, Y., Denker, J. S., and Solla, S. A. Optimal brain damage. In *Advances in neural information processing systems*, pp. 598–605, 1990.
- LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Louizos, C., Ullrich, K., and Welling, M. Bayesian compression for deep learning. In *Advances in Neural Information Processing Systems*, pp. 3288–3298, 2017.

- Mandt, S., Hoffman, M. D., and Blei, D. M. Stochastic gradient descent as approximate bayesian inference. *The Journal of Machine Learning Research*, 18(1):4873–4907, 2017.
- Martens, J., Sutskever, I., and Swersky, K. Estimating the hessian by back-propagating curvature. *arXiv preprint arXiv:1206.6464*, 2012.
- McDonald, R. and Schultheiss, P. Information rates of gaussian signals under criteria constraining the error spectrum. *Proceedings of the IEEE*, 52(4):415–416, 1964.
- Shannon, C. E. Coding theorems for a discrete source with a fidelity criterion. *IRE Nat. Conv. Rec*, 4(142-163):1, 1959.
- Silver, D., Schrittwieser, J., Simonyan, K., Antonoglou, I., Huang, A., Guez, A., Hubert, T., Baker, L., Lai, M., Bolton, A., et al. Mastering the game of go without human knowledge. *Nature*, 550(7676):354, 2017.
- Simonyan, K. and Zisserman, A. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- Ullrich, K., Meeds, E., and Welling, M. Soft weight-sharing for neural network compression. *arXiv preprint arXiv:1702.04008*, 2017.
- Vanhoucke, V., Senior, A., and Mao, M. Z. Improving the speed of neural networks on cpus. In *Proc. Deep Learning and Unsupervised Feature Learning NIPS Workshop*, volume 1, pp. 4. Citeseer, 2011.
- Wu, Y., Schuster, M., Chen, Z., Le, Q. V., Norouzi, M., Macherey, W., Krikun, M., Cao, Y., Gao, Q., Macherey, K., et al. Google’s neural machine translation system: Bridging the gap between human and machine translation. *arXiv preprint arXiv:1609.08144*, 2016.