Model compression as constrained optimization, with application to neural nets



Miguel Á. Carreira-Perpiñán

Electrical Engineering and Computer Science University of California, Merced

http://faculty.ucmerced.edu/mcarreira-perpinan

work with Yerlan Idelbayev and Arman Zharmagambetov

Outline

- Machine learning and the IoT: the need for model compression
- Direct compression via standard signal processing algorithms
- Model compression as constrained optimization
- The "learning-compression" (LC) algorithm
- Several examples of deep neural net compression:
 - 1. Low-rank
 - 2. Quantization
 - 3. Pruning
 - 4. Learning the ranks
 - 5. Additive combination of compressions
- Beyond compression:
 - Model compression as structure learning
 - Model compression as model compilation

Model compression

Model compression is an interesting problem:

- It involves several disciplines:
 - machine learning
 - optimization
 - signal processing
 - computer science
- It has real practical application.

We will focus mostly on compression of deep neural nets, but compressing other models is also important.

The rise of machine learning

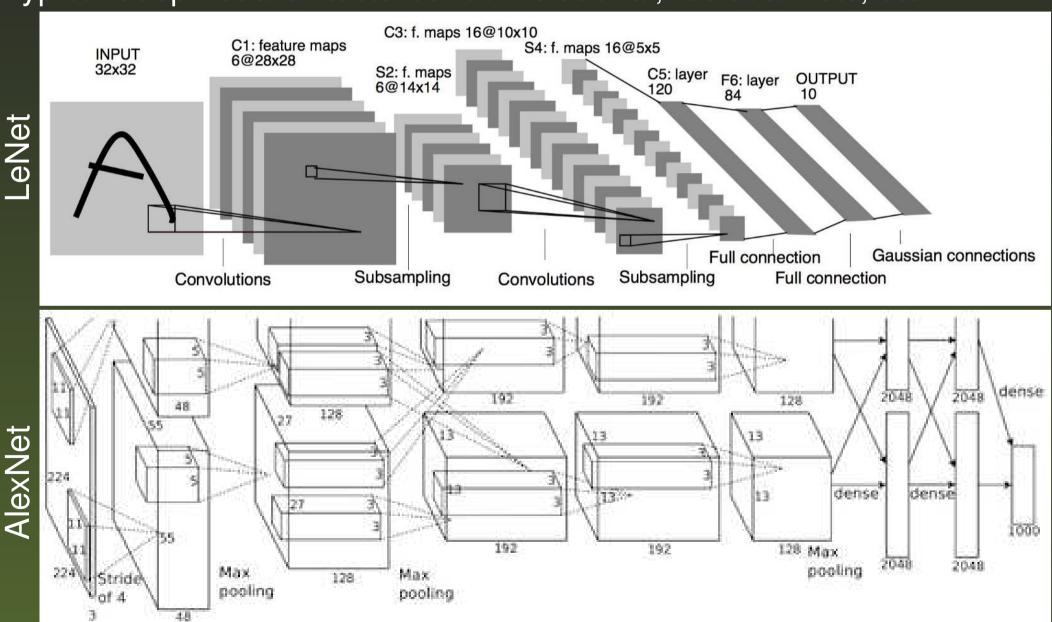
Machine learning, and deep learning in particular, has recently become practical and widespread, significantly pushing the state of the art in many applications:

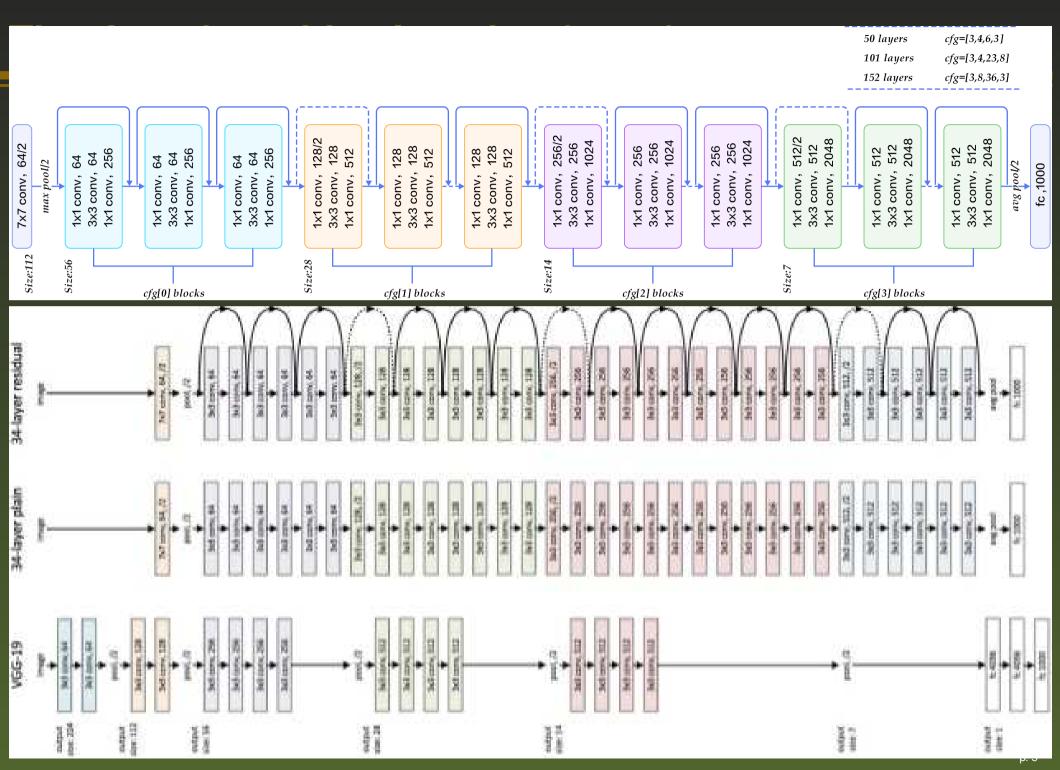
- Computer vision: object recognition, tracking, etc.
- Image and signal processing: super-resolution, style transfer, etc.
- Speech processing: speech recognition, speech synthesis, etc.
- NLP: dialog systems, image captioning, etc.
- Information retrieval in large databases of image, audio, text, etc.

Reasons:

- much data available publicly or privately
- sophisticated models & algorithms in the statistics & ML literature
- computing power provided by GPUs
- deep learning toolboxes accessible to non-experts. Tensorflow, PyTorch...

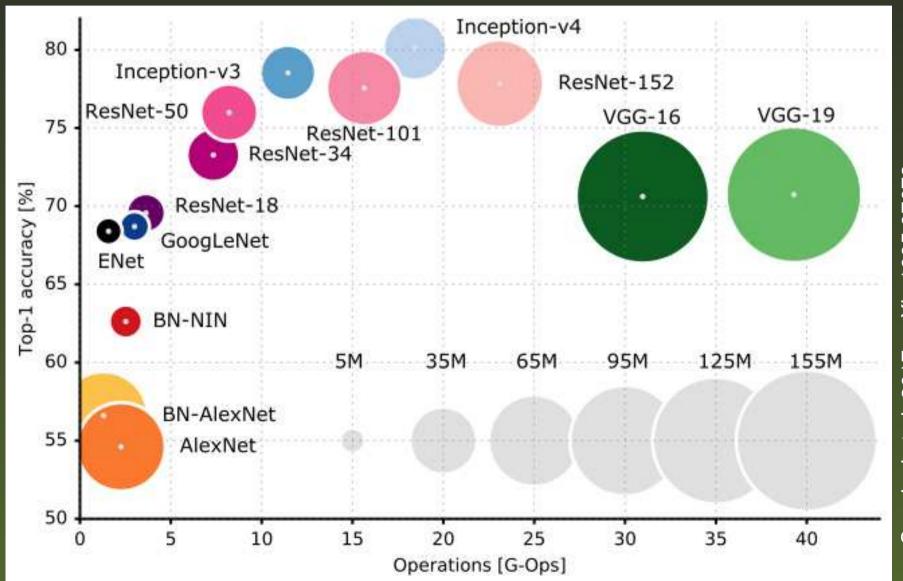
Typical deep net architectures: convolutional, residual nets, etc.





The rise of machine learning (cont.)

Downside: very good classification accuracy needs training large, deep neural nets on large datasets, with several GPUs over several days.



Canziani et al. 2017, arXiv:1605.07678

The rise of the Internet-of-Things (IoT)

Internet-of-Things (IoT) and other mobile devices have become pervasive in daily life and this trend is going to continue growing:

- smart phones
- smart cameras
- mobile robots
- wearables
- personal assistants, tablets
- sensors in medicine, ecology, smart buildings...

Their functionality can be significantly enhanced with intelligent, adaptive software for tasks in computer vision, speech, NLP, etc.

The rise of the Internet-of-Things (cont.)

Downside: IoT devices have stringent computing limitations:

- CPU speed
- memory
- bandwidth and connectivity
- battery life
- energy consumption
- other processes running in the device beyond the neural net.

Actual specs vary widely depending on the device and intended application, but are generally very limited compared to a workstation.

This places a limit in, say, the size or the energy consumption of the deep nets that can be deployed, and motivates the problem of model compression: taking a well-trained, reference model and compressing it into one that fits in the device (in memory, energy consumption, etc.).

Model compression and statistics

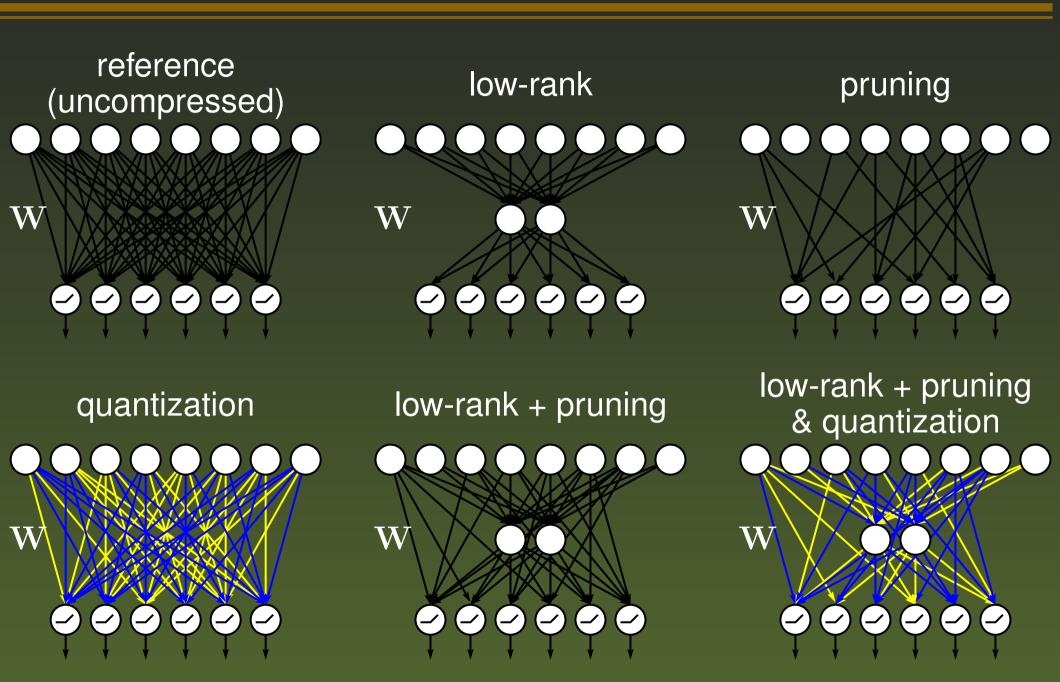
Although the model compression problem is motivated by a practical, engineering need, it is related to deeper questions about model selection, structure learning, effective number of parameters, regularization, etc.

Some forms of compression have a long history in statistics & ML, usually applied to linear models of the form y = Wx:

- * pruning = sparse learning (certain parameters w_{ij} are zero) Lasso, elastic net, etc.
- low-rank compression = parameter matrix W is low-rank reduced-rank regression (RRR)
- * quantization = learnable parameter sharing (certain groups of parameters w_{ij} take the same value) convolutional neural nets, tied models

In general, compression is a sophisticated form of regularization. Compressed models are less complex and may generalize better.

Model compression: examples



Neural net compression

Why not train a small model in the first place? Certainly possible, but:

- The desired size depends on the target device, so the training (which is not easy) must be repeated for each new device type.
- Training well a model without restrictions (the reference model) gives us a baseline: the best possible model for the available training data.
- For reasons not yet well understood, in practice it seems easier to achieve state-of-the-art performance by using a larger-than-needed model rather than a model of the right size: over-parameterization.

It seems that, for deep nets: 1) over-parameterized nets (more parameters than points, even without regularization) don't overfit; 2) training an over-parameterized net with SGD tends to minima with zero training error that generalize well; 3) the optimization is easier if the space has excess dimensions that allow us to take shortcuts towards a minimizer.

This leads to reference nets with many layers and many units per layer that can be significantly compressed.

A common approach at present is to train a good reference model and then compress it as needed for a particular device.

Neural net compression: related work

Many papers since around 2015, although neural net compression was already studied in the 1980–90s. Some common approaches:

- Direct compression: train a reference net, then compress its weights using standard signal compression techniques.

 Quantization (binarization, low precision), pruning weights/neurons, low-rank, etc.

 Simple and fast, but it ignores the loss function, so the accuracy deteriorates a lot for high compression rates.
- Embedding the compression mechanism in the SGD training.
 - remove neuron/weight values with low magnitude on the fly
 - binarize or round weight values in the backpropagation forward pass
 - Often heuristic, with no convergence guarantees.
- Other algorithms proposed for specific compression techniques.
- Multiple compression techniques can be combined (sequentially).

Over-parameterized nets can be considerably compressed even with simple techniques, but we seek optimal compression.

How to train the reference model

In machine learning, we typically learn by solving an optimization problem. We define and objective function (loss) $L(\mathbf{w})$, possibly s.t. constraints, given by the task, model and training set, e.g.:

- lacktriangle classification: cross-entropy $L(\mathbf{w}) = -\sum_{n=1}^{N} \sum_{k=1}^{K} y_{nk} f_k(\mathbf{x}_n; \mathbf{w})$
- regression: least-squares error $L(\mathbf{w}) = \sum_{n=1}^{N} \|\mathbf{y}_n \mathbf{f}(\mathbf{x}_n; \mathbf{w})\|^2$
- lacktriangle maximum likelihood $L(\mathbf{w}) = -\sum_{n=1}^{N} \log p(\mathbf{x}_n; \mathbf{w})$
- otc.

where $\{(\mathbf{x}_n, \mathbf{y}_n)\}_{n=1}^N$ is the labeled training set and $\mathbf{f}(\mathbf{x}; \mathbf{w})$ is the predictive function to be learnt (say, a neural net), with parameters \mathbf{w} .

We then find a (local) minimizer:

$$\min_{\mathbf{w}} L(\mathbf{w})$$

with a suitable optimization algorithm.

How to train the reference model (cont.)

For deep nets in particular:

- \bullet w = weights & biases in all layers of the net; the loss is nonconvex.
- Minimizing $L(\mathbf{w}) = \sum_{n=1}^N L_n(\mathbf{w})$ is typically achieved by stochastic gradient descent (SGD): $\mathbf{w}_{k+1} = \mathbf{w}_k \eta_k \sum_{n \in \mathcal{B}_k} \nabla L_n(\mathbf{w}_k), \ k = 0, 1 \dots$ Gradient step on a minibatch $\mathcal{B}_k \subset \{1, \dots, N\}$, typically with momentum. For convergence, the step sizes $\{\eta_k\}$ must satisfy the Robbins-Monro conditions (tend to zero at a certain rate). Practically, the user must carefully select the hyperparameters to ensure
 - Practically, the user must carefully select the hyperparameters to ensure relatively fast convergence (step sizes, momentum rate, minibatch size, possibly others), and stop training when a low enough loss is achieved. This is tricky. Other algorithms exist but generally SGD with well-tuned hyperparameters is best.
- GPUs are currently the best computing platform.
 Other platforms (e.g. CPUs) have not been as successful.
 Training times of days or even weeks for large problems.
- Deep learning toolboxes (with automatic differentiation) have made this process accessible to non-experts. Tensorflow, PyTorch, Theano, MXNet, Caffe...

Direct compression (via signal processing algorithms)

- Compression has been long studied in signal processing.
- Direct compression: simply take the reference model parameters and compress them using off-the-shelf compression algorithms, just as is done with compression of text, image or other signals:
 - Lossless: LZW, run-length / Huffman / arithmetic codes, etc. Generally won't achieve much compression.
 - ◆ Lossy: DCT (JPEG), wavelet, etc. (see below).
 Generally will achieve much more compression, depending on the error tolerated.
- Simple and practical, because we use existing algorithms with well-developed code:
 - variable-length codes: run Huffman's algorithm
 - pruning (zeroing): zero all weights smaller than a threshold
 - lacktriangle binarization into $\{-1,+1\}$: replace each weight w_i with $\mathrm{sgn}\,(w_i)$
 - lacktriangle quantization: run k-means on the weights $\mathbf{w} = \{w_1, \dots, w_P\} \subset \mathbb{R}$
 - low-rank: compute the singular value decomposition (SVD).

Direct compression (cont.)

The resulting compressed model \mathbf{w}^{DC} is optimal wrt the uncompressed, reference model $\overline{\mathbf{w}}$...

That is, among all possible compressed models, \mathbf{w}^{DC} has the lowest compression error (e.g. quadratic distortion $\|\mathbf{w}^{DC} - \overline{\mathbf{w}}\|^2$ for k-means or SVD).

... but it is not optimal wrt the loss $L(\mathbf{w})$.

That is, among all possible compressed models, \mathbf{w}^{DC} does not generally have the smallest loss.

This introduces a radical difference with the traditional view of signal compression. We still want to compress, but to minimize the loss, not the compression error.

How well does direct compression work?

- Quite well as long as we don't compress much, but the loss blows up as we compress more. With deep nets, surprisingly high compression rates are often possible because they are vastly over-parameterized in practice.
- In practice, we want to compress as much as possible, so we need an optimal way to do this that achieves the lowest loss for a given compression rate.

Traditional compression vs machine learning

How do we combine the following notions?

1. The traditional signal compression view:

- extensively studied
- many existing forms of compression
- usually carried out with an algorithm that is (near-)optimal wrt the compression error quantization by k-means, low-rank by SVD, etc.

2. The traditional machine learning view:

- * a model with a given architecture, e.g. a deep net
- a loss defined given a training set (e.g. cross-entropy) that we need to minimize over the model parameters
- many optimization algorithms specific for each loss and model gradient descent, SGD, variations of Newton's method, (L-)BFGS, EM algorithm...

Outline

- Machine learning and the IoT: the need for model compression
- Direct compression via standard signal processing algorithms
- Model compression as constrained optimization
- The "learning-compression" (LC) algorithm
- Several examples of deep neural net compression:
 - 1. Low-rank
 - 2. Quantization
 - 3. Pruning
 - 4. Learning the ranks
 - 5. Additive combination of compressions
- Beyond compression:
 - Model compression as structure learning
 - Model compression as model compilation

Desiderata for a good solution to the problem

Firstly, we need a precise, general mathematical definition of the problem of model compression that is amenable to numerical optimization techniques. Intuitively, we want to achieve the lowest loss possible for a given compression technique and compression rate.

Then, we would like an algorithm that:

- ❖ is generic, applicable to many losses and compression techniques so the user can easily try different types of compression on a given loss function without having to use a specialized algorithm for each case
- has optimality guarantees inasmuch as possible
- is easy to integrate in existing deep learning / ML toolboxes
- is efficient in training
- works on the native discrete space and can make use of effective discrete optimization solvers rather than relaxing/applying soft priors & truncating, for compression involving discrete variables: quantization, pruning, learning ranks.

An illustrative example: low-rank compression

Consider a linear regression problem with weight matrix $W_{d\times D}$:

$$\min_{\mathbf{W}} L(\mathbf{W}) = \sum_{n=1}^{N} \|\mathbf{y}_n - \mathbf{W}\mathbf{x}_n\|^2 \qquad \text{(the loss)}$$

We want $\operatorname{rank}(\mathbf{W}) \leq r \Leftrightarrow \mathbf{W} = \mathbf{U}\mathbf{V}^T$ with $\mathbf{U}_{d \times r}$, $\mathbf{V}_{D \times r}$, $r < \min(d, D)$ (the compression type). This results in (D+d)r parameters in (\mathbf{U}, \mathbf{V}) , which is smaller than the Dd parameters in \mathbf{W} if r is small enough.

Hence, the problem (called reduced-rank regression (RRR) in statistics) is:

$$\min_{\mathbf{U},\mathbf{V}} L(\mathbf{U},\mathbf{V}) = \sum_{n=1}^{N} \left\| \mathbf{y}_{n} - \mathbf{U}\mathbf{V}^{T}\mathbf{x}_{n} \right\|^{2}$$
 s.t. $\mathbf{U}^{T}\mathbf{U} = \mathbf{I}$

This could be solved in various ways, more or less convenient

- using gradient-based optimization, since the problem is differentiable
- lacktriangle using alternating optimization over ${f U}$ and ${f V}$
- in fact, there is a closed-form solution as an eigenproblem.

But all these solutions are specific to the form of the RRR problem, and do not apply generally to other losses, models or compressions.

An illustrative example: low-rank compression (cont.)

For example, imagine that we want to compress the matrix W so that, rather than having rank r, its elements are quantized to K values $C = \{c_1, \ldots, c_K\} \subset \mathbb{R}$. The problem would then be (quantization):

$$\min_{\mathbf{W},\mathcal{C}} L(\mathbf{W},\mathcal{C}) = \sum_{n=1}^{N} \|\mathbf{y}_n - \mathbf{W}\mathbf{x}_n\|^2$$
 s.t. $\{w_{ij}\}_{i,j=1}^{d,D} \subset \mathcal{C}$

Or, if we want to compress W so it uses binary values (binarization):

$$\min_{\mathbf{W}} L(\mathbf{W}, \mathcal{C}) = \sum_{n=1}^{N} \|\mathbf{y}_n - \mathbf{W}\mathbf{x}_n\|^2$$
 s.t. $\{w_{ij}\}_{i,j=1}^{d,D} \subset \{-1, +1\}$

None of the common approaches for RRR mentioned above work in these other cases; the objective function is not even differentiable.

An illustrative example: low-rank compression (cont.)

Rather than trying to design a different optimization algorithm for each combination of loss L and compression technique, we advocate the use of a single algorithm that can handle, with small adjustments, each combination.

Idea (for low-rank): introduce auxiliary variables $\mathbf{W} = \mathbf{U}\mathbf{V}^T$ and define the model compression problem as a constrained problem:

$$\min_{\mathbf{W},\mathbf{U},\mathbf{V}} L(\mathbf{W}) = \sum_{n=1}^{N} \|\mathbf{y}_n - \mathbf{W}\mathbf{x}_n\|^2 \quad \text{s.t.} \quad \mathbf{W} = \mathbf{U}\mathbf{V}^T, \quad \mathbf{U}^T\mathbf{U} = \mathbf{I}$$

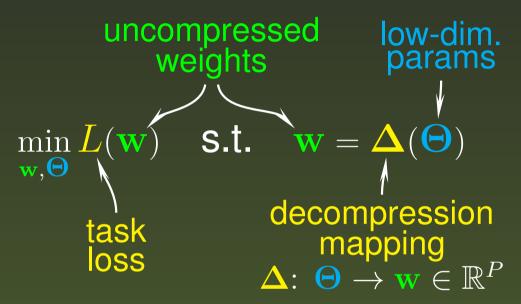
$$\text{compression part}$$

This separates the loss part from the compression part. The resulting "learning-compression" (LC) algorithm alternates training a regularized loss with compressing W via the SVD.

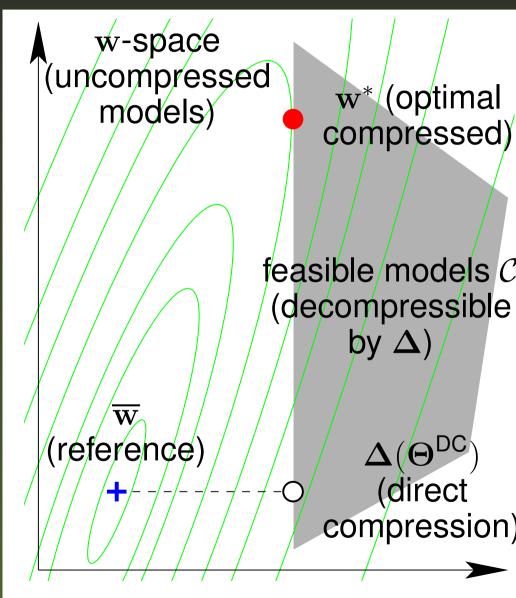
M. Á. Carreira-Perpiñán: Model compression as constrained optimization, with application to neural nets.

Part I: general framework. arXiv, 2017.

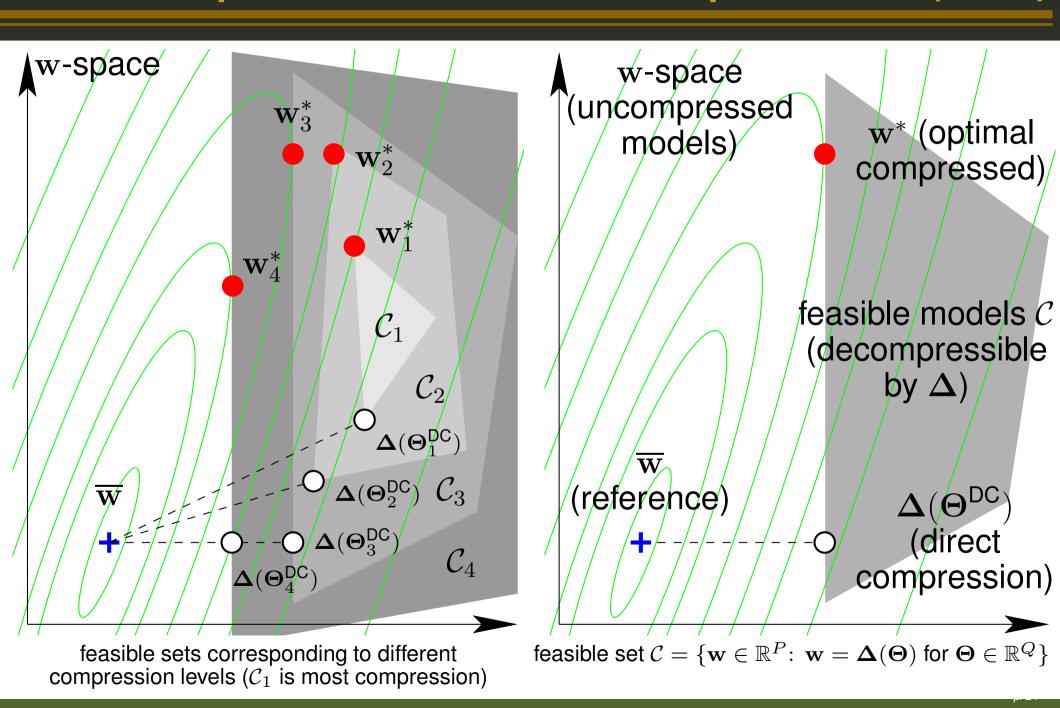
General formulation:



Compression and decompression are usually seen as algorithms, but here we regard them as mathematical mappings in parameter space. The details of the compression technique are abstracted in $\Delta(\Theta)$.



feasible set $\mathcal{C} = \{\mathbf{w} \in \mathbb{R}^P \colon \, \mathbf{w} = \mathbf{\Delta}(\mathbf{\Theta}) \; ext{for } \mathbf{\Theta} \in \mathbb{R}^Q \}$



The input to the problem $\min_{\mathbf{w}, \mathbf{\Theta}} L(\mathbf{w})$ s.t. $\mathbf{w} = \Delta(\mathbf{\Theta})$ is:

- a machine learning model with parameters w of dimension P
 e.g. a neural net with weights arranged as a feedforward architecture
- \clubsuit a machine learning task (dataset and loss function): $L(\mathbf{w})$ e.g. classification on MNIST via the cross-entropy
- \spadesuit a compression technique: given some low-dimensional parameters Θ (of dimension Q), we can recover some (not all) uncompressed parameters \mathbf{w} via the decompression mapping: $\mathbf{w} = \Delta(\Theta)$ e.g. low-rank, quantization, pruning...
- The compression ratio will be related to the dimensions of w and Θ .
- The solution of the problem is the numerical values of an optimal Θ^* . This implicitly corresponds to some uncompressed model $\mathbf{w}^* = \Delta(\Theta^*)$, but we never make it explicit: the implementation of the model uses Θ^* , not \mathbf{w}^* (both in terms of storage and of inference, i.e., applying the model to an input). Ex: compute $(\mathbf{U}^*(\mathbf{V}^{*T}\mathbf{x}))$ rather than $\mathbf{W}^*\mathbf{x}$.

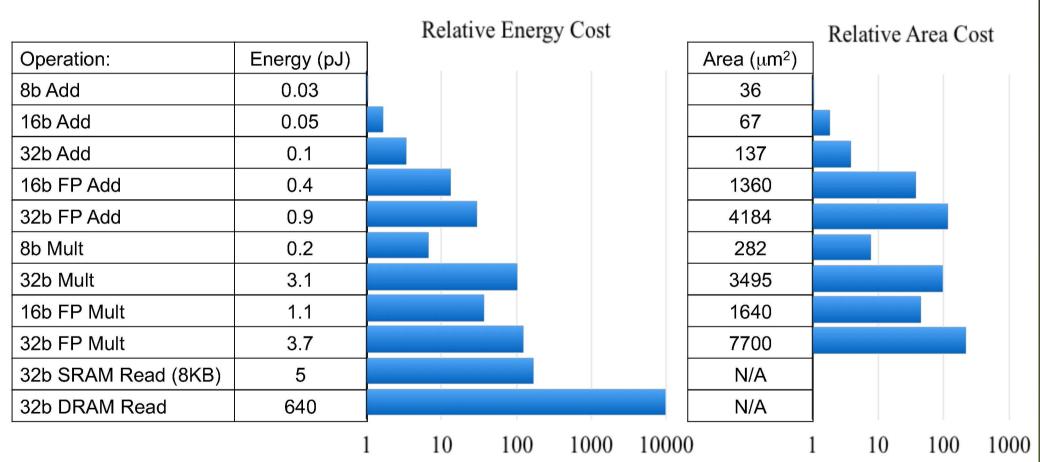
We can make the "model compression as constrained optimization" framework more flexible by adding a penalty term on Θ :

$$\min_{\mathbf{w}, \mathbf{\Theta}} L(\mathbf{w}) + \lambda C(\mathbf{\Theta})$$
 s.t. $\mathbf{w} = \Delta(\mathbf{\Theta})$

This has several uses:

- $C(\Theta)$ can be the runtime or energy of the inference pass in a neural net (based on the number of arithmetic and data movement operations).
 - The size of Θ (memory storage) is given by its dimension and the constraint on $\Delta(\Theta)$.
 - The tradeoff between accuracy L and cost C is controlled by λ .
- It helps to do model selection efficiently, e.g. learning automatically the rank of each layer of a neural net.





Energy numbers are from Mark Horowitz "Computing's Energy Problem (and what we can do about it)", ISSCC 2014

Area numbers are from synthesized result using Design Compiler under TSMC 45nm tech node. FP units used DesignWare Library.

Outline

- Machine learning and the IoT: the need for model compression
- Direct compression via standard signal processing algorithms
- Model compression as constrained optimization
- The "learning-compression" (LC) algorithm
- Several examples of deep neural net compression:
 - 1. Low-rank
 - 2. Quantization
 - 3. Pruning
 - 4. Learning the ranks
 - 5. Additive combination of compressions
- Beyond compression:
 - Model compression as structure learning
 - Model compression as model compilation

The "learning-compression" (LC) algorithm

Use a penalty method (e.g. quadratic penalty, augmented Lagrangian). Minimize:

$$Q(\mathbf{w}, \mathbf{\Theta}; \mu) = L(\mathbf{w}) + \frac{\mu}{2} ||\mathbf{w} - \mathbf{\Delta}(\mathbf{\Theta})||^2$$

as $\mu \to \infty$, using alternating optimization over w and Θ :

- \bigstar L ("learning") step: $\min_{\mathbf{w}} L(\mathbf{w}) + \frac{\mu}{2} ||\mathbf{w} \Delta(\Theta)||^2$ independent of the compression technique
- \diamond C ("compression") step: $\min_{\Theta} \|\mathbf{w} \Delta(\Theta)\|^2$ independent of the loss and dataset.

Generic: one algorithm, many compression techniques.

To change the compression type, we just need to change the C step.

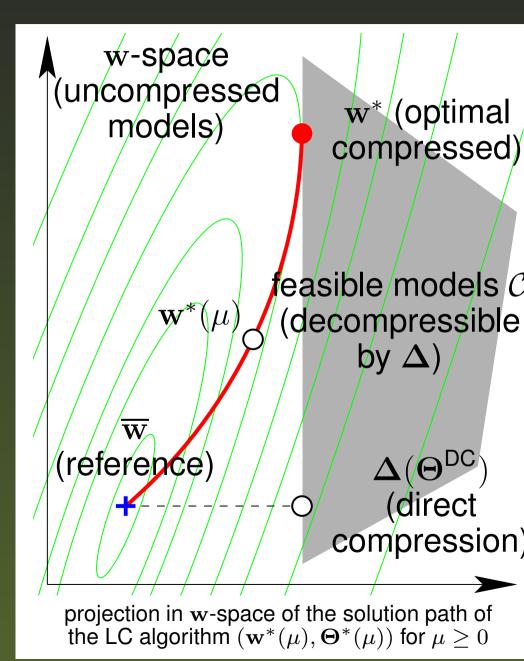
At all times, the LC algorithm keeps two copies of the parameters:

- An uncompressed one in w, updated in the L step.
- ◆ A compressed one in ⊖, updated in the C step. This is associated with an uncompressed $\Delta(\Theta)$. This always equals the "direct compression" of the current w (orthogonal projection of w on the feasible set of compressed models).

The "learning-compression" (LC) algorithm (cont.)

The LC algorithm is a homotopy method. The minimizers of Q trace a path $(\mathbf{w}^*(\mu), \mathbf{\Theta}^*(\mu))$ for $\mu \geq 0$:

- * Beginning of the path $(\mu \to 0^+)$: direct compression $(\overline{\mathbf{w}}, \boldsymbol{\Theta}^{DC})$ (training the reference model, then compressing its weights).
- * End of the path $(\mu \to \infty)$: local solution $(\mathbf{w}^*, \mathbf{\Theta}^*)$ of the constrained problem.



LC algorithm pseudocode, aug. Lagrangian version

input training data and model with parameters (weights) w

$$\mathbf{w} \leftarrow \overline{\mathbf{w}} = \operatorname{arg\,min}_{\mathbf{w}} L(\mathbf{w})$$

reference model

$$oldsymbol{\Theta} \leftarrow oldsymbol{\Theta}^{\mathsf{DC}} = oldsymbol{\Pi}(\overline{\mathbf{w}}) = rg\min_{oldsymbol{\Theta}} \|\overline{\mathbf{w}} - oldsymbol{\Delta}(oldsymbol{\Theta})\|^2$$

compress reference model

$$\lambda \leftarrow 0$$

for
$$\mu = \mu_0 < \mu_1 < \cdots < \infty$$

$$\mathbf{w} \leftarrow \operatorname{arg\,min}_{\mathbf{w}} L(\mathbf{w}) + \frac{\mu}{2} \|\mathbf{w} - \mathbf{\Delta}(\mathbf{\Theta}) - \frac{1}{\mu} \mathbf{\lambda} \|^2$$

L step: learn model

$$\mathbf{\Theta} \leftarrow \mathbf{\Pi}(\mathbf{w} - \frac{1}{u}\boldsymbol{\lambda}) = \arg\min_{\mathbf{\Theta}} \|\mathbf{w} - \frac{1}{u}\boldsymbol{\lambda} - \boldsymbol{\Delta}(\mathbf{\Theta})\|^{2}$$

C step: compress mode

$$\lambda \leftarrow \lambda - \mu(\mathbf{w} - \Delta(\Theta))$$

Lagrange multipliers

if $\|\mathbf{w} - \mathbf{\Delta}(\mathbf{\Theta})\|$ is small enough then exit the loop

return w, Θ

Runtime: a few times slower than training the reference model.

LC algorithm: convergence

The "model compression as constrained optimization" framework includes a large variety of problems, usually nonconvex, whose nature can be continuous (e.g. low-rank, pruning with ℓ_1) or combinatorial (e.g. quantization, pruning with ℓ_0).

Convergence of the LC algorithm to a (local) minimizer of the constrained problem can be established for the easier cases (smooth convex) but not in general for the harder cases (nonconvex, discrete). These are the more interesting ones in practice, and the ones the LC algorithm is designed to handle.

We suspect that the partial separability of Q over w and Θ , and the fact that the C (compression) step is often exactly solvable, may afford stronger results, but this is an open question.

LC algorithm: convergence (cont.)

Smooth case: contin. differentiable loss & decompression mapping.

Theorem 1 (for the QP $Q(\mathbf{w}, \Theta; \mu)$). Given a positive increasing sequence $(\mu_k) \to \infty$, a nonnegative sequence $(\tau_k) \to 0$, and a starting point (\mathbf{w}^0, Θ^0) , suppose the QP method finds an approximate minimizer (\mathbf{w}^k, Θ^k) of $Q(\mathbf{w}^k, \Theta^k; \mu_k)$ that satisfies $\|\nabla_{\mathbf{w}, \Theta} Q(\mathbf{w}^k, \Theta^k; \mu_k)\| \le \tau_k$ for $k = 1, 2, \ldots$ Then, $\lim_{k \to \infty} \left((\mathbf{w}^k, \Theta^k) \right) = (\mathbf{w}^*, \Theta^*)$, which is a KKT point for the constrained problem, and its Lagrange multiplier vector has elements $\boldsymbol{\lambda}_i^* = \lim_{k \to \infty} \left(-\mu_k \left(\mathbf{w}_i^k - \boldsymbol{\Delta}(\Theta_i^k) \right) \right), i = 1, \ldots, P$.

That is, there exists a continuous path $(\mathbf{w}^*(\mu), \mathbf{\Theta}^*(\mu))$ that converges to a local stationary point (typically a minimizer) of the constrained problem and it can be loosely followed by approximately (but accurately enough) optimizing Q as $\mu \to \infty$.

Does alternating optimization of Q for fixed μ over \mathbf{w} and $\mathbf{\Theta}$ converge to a solution $(\mathbf{w}^*(\mu), \mathbf{\Theta}^*(\mu))$?

- \diamond If Q is convex, yes quite generally.
- ❖ If Q is nonconvex: convergence results are complex and more restrictive. One simple case where convergence occurs is if there is a unique minimizer along each block (w and Θ).

LC algorithm: convergence (cont.)

Nonsmooth or discrete case.

- * These often are NP-hard problems, e.g. quantization (with adaptive assignments), binarization and ℓ_0 pruning.
- Convergence results are an open problem.
- With mixed continuous-discrete problems, the LC algorithm has the flavor of k-means; it stops at a "local optimum" in the sense that no further progress is possible.
- Empirically, the LC algorithm is competitive with specialized approaches for quantization, pruning or low-rank compression.
- It seems able to explore the discrete space (of weight assignments to codebook entries, or of the weight subset to be pruned) while driving down the loss.

LC algorithm: solving the L step

The L step always takes the same form regardless of the compression technique:

$$\min_{\mathbf{w}} L(\mathbf{w}) + \frac{\mu}{2} \|\mathbf{w} - \mathbf{\Delta}(\mathbf{\Theta})\|^{2}$$

- This is the original loss on the uncompressed weights \mathbf{w} , regularized with a quadratic term on the weights (which "decay" towards a validly compressed model $\Delta(\Theta)$).
- It can be optimized in the same way as the reference net. Simply add $\mu(\mathbf{w} \Delta(\Theta))$ to the gradient.
- * With large datasets we typically use SGD. We clip the learning rates so they never exceed $\frac{1}{\mu}$ to avoid oscillations as $\mu \to \infty$.
- Straightforward to integrate into existing deep learning frameworks.

The L step behaves like a regular, machine-learning training of the original model. The compression technique appears only in the C step.

LC algorithm: solving the C step

Solving the C step means optimally compressing the current weights w:

$$\min_{\mathbf{\Theta}} \|\mathbf{w} - \mathbf{\Delta}(\mathbf{\Theta})\|^2$$

- * Fast: it does not require the dataset (which appears in $L(\mathbf{w})$).

 Runtime usually negligible compared to the L step.
- Find the closest compressed model $\mathbf{v} = \mathbf{\Pi}(\mathbf{w})$ to \mathbf{w} : $\min_{\mathbf{e},\mathbf{v}} \|\mathbf{w} \mathbf{v}\|^2$ s.t. $\mathbf{v} = \mathbf{\Delta}(\mathbf{e})$.
- The solution depends on the choice of the decompression mapping Δ , and is known for many common compression techniques:
 - binarization: binarize each weight
 - lack quantization in general: k-means
 - pruning: zero all but top weights

- closest matrix of type X to W(X = low-rank, circulant, Kronecker...)
- ... and combinations thereof.
- \bullet If having a penalty term: $\min_{\bullet} \|\mathbf{w} \mathbf{\Delta}(\bullet)\|^2 + \lambda C(\bullet)$.

The C step behaves like a regular compression of a signal and simply requires calling the corresponding compression subroutine as a black box. The loss and training set appear only in the L step.

Outline

- Machine learning and the IoT: the need for model compression
- Direct compression via standard signal processing algorithms
- Model compression as constrained optimization
- The "learning-compression" (LC) algorithm



Several examples of deep neural net compression:

- 1. Low-rank
- 2. Quantization
- 3. Pruning
- 4. Learning the ranks
- 5. Additive combination of compressions
- Beyond compression:
 - Model compression as structure learning
 - Model compression as model compilation

1. Low-rank compression: C step (unpublished)

The matrix of parameters should be low-rank. In a deep neural net, this is applied to all (or some) layers, so that layer k's weight matrix \mathbf{W}_k (of $a_k \times b_k$) should be low-rank. If we set the maximum rank of layer k to r_k :

- lacktriangle Decompression mapping: $\mathbf{W}_k = \mathbf{U}_k \mathbf{V}_k^T$ with $\mathbf{U}_{a_k \times r_k}$ and $\mathbf{V}_{b_k \times r_k}$, so $m{\Theta} = \{\mathbf{U}_k, \mathbf{V}_k\}_{k=1}^K$.
- \diamond Solution: compute the SVD of \mathbf{W}_k and preserve its leading r_k s.v.'s.
- Fquivalent to a new neural net architecture that adds an internal linear layer of $r_k \le a_k, b_k$ units to layer k, effectively reducing its number of units to r_k . This is quite a drastic alteration of the model and usually is not competitive with quantization or pruning.
- * Computationally, it is more efficient to eliminate the constraints and apply the chain rule to minimize the loss $L(\{\mathbf{U}_k, \mathbf{V}_k\}_{k=1}^K)$ directly, rather than use the LC algorithm.

Much better: optimize over the rank r_k of each layer, as well as over $\{U_k, V_k\}_{k=1}^K$. This can also be done with the LC algorithm (see later), 38

2. Quantization with adaptive codebook: C step

M. Á. Carreira-Perpiñán and Y. Idelbayev: *Model compression as constrained optimization, with application to neural nets. Part II: quantization.* arXiv, 2017.

In scalar quantization, each real-valued weight w_i must equal one of the K values in a codebook $\mathcal{C} = \{c_1, \dots, c_K\} \subset \mathbb{R}$.

- \bullet Compressed net: codebook (K floats) + $\lceil \log_2 K \rceil$ bits per weight.
- The decompression mapping Δ is a table lookup $w_i = c_{\kappa(i)}$ where $\kappa: \{1, \ldots, P\} \to \{1, \ldots, K\}$ is a discrete mapping that assigns each weight to one codebook entry.
- \bullet Rewriting $\kappa()$ using binary assignment variables $\mathbf{Z} \in \{0,1\}^{P \times K}$:

$$\min_{\mathcal{C}, \kappa} \sum_{p=1}^{P} \|w_i - c_{\kappa(i)}\|^2 \iff \min_{\mathcal{C}, \mathbf{Z}} \sum_{i,k=1}^{P,K} \mathbf{z}_{ik} \|w_i - \mathbf{c}_k\|^2 \text{ s.t. } \begin{cases} \sum_{k=1}^{K} \mathbf{z}_{ik} = 1\\ i = 1, \dots, P. \end{cases}$$

- So $\Theta = \{C, \mathbb{Z}\}$, and the decompression mapping is $w_i = \sum_{k=1}^K z_{ik} c_k$.
- This is the squared distortion problem. With scalars, it is solvable exactly in polynomial time. Or, approximate solution by k-means.

Alternating optimization over $\mathcal C$ and $\mathbf Z$, converges to a local minimum in a finite number of iterations.

2. Quantization with fixed codebook: C step

Sometimes (e.g. because of fast hardware arithmetic), it is convenient to use a predetermined codebook, often implicit (no need to store it), so $\bigcirc = \mathbb{Z}$, or equivalently the decompression mapping is $w_i = \theta_i$ s.t. $\theta_i \in \mathcal{C}$:

- \Leftrightarrow Binarization: $C = \{-1, +1\}$.
- **Ternarization:** $C = \{-1, 0, +1\}.$
- * Powers of 2: $C = \{0, \pm 1, \pm 2^{-1}, \dots, \pm 2^{-C}\}.$

The exact solution of the C step is to assign each weight to its closest codebook entry (e.g. binarization: $w_i \to \operatorname{sgn}(w_i)$).

A partly adaptive codebook can be obtained by learning a scale factor a > 0, e.g. $C = \{-a, +a\}$ or $\{-a, 0, +a\}$. The C step for these can also be solved exactly (e.g. binarization: $a = \frac{1}{P} \sum_{i=1}^{P} |w_i|, w_i \to a \operatorname{sgn}(w_i)$).

However, all these variations are generally worse, in terms of the compression-loss tradeoff, than using an adaptive codebook.

2. Quantization with fixed assignments: C step

The constraints have the form $w_i = c_{\kappa(i)}$ for each weight w_i , where the assignment mapping $\kappa()$ is fixed, so $\Theta = \mathcal{C}$. Equivalently, the weights associated with the kth centroid must take the same value:

$$w_i = c_k \qquad \forall i \in \{1, \dots, P\} : \kappa(i) = k$$

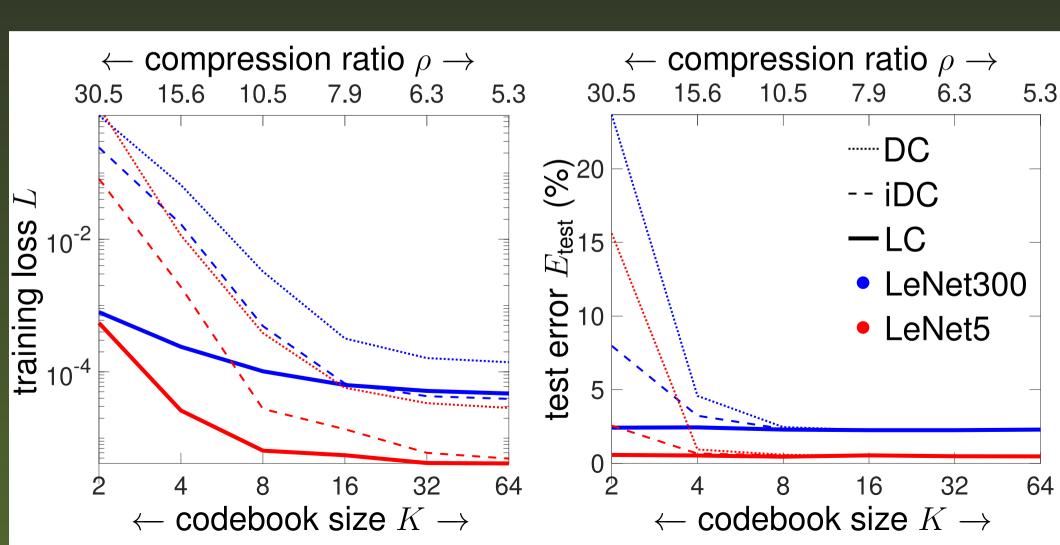
The C step has an exact solution where c_k equals the mean of the weights associated with it.

This is the same as weight sharing, a widely successful technique in convolutional neural nets, where we have a priori information about what good assignments should be (homogeneous filters). It is also known as parameter tying in speech recognition using HMMs with Gaussian mixtures.

Computationally, it is more efficient to eliminate the constraints and apply the chain rule to minimize the loss $L(\mathcal{C})$ directly over \mathcal{C} , rather than using the LC algorithm.

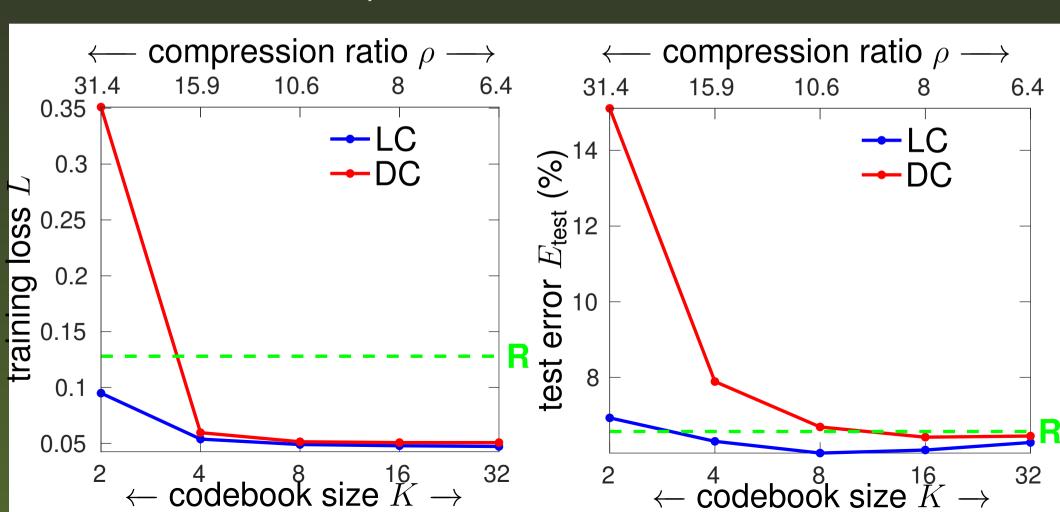
2. Quantization with adaptive codebook: experiments

More compression for the same target loss than other algorithms. LeNet neural nets on MNIST dataset: nearly no error degradation using K=2 (1 bit/weight, compression ratio $\times 30.5$). Error-compression tradeoff curves:



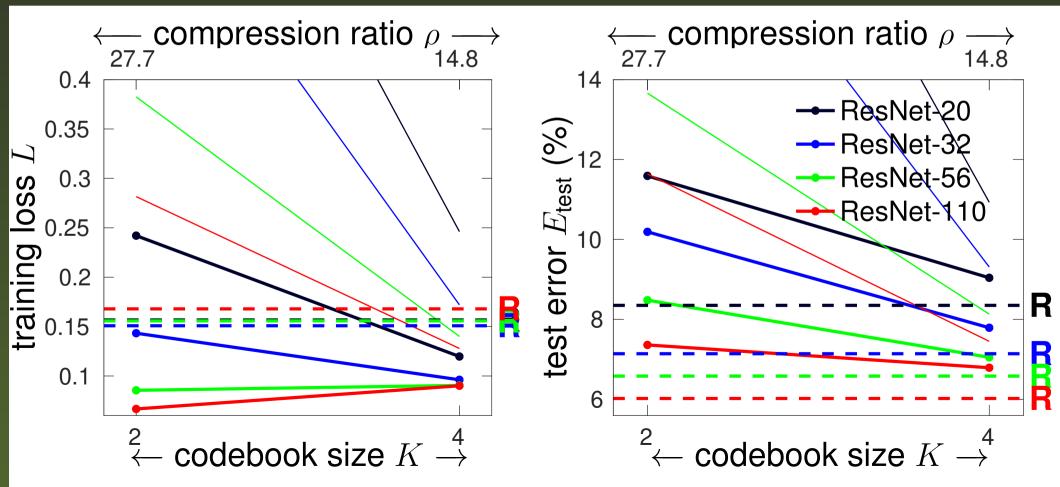
2. Quantization w/adaptive codebook: experiments (con

VGG-16 neural net on CIFAR10 dataset (16 layers, 15M weights, 57MB storage): nearly no error degradation using K=2 (1 bit/weight, compression ratio $\times 30.5$), where it would take only 1.85MB and fit on the L2 cache of modern processors.



2. Quantization w/adaptive codebook: experiments (con

ResNet neural nets on CIFAR10 dataset (20,32,56,110 layers / 0.27M,0.46M,0.85M, 1.7M weights, resp.): nearly no error degradation using K=4, some degradation using K=2 (1 bit/weight, compression ratio \times 26). ResNets are much leaner models (achieving state-of-the-art classification with a much smaller number of weights), and harder to compress. LC algorithm: —, DC: —, reference net (\mathbf{R}): ---.



3. Pruning: C step

M. Carreira-Perpiñán & Y. Idelbayev: "Learning-Compression" algorithms for neural net pruning, CVPR 2018

In pruning, the P weights are real-valued but we can have at most κ nonzero weights.

- Compressed net: store nonzero weights (and their indices).
- * Decompression mapping: $\Delta(\theta) = \theta$ s.t. $\|\theta\|_0 \le \kappa$.
 Other formulations possible: penalty instead of a constraint, or other sparsifying norm such as ℓ_1 .
- \diamond Solution: leave as is the top- κ weights (in magnitude), zero the rest. In general, some king of thresholding.
- With a deep net, this automatically finds the optimal number of weights to prune in each layer. We don't need a per-layer pruning hyperparameter κ_i .

This is like magnitude pruning, but weights are not irrevocably removed. They are marked as currently pruned in the C step (via θ), and this is interlaced with the L step (which drives \mathbf{w} towards θ). So the set of pruned weights changes during LC iterations as we converge.

Magnitude pruning is a simple approach where we remove all but the top- κ weights (in magnitude) of the reference net and retrain the remaining κ weights.

p. 45

3. Pruning: C step (cont.)

We can limit the number of nonzeros in w via an ℓ_0/ℓ_1 cost $C(\mathbf{w})$:

Constraint form:
$$\min_{\mathbf{w}, \boldsymbol{\theta}} L(\mathbf{w})$$
 s.t. $C(\boldsymbol{\theta}) \le \kappa, \mathbf{w} = \boldsymbol{\theta}$

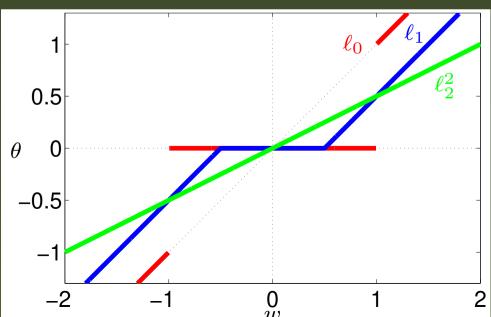
$$\Rightarrow \boldsymbol{\theta}^* = \Pi_C^{\le}(\mathbf{w}; \kappa) = \arg\min_{\boldsymbol{\theta}} \|\mathbf{w} - \boldsymbol{\theta}\|^2 \text{ s.t. } C(\boldsymbol{\theta}) \le \kappa$$

Penalty form:
$$\min_{\mathbf{w}, \boldsymbol{\theta}} L(\mathbf{w}) + \alpha C(\boldsymbol{\theta})$$
 s.t. $\mathbf{w} = \boldsymbol{\theta}$

$$\Rightarrow \boldsymbol{\theta}^* = \Pi_C^+(\mathbf{w}; \frac{2\alpha}{\mu}) = \arg\min_{\boldsymbol{\theta}} \|\mathbf{w} - \boldsymbol{\theta}\|^2 + \frac{2\alpha}{\mu} C(\boldsymbol{\theta})$$

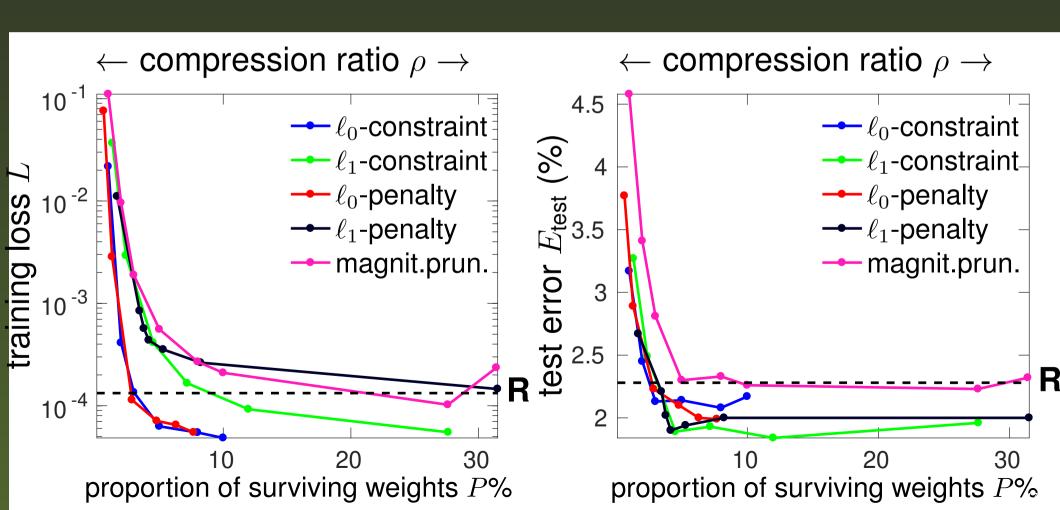
The solution of all cases is in an elementwise thresholding operator that computes θ_i from w_i for each weight.

$C(\mathbf{w})$	Constraint form † $oldsymbol{ heta}^* = \Pi^{\leq}_C(\mathbf{w};\kappa)$	Penalty form $oldsymbol{ heta}^* = \Pi^+_C(\mathbf{w}; rac{2lpha}{\mu})$
ℓ_0 : $\ \mathbf{w}\ _0$	$w_i \cdot I(w_i > \eta_0)$	$w_i \cdot I\left(w_i > \sqrt{\frac{2\alpha}{\mu}}\right)$
$\ell_1 \colon \left\ \mathbf{w} ight\ _1$	$(w_i - \operatorname{sgn}(w_i) \eta_1) \cdot I(w_i > \eta_1)$	$\left(w_i - \operatorname{sgn}\left(w_i\right) \frac{\alpha}{\mu}\right) \cdot I\left(w_i > \frac{\alpha}{\mu}\right)$
ℓ_2^2 : $\ \mathbf{w}\ _2^2$	$\sqrt{\kappa} w_i / {\ \mathbf{w} \ }_2$	$w_i/(1+\frac{2\alpha}{\mu})$
†Each formula applies if $C(\mathbf{w}) > \kappa$, otherwise $\theta = \mathbf{w}$.		



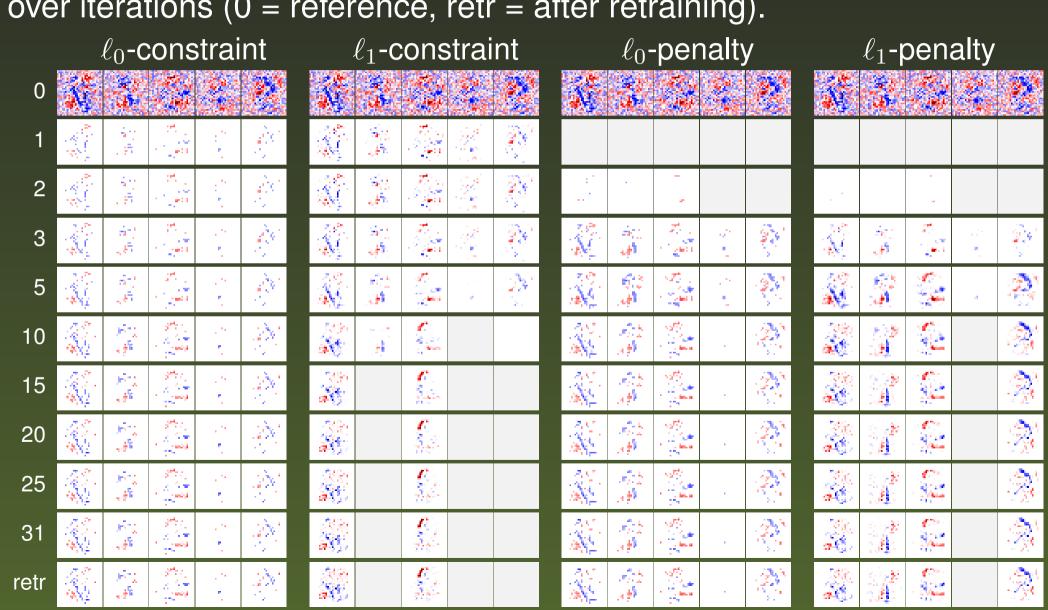
3. Pruning: experiments

More compression for the same target loss than other algorithms. LeNet300 neural net on MNIST dataset (3 layers, 266k weights): nearly no error degradation using P=2% of the weights. Magnitude pruning removes all but the largest P% weights and retrains them.



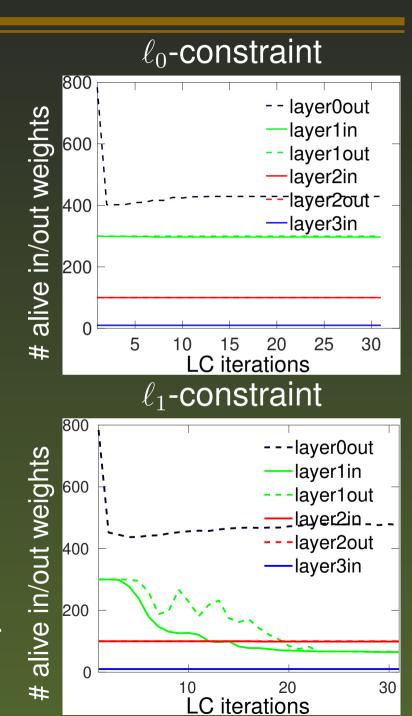
3. Pruning: experiments (cont.)

Weight vector of selected first-layer neurons for LeNet300 with $P \approx 5\%$ over iterations (0 = reference, retr = after retraining).



3. Pruning: experiments (cont.)

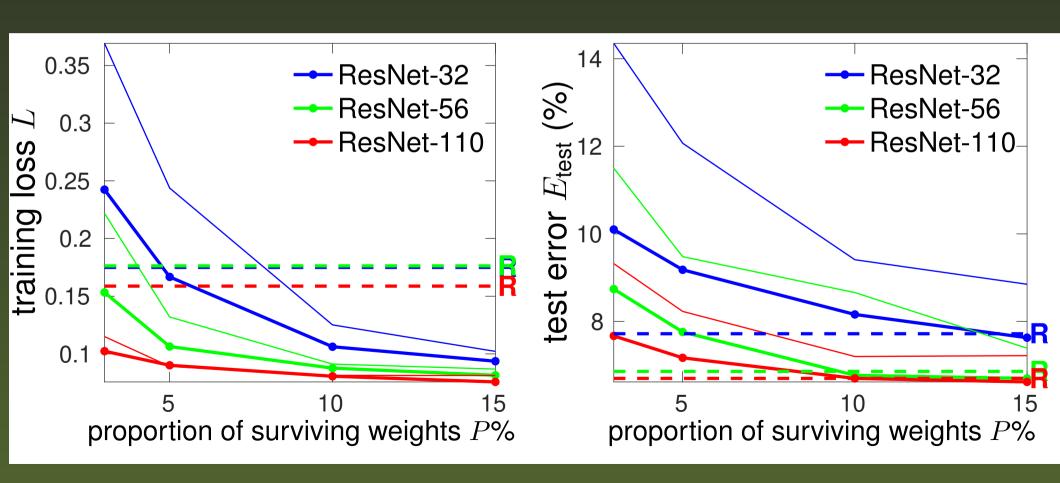
- The final weight vectors often segment the image into negative and positive regions reminiscent of center-surround receptive fields, but these regions are sparse rather than compact.
- * The algorithm prunes weights, not neurons, but we observe aggressive neuron pruning, particularly in the first layer. The original LeNet300 architecture 784-300-100-10 becomes 400-64-99-10 with similar error (for the ℓ_1 -constraint).
- Hence, the algorithm might be useful to do feature selection and determine the optimal number of neurons in each layer automatically (structure learning).



3. Pruning: experiments (cont.)

ResNet neural nets on CIFAR10 dataset (32, 56, 110 layers / 0.46M, 0.85M, 1.7M weights, resp.): we can reach $P \approx 5$ –10% surviving weights with no error degradation.

LC algorithm (ℓ_0 -constraint): —, magnitude pruning: —, reference net (**R**): ---.



4. Learning the ranks (unpublished) results

5. Additive combination of compressions (unpublished) results

Quantization vs pruning

Our common framework (LC algorithm) allows us to compare different compression techniques.

Assume:

- \diamond A neural net with P_1 multiplicative weights and P_0 biases.
- We don't compress biases since they are more critical and $P_0 \ll P_1$.
- A float takes b = 32 bits storage.

Compression ratio
$$\rho = \frac{\text{\#bits(reference)}}{\text{\#bits(compressed)}} = \frac{(P_1 + P_0)b}{\text{\#bits(compressed)}}.$$

Quantization (codebook of K entries):

#bits(compressed) =
$$P_1\lceil \log_2 K \rceil + (P_0 + K)b \Rightarrow \rho(K) \approx \frac{b}{\log_2 K}$$
.

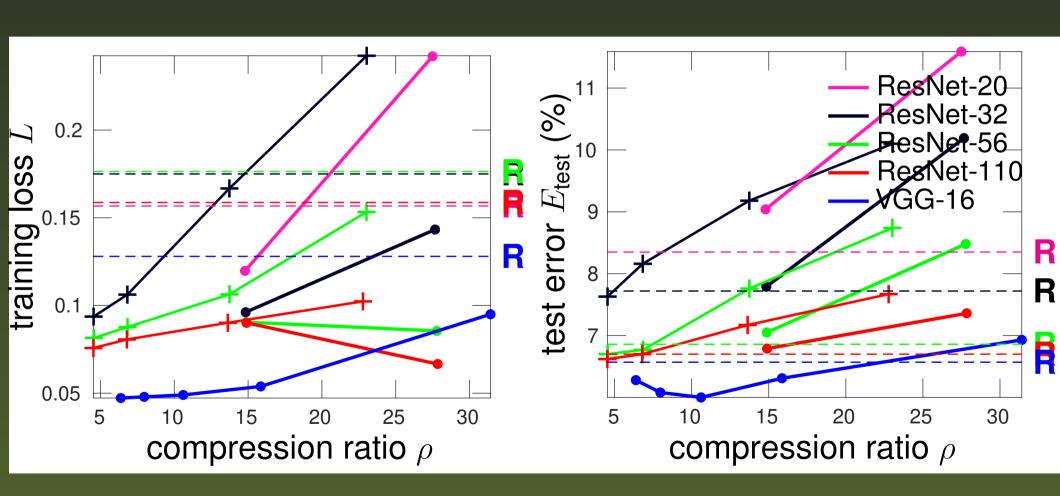
Pruning (P% surviving weights) in coordinate list format (which stores (row,column,value) per nonzero weight) and $b_{idx} = \#bits(row,column)$:

#bits(compressed) =
$$(PP_1 + P_0)(b + b_{idx}) \Rightarrow \rho(P) \approx \frac{1}{P} \frac{b}{(b+b_{idx})}$$
.

Quantization vs pruning (cont.)

Quantization beats pruning in terms of storage.

Quantization: ●, pruning: +, reference net (**R**): ---. All results using the LC algorithm.



Outline

- Machine learning and the IoT: the need for model compression
- Direct compression via standard signal processing algorithms
- Model compression as constrained optimization
- The "learning-compression" (LC) algorithm
- Several examples of deep neural net compression:
 - 1. Low-rank
 - 2. Quantization
 - 3. Pruning
 - 4. Learning the ranks
 - 5. Additive combination of compressions



Beyond compression:

- Model compression as structure learning
- Model compression as model compilation

Model compression as structure learning (unpublished)

Compression is a sophisticated form of regularization:

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

$$\min Q(\mathbf{w}, \mathbf{\Theta}; \mu) = L(\mathbf{w}) + \frac{\mu}{2} \|\mathbf{w} - \Delta(\mathbf{\Theta})\|^2$$

- The regularization term " $\|\mathbf{w} \Delta(\Theta)\|^2$ " is itself parameterized. Cf. weight decay / ridge regression.
- The parameters w must take a particular structure which is itself learned (which weights are nonzero, which weights are equal, etc.).

We could use this to achieve better generalization rather than (just) compression. In fact, a better reference model (lower test error) is often achieved by applying a small amount of compression to the original reference (slightly smaller K, P, r).

So we should always use compression as a postprocessing step to get a better model (not necessarily for actual compression).

But we can take this further...

Model compression as structure learning (cont.)

Traditionally, the structure of the model (e.g. number of layers and hidden units) is designed by hand and optimized via model selection. The framework above suggests a different approach:

- 1. Define a "blank slate" model, large and over-parameterized.

 Ex: a neural net with many layers & units.
- 2. Learn its structure and parameter values using the LC algorithm.
 - removing connections (pruning)
 - forcing subsets of parameters to be equal (quantization)
 - thinning layers (low-rank)

The LC algorithm can also learn the per-layer sizes (codebook sizes, number of nonzeros, ranks...).

- 3. The "total size" (effective # parameters or compression level) is controlled by a hyperparameter λ , which could be cross-validated.
- Indeed, particular cases of these are well-known in statistics (usually applied with linear models):
 - Sparsity: Lasso and related models.
 - Low-rank: reduced-rank regression (RRR).

Model compression as model compilation (unpublished)

Traditional compilers (invented in the 1950s): a compiler translates (highly efficiently) a source program into a binary executable:

- Input: program = a description of an algorithm in a high-level programming language. Handled by the compiler's front end.
 Ex: a C program for Quicksort.
- Output: binary code = an implementation of the program that runs in a specific, target system. Handled by the compiler's back end. Ex: instruction set, memory structure (register, caches, DRAM), CPU/GPU, etc.

This is much simpler than having to write a separate assembly code for each target system, and makes software development practical and scalable. Write M+N compilers (M front ends and N back ends) rather than MN compilers. The compiler creates efficient code (in execution runtime, memory size, etc.) by using various "optimizations".

Loop unrolling, inlining, constant folding, algebraic simplification, common subexpression elimination, etc.

Crucially, a compiler outputs binary code that is correct.

A binary code for Quicksort will correctly sort any input array, on any target system.

Model compression as model compilation (cont.)

Machine learning compilers: a compiler translates (highly efficiently) a source program into a binary executable:

- Input: program = a description of a ML model f(x; w) in a high-level programming language, i.e., numerical values for the parameters w and a description of how to compute y = f(x; w) (inference). Ex: the weight values and architecture (layers, connectivity, nonlinearities, etc.) of a deep neural net. And a loss function L(w) which quantifies the goodness of a specific w on a training set.
- Output: binary code = an implementation of the program $f(x; w^*)$ that runs in a specific system. We may have $w \neq w^*$.

 Ex: instruction set, memory structure (register, caches, DRAM), CPU/GPU, multiprocessor, etc.

The compilation is inexact: the binary code implements the model for a different parameter vector **w*** (the compiler can alter the source code). But **w*** should be picked optimally so that it minimizes the loss subject to system constraints and memory/runtime/energy requirements. Much can be gained in terms of efficiency by slightly increasing the loss.

Model compression as model compilation (cont.)

The compilation process should be defined as a numerical optimization problem involving the system, and should be generic and modular, to handle arbitrary combinations of source program and target system. This may be done via the "model compression as constrained optimization" framework and the LC algorithm:

$$\min_{\mathbf{w}, \mathbf{\Theta}} L(\mathbf{w}) + \lambda C(\mathbf{\Theta})$$
 s.t. $\mathbf{w} = \Delta(\mathbf{\Theta})$

- ♦ w = high-level source program, defines the model
- \bullet \bullet = low-level program (binary code), defines the model $\Delta(\bullet)$ that is actually implemented in the target system
- $\star L(\mathbf{w})$ ensures program correctness in an optimal sense
- $C(\Theta)$: cost of the low-level program, system-dependent Runtime, memory, energy...
- Δ (Θ) is defined by the target system constraints Memory and processors available, instruction set...

In a traditional (correct) compiler, $\mathbf{w} = \boldsymbol{\Delta}(\boldsymbol{\Theta})$ (lossless compilation). $_{\scriptscriptstyle p.60}$

Model compression as model compilation (cont.)

The L step ensures the final program $\mathbf{w}^* = \Delta(\mathbf{\Theta}^*)$ optimizes the loss.

The C step does the compilation of the current program w in an inexact way. In addition to compression techniques (quantization, pruning, etc.), it should use standard compiler techniques to make best use of the system architecture, e.g.:

- \diamond layout and reshaping (e.g. storing Θ vs Θ^T , tensors...)
- memory allocation (registers, caches, DRAM, etc.)
- scheduling of threads (e.g. for a matrix-vector product in a GPU)
- choice of machine instructions (e.g. loop vs vectorized)
- loop unrolling, inlining, constant folding, algebraic simplification, common subexpression elimination, etc.

The cost C may be estimated by a model (e.g. number of arithmetic and data movement operations) or by an actual binary code generation whose runtime is measured on a test dataset.

Conclusion

- A precise mathematical formulation of neural net compression as a constrained optimization problem.
- A generic way to solve it, the learning-compression (LC) algorithm.

 A single algorithm handles any combination of compression technique and loss function.
 - → The compression technique appears as a black-box subroutine in the C step, independent of the loss and dataset.
 We can try a different type of compression by simply calling its subroutine.
 - The loss and neural net training by SGD appear in the L step, independent of the compression technique.
 - As if we were training the original, reference net with a weight decay.
- Easy to implement in deep learning toolboxes.
- Convergence to local optimum under some assumptions.
- Very effective in practice: not much slower than training the reference model, and achieves more compression for the same target loss than direct compression and other approaches.
- Extensions: structure learning, "model compiler".

References and code

General framework:

M. Á. Carreira-Perpiñán: Model compression as constrained optimization, with application to neural nets. Part I: general framework. arXiv:1707.01209, Jul. 5, 2017.

We are developing this framework for number of compression types:

- C-P and Y. Idelbayev: Model compression as constrained optimization, with application to neural nets. Part II: quantization. arXiv:1707.04319, Jul. 13, 2017.
- ❖ C-P and Y. Idelbayev: "Learning-Compression" algorithms for neural net pruning. CVPR, 2018.
- C-P and Y. Idelbayev: Model compression as constrained optimization, with application to neural nets. Part III: pruning. arXiv, 2019.
- C-P and Y. Idelbayev: Model compression as constrained optimization, with application to neural nets. Part IV: low-rank. arXiv, 2019.
- C-P and Y. Idelbayev: Model compression as constrained optimization, with application to neural nets. Part V: additive combinations. arXiv, 2019.
- C-P and A. Zharmagambetov: Fast model compression. BayLearn, 2018.
- ❖ Y. Idelbayev and C-P: Low-rank compression of neural nets: learning the rank of each layer. Submitted.
- * Y. Idelbayev and C-P: More general and effective model compression via an additive combination of compressions. Submitted.

Partly supported by NSF award IIS-1423515, NVIDIA GPU donations and MERCED cluster (NSF grant ACI-1429783).