Low Rank Compression of Neural Networks: LC Algorithms and Open-source Implementation

Yerlan Idelbayev Advisor: Miguel Á. Carreira-Perpiñán Department of CSE, University of California, Merced http://eecs.ucmerced.edu



The code is available at:

 $\tt https://github.com/UCMerced-ML/LC-model-compression$

Machine learning and neural networks

Most of this talk applies to ML models in general, however, we focus on neural networks

Neural networks have established state-of-the-art performance nearly in every machine learning task:

- Natural Language Processing (NLP): dialog systems, translators...
- Computer vision: image and video recognition, classification...
- Speech processing: speech-to-text, audio synthesis...
- Various signal enhancement tasks (photo, audio, video)

The models trained on these tasks have significant practical importance

Some famous neural network use cases



Translators

≡ Google Translate								
TURKISH		ENGLISH						
o bir doktor			\times					
\$ •D								
he is a doctor								
•)		٥						

The images are obtained from official websites or blogposts of the services.

Current mobile devices contain dozens of neural networks, some of them running non-stop!

The improvement of NN performance





Why? The improvements are attributed to several factors:

- more data and more compute
- better software with frameworks like PyTorch, TensorFlow, MxNet
- better algorithms and vast amount of collected empirical knowledge

The deployment of NN

Two practical regimes of deployment

Cloud deployment:

- access to powerful hardware
- the hardware can be chosen
- high-end hardware is typically allocated for training only
- running at scale is expensive

Edge and mobile deployment:

- hardware can not be chosen
- stringent constraints:
 - processing capabilities (CPU, GPU)
 - power consumption and battery life
 - SLA for responsiveness (for usability)
 - for mobile devices: most hardware is old and low-end

A study from Facebook over its mobile users

Most of the mobile devices in use are old, in 2018 only

25% of users had CPUs designed later than 2013.



Most of the mobile devices in use are low-end:



The figures are obtained from Wu et al. (doi: 10.1109/HPCA.2019.00048)

Model compression

There is a mismatch between requirements of current deep models and capabilities of end-user hardware.

Research question:

• How to obtain a smaller (compressed) neural network that has as close performance (e.g., accuracy) as possible to the original large model?

We will formulate the compression problem as a constrained optimization and give an efficient algorithm based on solid optimization principles (the Learning-Compression algorithm)

In this talk we show the application of LC for the problems of low-rank compression of neural networks.

Model compression: detour on neural networks



figure from Carreira-Perpiñán and Weiran Wang, arxiv:1212.5921

A neural network with *K* layers is a computational graph comprised of *K* linear and non-linear transforms applied to input \mathbf{x} :

 $f(\mathbf{x}; \mathbf{w}) = \sigma(\mathbf{W}_K \dots \sigma(\mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x})))$

The weights $\mathbf{w} = {\mathbf{W}_1, \dots, \mathbf{W}_K}$ are trained on a dataset of input-output pairs (\mathbf{x}, \mathbf{y}) to make the network output $f(\mathbf{x}; \mathbf{w})$ closer to the true output \mathbf{y} :

regression:
$$\begin{split} \min_{\mathbf{w}} & L(\mathbf{w}) = \sum_{\mathbf{x}, \mathbf{y}} \|\mathbf{y} - f(\mathbf{x}; \mathbf{w})\|^2 \\ \text{classification:} & \min_{\mathbf{w}} & L(\mathbf{w}) = \sum_{\mathbf{x}, \mathbf{y}} \text{CrossEntropy}(\mathbf{y}, f(\mathbf{x}; \mathbf{w})) \end{split}$$

Typically, the networks are optimized using the stochastic gradient descent (SGD): a procedure that uses an estimate of the true gradient $(\frac{\partial}{\partial \mathbf{w}} L(\mathbf{w}))$ computed on training mini-batches.



Images are from the slides of Miguel Á. Carreira-Perpiñán

The LC algorithm: general formulation



Compression details are abstracted in $\Delta(\Theta)$:

- low-rank: $\Delta(\Theta) = \mathbf{U}\mathbf{V}^T$ where $\Theta = \{\mathbf{U}, \mathbf{V}\}$
- pruning: $\Delta(\Theta) = \Theta$ s.t. $\|\Theta\|_0 \le \kappa$



The LC algorithm (cont.)

Reformulate using penalty methods and optimize the following while driving $\mu \to \infty$:

$$\min_{\mathbf{w},\boldsymbol{\Theta}} \quad L(\mathbf{w}) + \lambda C(\boldsymbol{\Theta}) + \frac{\mu}{2} \|\mathbf{w} - \boldsymbol{\Delta}(\boldsymbol{\Theta})\|^2$$

Apply alternating optimization wrt \mathbf{w} and $\mathbf{\Theta}$, which gives the (LC) algorithm:

• Learning (L) step:

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

- This is a regular training of the model, but with a quadratic regularization term.
- L step is independent of compression mechanism.
- We will use SGD and standard NN software
- Compression (C) step:

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

- For $\lambda = 0$ the C step has a form of optimal projection of the weights
- Does not involve the dataset (no *L* term).
- Many well studied cases with fast solutions

The LC algorithm: pseudocode

 $\begin{array}{ll} \begin{array}{ll} \begin{array}{l} \text{input} \text{ training data and model with parameters w} \\ \hline \mathbf{w} \leftarrow \overline{\mathbf{w}} = \arg\min_{\mathbf{w}} L(\mathbf{w}) & \text{pretrained model} \\ \hline \mathbf{\Theta} \leftarrow \mathbf{\Theta}^{\mathsf{DC}} = \mathbf{\Pi}(\overline{\mathbf{w}}) & \text{init compression} \\ \hline \boldsymbol{\beta} \leftarrow \mathbf{0} & \text{init compression} \\ \hline \boldsymbol{\beta} \leftarrow \mathbf{0} & \text{init compression} \\ \hline \mathbf{for} \ \mu = \mu_0 < \mu_1 < \cdots < \infty & \text{w} \leftarrow \arg\min_{\mathbf{w}} L(\mathbf{w}) + \frac{\mu}{2} \|\mathbf{w} - \mathbf{\Delta}(\mathbf{\Theta}) - \frac{1}{\mu} \boldsymbol{\beta}\|^2 & \text{L step} \\ \hline \mathbf{\Theta} \leftarrow \arg\min_{\mathbf{\Theta}} \frac{\mu}{2} \|\mathbf{w} - \frac{1}{\mu} \boldsymbol{\beta} - \mathbf{\Delta}(\mathbf{\Theta})\|^2 + \lambda C(\mathbf{\Theta}) & \text{C step} \\ \hline \boldsymbol{\beta} \leftarrow \boldsymbol{\beta} - \mu(\mathbf{w} - \mathbf{\Delta}(\mathbf{\Theta})) & \text{if } \|\mathbf{w} - \mathbf{\Delta}(\mathbf{\Theta})\| \text{ is small enough then exit the loop} \\ \hline \mathbf{return} \ \mathbf{w}, \ \mathbf{\Theta} \end{array}$

Low-rank compression with rank selection



low-rank

We replace a matrix W with rank-r matrix

- Such a low-rank matrix can be written as the product UV^T
 - For small values of r this reduces FLOPs and storage
 - Low-rank decomposition can achieve speed-up on any hardware (unlike weight pruning)
- If ranks are known, training the low-rank models is not hard: use SGD
- Selecting the right rank for each layer is challenging for *K*-layer network, with *R* ranks per layer there are *R*^{*K*} combinations to try

We use the LC algorithm to jointly learn both weights and ranks.

Low-rank compression: formulation of the problem

For a *K*-layer deep net with weights $\mathbf{W} = (\mathbf{W}_1, \dots, \mathbf{W}_K)$. Our compression problem is defined as:

task loss compression cost (e.g., FLOPs)

$$min_{\mathbf{W}} L(\mathbf{W}) + \lambda C(\mathbf{W})$$

s.t. $r_k = \operatorname{rank}(\mathbf{W}_k) \le R_k, \ k = 1, \dots, K.$
max possible rank of layer k

- $\lambda C(\mathbf{W})$ controls the tradeoff between model performance and compression
- mixed-integer optimization
- the formulation does not allow to immediately use the LC algorithm

Low-rank compression: cost function

The compression cost C measures the model size/FLOPs and is defined as

$$C(\mathbf{W}) = C(\mathbf{r} = \{r_1, \dots, r_K\}) = C_1(r_1) + \dots + C_K(r_K).$$

Here, $C_k(r_k)$ is cost of compressing layer k using rank r_k . Assume, \mathbf{W}_k is $m \times n$ matrix:

- size: $C_k(r_k) = m \times r_k + r_k \times n$
- FLOPs: $C_k(r_k) = (m \times r_k + r_k \times n) \times M$ where *M* is #times filter (layer) has been applied

Low-rank compression: optimization

We introduce auxiliary variables Θ_k with a constraint of $\mathbf{W}_k = \Theta_k$ for $k = 1, \dots, K$:

$$\min_{\mathbf{V},\Theta,\mathbf{r}} \quad L(\mathbf{W}) + \lambda \, C(\mathbf{r})$$

s.t. $\mathbf{W}_k = \boldsymbol{\Theta}_k,$ $r_k = \operatorname{rank}(\boldsymbol{\Theta}_k) \le R_k$ $k = 1, \dots, K.$

and then apply penalty to equality constraints (with $\mu
ightarrow \infty$)

$$\min_{\mathbf{W}, \boldsymbol{\Theta}, \mathbf{r}} \quad L(\mathbf{W}) + \lambda C(\mathbf{r}) + \frac{\mu}{2} \sum_{k=1}^{K} \|\mathbf{W}_k - \boldsymbol{\Theta}_k\|^2$$
s.t. $r_k = \operatorname{rank}(\boldsymbol{\Theta}_k) \le R_k, \ k = 1, \dots, K.$

Low-rank compression: optimization (cont.)

After our simple manipulations, we obtained the following:

$$\min_{\mathbf{W},\Theta,\mathbf{r}} \quad L(\mathbf{W}) + \lambda C(\mathbf{r}) + \frac{\mu}{2} \sum_{k=1}^{K} \|\mathbf{W}_k - \mathbf{\Theta}_k\|^2$$

s.t. $r_k = \operatorname{rank}(\mathbf{\Theta}_k) \le R_k, \ k = 1, \dots, K.$

Now, the alternation over W and $\{\Theta, \mathbf{r}\}$ gives us the LC algorithm:

L step:
$$\min_{\mathbf{W}} L(\mathbf{W}) + \frac{\mu}{2} \sum_{k=1}^{K} \|\mathbf{W}_k - \mathbf{\Theta}_k\|^2$$

same old L step, will be handled by SGD

C step:
$$\min_{\Theta, \mathbf{r}} \frac{\mu}{2} \sum_{k=1}^{K} \|\mathbf{W}_k - \mathbf{\Theta}_k\|^2 + \lambda C(\mathbf{r}) \text{ s.t. } r_k \leq R_k \quad \forall k = 1, \dots, K$$

C step separates over layers into smaller subproblems

Low-rank compression: C step

$$\min_{\boldsymbol{\Theta},\mathbf{r}} \frac{\mu}{2} \sum_{k=1}^{K} \|\mathbf{W}_k - \boldsymbol{\Theta}_k\|^2 + \lambda \sum_{k=1}^{K} C(r_k), \quad \text{s.t.} \quad r_k \le R_k, \quad k = 1, \dots, K$$

The C step problem separates over layers k = 1, ..., K into:

$$\min_{\Theta_k, r_k} \frac{\mu}{2} \|\mathbf{W}_k - \Theta_k\|^2 + \lambda C(r_k) \quad \text{s.t.} \quad r_k = \operatorname{rank}(\Theta_k) \le R_k$$

Solution: If the rank of decomposition (r_k) is fixed, this reduces to best r_k -rank approximation problem of a matrix. The solution of "best rank" approximation is given by Eckhart-Young theorem and computed using SVD. However, we have several ranks to check: $r_k = 1 \dots R_k$: we simply enumerate over every possible rank and chose the one corresponding to lowest loss. No need to compute multiple SVDs, one is sufficient.

Behavior: We automatically select the best rank for our overall compression cost.

Low-rank compression: experiments on CIFAR10 networks CIFAR10: 10 classes, 32 × 32 images:





50K in train set, 10K in test set

Error-compression space of test error, inference FLOPs and number of parameters (ball size for each net). Different networks have different colors. **B** — reference network.

Results of our algorithm over different λ values for a given network span a curve, shown as connected circles

Other published results using low-rank are shown as isolated circles; filter pruning results shown as isolated squares.

Low-rank: AlexNet on ImageNet experiments

The ImageNet task asks to classify a colored image into one of the 1000 classes. AlexNet is a convolutional network with 60M parameters

	MFLOPs	top-1	top-5	$ ho_{FLOPs}$
Caffe-AlexNet [11]	724	42.70	19.80	1.00
Kim et al. [13], Tucker	272	n/a	21.67	2.66
Tai et al. [21], scheme 2	185	n/a	20.34	3.90
Wen et al. [23], scheme	1 269	n/a	20.14	2.69
Kim et al. [12], scheme 2	2 272	43.40	20.10	2.66
Yu et al. [28], filter prun.	232	44.13	n/a	3.12
Li et al. [17], filter prun.	334	43.17	n/a	2.16
Ding et al. [3], filter prun	. 492	43.83	20.47	1.47
our scheme 1, $\lambda = 0.20$	231	41.56	18.72	3.13
our scheme 2, $\lambda = 0.20$	152	41.03	18.23	4.78

Compression with our algorithm vs published work using low-rank methods and structured pruning. ρ_{FLOPs} — reduction in FLOPs.

Can we apply low-rank to optimize inference time?

Historically, low-rank was used to reduce sizes and FLOPs of the models. But:

- fewer FLOPs not necessarily mean faster runtime!
- Can we select the ranks per each layer to minimize on-device runtime? (requires on-device measurements)

Hard problem There are combinatorial number of ranks and corresponding on-device measurements. We tackle it by

- building an accurate and fast to compute runtime model
- formulating a suitable optimization problem
- and giving an efficient optimization algorithm based on Learning-Compression framework

Our target device :



Jetson N	lano
CPU	4-core ARM Cortex-A57, 1.4 GHz
GPU	128 CUDA cores at 0.9 GHz
RAM	4 GB 64-bit LPDDR4, 1.6 GHz
OS	Ubuntu 18.04.5 LTS
Kernel	GNU/Linux 4.9.140-tegra
Storage	128 GB microSDXC memory card

Inference-targeted low-rank: Device runtime model

Let's define the runtime $\mathcal{R}(\mathbf{W})$ as the time to process a single image through a *K*-layer net with weights $\mathbf{W} = {\mathbf{W}_1, \dots, \mathbf{W}_K}$.

- runtime is function of layer's ranks
- · runtime can be directly measured on device
- assuming R ranks per layer, there are R^K different configurations to measure

We model the runtime as the sum of inferences through individual layers:

$$\mathcal{R}(\mathbf{W}) = \mathcal{R}(\mathbf{r}) = \mathcal{R}_1(r_1) + \mathcal{R}_2(r_2) + \dots + \mathcal{R}_K(r_K).$$
(1)

- only need to measure R different rank configurations for each of the K layers
- total required measurements: $R \times K$ (vs R^K)

Inference-targeted low-rank: Device runtime model (cont.)

Even RK on-device measurements are time consuming and noisy, thus:

- for each layer we run measurements for equally spaced set of ranks (e.g., r = 1, 10, 20, ...)
- fit ℓ_1 regression on the measurements to interpolate and reduce noise



Inference-targeted low-rank: Problem formulation

Given a *K*-layer net with weights $\mathbf{W} = {\mathbf{W}_1, \dots, \mathbf{W}_K}$ trained on the loss *L* (e.g., cross-entropy), we formulate the following device-dependent rank selection problem:

$$\min_{\mathbf{W},\mathbf{r}} \quad L(\mathbf{W}) + \lambda \mathcal{R}(\mathbf{r})$$
s.t. rank $(\mathbf{W}_k) = r_k, \quad k = 1, \dots, K.$
(2)

Here, the term $\lambda \mathcal{R}(\mathbf{r})$ controls the tradeoff between on-device inference speed and model loss.

Inference-targeted low-rank: Experiments



optimizing for latency beats optimizing for FLOPs in terms of final latency (duh?)

Low-rank: What happens with non-matrix weights?

Weights do not necessarily come as matrices.

For example weights of convolutional layers are typically stored as NCHW or NHWC tensors.

To apply low-rank, we reshape the tensors into matrices!



*This is known as matricization in tensor algebra.

More on reshapes: Efficient implementation

Some of the reshapes give a rise to efficient low-rank schemes.



Which reshapes are the best? How to select them optimally?

· Historically a single fixed scheme was used throughout the NN for the compression

- In our CVPR2020 work, we used fixed scheme throughout.
- This is suboptimal!
- Can we select the best scheme per each layer?
- The problem involves selecting ranks as well.

Hard problem. There are a combinatorial number of configurations of ranks and schemes.

However, learning-compression algorithm can solve it too.

Low-rank with scheme selection: Problem formulation

Given a *K*-layer net with weights $\mathbf{W} = \{\mathcal{W}_1, \dots, \mathcal{W}_K\}$ trained on the loss *L* (e.g., cross-entropy), we formulate the following rank and scheme selection problem:

$$\min_{\mathbf{W}, \Theta, \mathbf{r}, \mathbf{s}} \quad L(\mathbf{W}) + \lambda C(\Theta, \mathbf{r})$$
s.t. $\Theta_k = \mathcal{R}(\mathcal{W}_k, s_k),$

$$\operatorname{rank}(\Theta_k) = r_k, \quad \forall k = 1, \dots, K$$

$$(3)$$

Here, the term $\lambda C(\Theta, \mathbf{r})$ controls the amount of compression and can target a specific cost of interest like FLOPs or storage.

Low-rank with scheme selection: Deriving the L and C steps Let us apply a penalty method and obtain an equivalent formulation (with $\mu \rightarrow \infty$):

$$\min_{\mathbf{W},\mathbf{\Theta},\mathbf{r},\mathbf{s}} L(\mathbf{W}) + \lambda C(\mathbf{\Theta},\mathbf{r}) + \frac{\mu}{2} \sum_{k=1}^{K} \|\mathbf{\Theta}_{k} - \mathcal{R}(\mathcal{W}_{k},s_{k})\|_{F}^{2}$$
s.t. rank $(\mathbf{\Theta}_{k}) = r_{k}, \quad \forall k = 1, \dots, K.$
(4)

Apply alternating optimization over variables W and $\{\Theta, \mathbf{r}, \mathbf{s}\}$:

• The step over W, which we call a learning (L) step, has the form of:

$$\min_{\mathbf{W}} \quad L(\mathbf{W}) + \frac{\mu}{2} \sum_{k=1}^{K} \left\| \boldsymbol{\Theta}_k - \mathcal{R}(\mathcal{W}_k, s_k) \right\|_F^2.$$

• The step over $\{\Theta, \mathbf{r}, \mathbf{s}\}$, which we call a compression (C) step, has the form of:

$$\min_{\boldsymbol{\Theta}, \mathbf{r}, \mathbf{s}} \quad \lambda C(\boldsymbol{\Theta}, \mathbf{r}) + \frac{\mu}{2} \sum_{k=1}^{K} \left\| \boldsymbol{\Theta}_{k} - \mathcal{R}(\mathcal{W}_{k}, s_{k}) \right\|_{F}^{2}$$

Low-rank with scheme selection: Solution of the C step

Due to the layerwise separability of cost function, the C-step problem separates over the layers into K smaller problems:

$$\min_{\boldsymbol{\Theta}_{k}, r_{k}, s_{k}} \quad \lambda C(\boldsymbol{\Theta}_{k}, r_{k}) + \frac{\mu}{2} \|\boldsymbol{\Theta}_{k} - \mathcal{R}(\mathcal{W}_{k}, s_{k})\|_{F}^{2}$$
s.t. rank $(\boldsymbol{\Theta}_{k}) = r_{k}.$

$$(5)$$

Solution:

- For a fixed scheme s_k the solution is known in closed form for multiple costs C
 - For storage and FLOPs, the solution involves SVD and enumeration [8]
 - For nuclear-norm cost, the solution involves singular value shrinkage [1]
- Therefore, to find global solution, we iterate over possible schemes and re-use steps for fixed scheme.

Low-rank with scheme selection: Experiments LeNet5 on MNIST VGG16 on CIFAR10 2 scheme 1 scheme 1 ---- scheme 2 scheme 2 7.5 scheme 3 ---- scheme 3 1.5 % 1.5 % -- ours -- ours 7 6.5 **R** [29] [16] [<mark>6</mark>] 6 0.5 0.5 50 100 200 0 150**MFLOPs MFLOPs** Insights:

- · selecting the scheme along with ranks outperforms rank selection
- the scheme 2 is selected more often (as optimal) per each layer

The software

The LC algorithm: pseudocode and software implementation

```
class LCAlgorithm():
    # Housekeeping code skipped
    def run(self):
        self.mu = 0
        self.c_step(step_number=0)
        for i, mu in enumerate(self.mu_schedule):
            self.mu = mu
            self.l_step(i) # user defined
        self.l_step(i) # user defined
        self.c_step(i) # library call
        self.multipliers_step()
```

The LC software

- Written in python using NumPy and PyTorch
- L step We hand off the L step to the user through the lambda functions.

```
def my_l_step(model, lc_penalty, args**):
    # ...
    loss = model.loss(out_, target_) + lc_penalty()
    loss.backward()
    optimizer.step()
    # ...
```

• C step Many compression are implemented, and you can chose any. If desired, you can add your own compression too by extending the CompressionTypeBase: class.

```
class ScaledBinaryQuantization(CompressionTypeBase):
    def compress(self, data):
        a = np.mean(np.abs(data))
        quantized = 2 * a * (data > 0) - a
        return quantized
```

Туре	Forms
Quantization	Adaptive Quantization into $\{c_1, c_2, \dots c_K\}$ Binarization into $\{-1, 1\}$ and $\{-c, c\}$ Ternarization into $\{-c, 0, c\}$
Pruning	$ \begin{array}{l} \ell_0 \text{-constraint (s.t., } \ \mathbf{w}\ _0 \leq \kappa) \\ \ell_1 \text{-constraint (s.t., } \ \mathbf{w}\ _0 \leq \kappa) \\ \ell_0 \text{-penalty } (\alpha \ \mathbf{w}\ _0) \\ \ell_1 \text{-penalty } (\alpha \ \mathbf{w}\ _1) \end{array} $
Low-rank	Low-rank compression to a given rank Low-rank with <i>automatic</i> rank selection for FLOPs reduction Low-rank with <i>automatic</i> rank selection for storage compression Low-rank with <i>device-targeted</i> compression
Additive Combinations	Quantization + Pruning Quantization + Low-rank Pruning + Low-rank Quantization + Pruning + Low-rank

The LC software: currently implemented compressions

The LC software: the ease of exploration

• Mix-and-match through compression tasks. We structured the software in such way that any compression can be applied to any part of the model, and you can mix them as well!

For example, the following semantics:

Translates into the following code:

```
from lc.torch import ParameterTorch as P, AsVector, AsIs
compression_tasks = {
    P([11.weight, 13.weight]): (AsVector, AdaptiveQuantization(k=6)),
    P(12.weight): (AsIs, LowRank(target_rank=3))
}
```

Example: Apples-to-apples comparison between compressions



See more on it in [9]

Example: Additive compressions to achieve smallest AlexNets

The frameworks and software allows easy exploration of new compressions. For example, how about additive combination of quantization and pruning?



Conclusion and future work

- We have presented the Learning-Compression algorithm: and showed how it can be applied to the problems low-rank compression for NNs.
- The resulting algorithms require alternation of two steps: learning (L) step involving model loss and dataset and compression (C) step involving actual compression problem.
- Our trained models achieve considerable amount of compression and are on-par or better than other approaches.
- During the course of my PhD studies, we have extended the LC framework to include following NN compressions:
 - quantization (arxiv, CVPR2021)
 - pruning (CVPR2018)
 - low-rank (CVPR2020, DCC2021, ICASSP2021, ICIP2021)
 - additive combinations of above (ECML2021)
 - comparison of pruning, quantization, and low-rank (IJCNN2021)
- Many extensions are possible ...

Future work: New models

My work is primarily about neural networks for vision, but, other models are of interest:

- language models like LSTMs, Transformers:
 - a particular interest is in compression of (tall) softmax layers
 - heteregenous structure of such models imply selection of "best" compressions
- recommendation systems

General ML models:

- probabilistic models (e.g., MMs and HMMs)
- kernel methods (for both classification and regression)

A particular challenge with these are a special constraints: e.g., with HMM's Gaussian emission matrices, weight must be psd

Future work: Cost-based selection mechanisms

We have demonstrated the power of compression cost function $C(\mathbf{W})$ that allows to select best ranks for size, FLOPs and runtime.

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

There are a lot of extensions on this front:

- Design of new cost-functions targeting:
 - power-consumption in the data center/mobile-device/etc
 - layout and placement of model (e.g., store W or W^T, run everything on CPU or GPU?)
 - utilization (cache misses, transfer speeds, loading times)
- Hardware-Compression co-design:
 - design hardware features optimized for specific compressions/models

Future work: Single entry for the world of model compression

Model compression is hard, but our LC framework, its extensions and software based on it allow us to drastically simplify it for users:

- Users only need to provide an L-step implementation (SGD-based) and choose required compressions from the library
- Not harder than training the original model

Due to its simplicity and competitive results, we believe it can be a single entry point to model compression:

• as PyTorch or TensorFlow but for model compression

To do so, the following is required:

- include new models/selection mechanisms/compressions/results
- advertise/popularize among community members
- include out-of-the-box support for "standard" features (framework bindings, pipelines, etc)
- create tutorials/workshops
- look for an industry support

The end. Thank you!

References

- [1] J.-F. Cai, E. J. Candès, and Z. Shen. A singular value thresholding algorithm for matrix completion. SIAM J. Optimization, 20(4):1956–1982, 2010.
- [2] Y. Choi, M. El-Khamy, and J. Lee. Towards the limit of network quantization. In Proc. of the 5th Int. Conf. Learning Representations (ICLR 2017)
- [3] X. Ding, G. Ding, Y. Guo, J. Han, and C. Yan. Approximated oracle filter pruning for destructive CNN width optimization. In Proc. of the 36th Int. Conf. Machine Learning (ICML 2019)
- [4] S. Han, H. Mao, and W. J. Dally. Deep compression: Compressing deep neural networks with pruning, trained quantization and Huffman coding. In Proc. of the 4th Int. Conf. Learning Representations (ICLR 2016)
- [5] Y. He, X. Zhang, and J. Sun. Channel pruning for accelerating very deep neural networks. In Proc. 16th Int. Conf. Computer Vision (ICCV'17)
- [6] Y. He, P. Liu, Z. Wang, Z. Hu, and Y. Yang. Filter pruning via geometric median for deep convolutional neural networks acceleration. In Proc. of the 2019 IEEE Conf. Computer Vision and Pattern Recognition (CVPR'19)
- [7] Z. Huang and N. Wang. Data-driven sparse structure selection for deep neural networks. In Proc. 15th European Conf. Computer Vision (ECCV'18)
- [8] Y. Idelbayev and M. Á. Carreira-Perpiñán. Low-rank compression of neural nets: Learning the rank of each layer. In Proc. of the 2020 IEEE Computer Society Conf. Computer Vision and Pattern Recognition (CVPR'20)
- [9] Y. Idelbayev and M. Á. Carreira-Perpiñán. An empirical comparison of quantization, pruning and low-rank neural network compression using the LC toolkit. In Int. J. Conf. Neural Networks (IJCNN'21).
- [10] Y. Idelbayev and M. Á. Carreira-Perpiñán. More general and effective model compression via an additive combination of compressions. In Proc. of the 32nd European Conf. Machine Learning (ECML-21)
- [11] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell. Caffe: Convolutional architecture for fast feature embedding. arXiv:1408.5093 [cs.CV], June 20 2014.
- [12] H. Kim, M. U. K. Khan, and C.-M. Kyung. Efficient neural network compression. In Proc. of the 2019 IEEE Computer Society Conf. Computer Vision and Pattern Recognition (CVPR'19)
- [13] Y.-D. Kim, E. Park, S. Yoo, T. Choi, L. Yang, and D. Shin. Compression of deep convolutional neural networks for fast and low power mobile applications. In Proc. of the 4th Int. Conf. Learning Representations (ICLR 2016)
- [14] C. Leng, H. Li, S. Zhu, and R. Jin. Extremely low bit neural network: Squeeze the last bit out with ADMM. In Proc. of the 32nd AAAI Conference on Artificial Intelligence (AAAI 2018)
- [15] C. Li and C. J. R. Shi. Constrained optimization based low-rank approximation of deep neural networks. In Proc. 15th European Conf. Computer Vision (ECCV'18)
- [16] H. Li, A. Kadav, I. Durdanovic, and H. P. Graf. Pruning filters for efficient ConvNets. In Proc. of the 5th Int. Conf. Learning Representations (ICLR 2017)
- [17] J. Li, Q. Qi, J. Wang, C. Ge, Y. Li, Z. Yue, and H. Sun. OICSR: Out-in-channel sparsity regularization for compact deep neural networks. In 2019 IEEE Conf. Computer Vision and Pattern Recognition (CVPR'19)
- [18] M. Lin, R. Ji, Y. Wang, Y. Zhang, B. Zhang, Y. Tian, and L. Shao. HRank: Filter pruning using high-rank feature map. In Proc. of the 2020 IEEE Computer Society Conf. Computer Vision and Pattern Recognition (CVPR'20)
- [19] S. Lin, R. Ji, C. Yan, B. Zhang, L. Cao, O. Ye, F. Huang, and D. Doermann. Towards optimal structured CNN pruning via generative adversarial learning. In Proc. of the 2019 IEEE Computer Society Conf. Computer Vision and Pattern Recognition (CVPR'19)
- [20] M. Rastegari, V. Ordonez, J. Redmon, and A. Farhadi. XNOR-net: ImageNet classification using binary convolutional neural networks. In Proc. 14th European Cont. Computer Vision (ECCV'16)
- [21] C. Tai, T. Xiao, Y. Zhang, X. Wang, and W. E. Convolutional neural networks with low-rank regularization. In Proc. of the 4th Int. Conf. Learning Representations (ICLR 2016)
- [22] F. Tung and G. Mori. CLIP-Q: Deep network compression learning by in-parallel pruning-quantization. In Proc. of the 2018 IEEE Computer Society Conf. Computer Vision and Pattern Recognition (CVPR'18)
- [23] W. Wen, C. Xu, C. Wu, Y. Wang, Y. Chen, and H. Li. Coordinating filters for faster deep neural networks. In Proc. 16th Int. Conf. Computer Vision (ICCV'17)
- [24] J. Wu, C. Leng, Y. Wang, Q. Hu, and J. Cheng. Quantized convolutional neural networks for mobile devices. In Proc. of the 2016 IEEE Computer Society Conf. Computer Vision and Pattern Recognition (CVPR'16)
- [25] X. Xu, Y. Ding, S. X. Hu, M. Niemier, J. Cong, Y. Hu, and Y. Shi. Scaling for edge inference of deep neural networks. Nature Electronics, 1:216–222, Apr. 17 2018.
- [26] H. Yang, S. Gui, Y. Zhu, and J. Liu. Automatic neural network compression by sparsity-quantization joint learning: A constrained optimization-based approach. In Proc. of the 2020 IEEE Computer Society Cont. Computer Vision and Pattern Recognition (CVPR20)
- [27] J. Ye, X. Lu, Z. Lin, and J. Wang. Rethinking the smaller-norm-less-informative assumption in channel pruning of convolution layers. In Proc. of the 6th Int. Conf. Learning Representations (ICLR 2018)
- [28] R. Yu, A. Li, C.-F. Chen, J.-H. Lai, V. I. Morariu, X. Han, M. Gao, C.-Y. Lin, and L. S. Davis. NISP: Pruning networks using neuron importance score propagation. In Proc. of the 2018 IEEE Computer Society Conf. Computer Vision and Pattern Recognition (CVPR'18)
- [29] C. Zhao, B. Ni, J. Zhang, Q. Zhao, W. Zhang, and Q. Tian. Variational convolutional neural network pruning. In Proc. of the 2019 IEEE Computer Society Conf. Computer Vision and Pattern Recognition (CVPR'19)
- [30] S. Zhou, Z. Ni, X. Zhou, H. Wen, Y. Wu, and Y. Zou. Dorefa-net: Training low bitwidth convolutional neural networks with low bitwidth gradients. arXiv:1606.06160, July 17 2016.
- [31] Z. Zhuang, M. Tan, B. Zhuang, J. Liu, Y. Guo, Q. Wu, J. Huang, and J. Zhu. Discrimination-aware channel pruning for deep neural networks. In Advances in Neural Information Processing Systems (NEURIPS), volume 31, pages 815–861.

Additional slides

Easy exploration of compressions

Having an L-step implementation (you only need one), definition of compression is very simple:

quantize each layer with separate codebooks

```
compression_tasks = {
  Param(l1.weight): (AsVector, AdaptiveQuantization(k=2)),
  Param(l2.weight): (AsVector, AdaptiveQuantization(k=2)),
  Param(l3.weight): (AsVector, AdaptiveQuantization(k=2))
}
```

	<pre>compression_tasks = { Param([l1.weight, l2.weight, l3.weights]):</pre>
prune all but 5%	(AsVector, ConstraintLOPruning(kappa=13310)) # 13310 = 5%
	}

prune first layer, low-rank to second, quantize third

```
compression_tasks = {
  Param(l1.weight): (AsVector, ConstraintLOPruning(kappa
  =5000)),
  Param(l2.weight): (AsIs, LowRank(target_rank=10))
  Param(l3.weight): (AsVector, AdaptiveQuantization(k=2))
}
```

Low-rank: experiments on CIFAR10 (selected ranks)



Ranks and FLOPs over layers for VGG16 and ResNet20 using our method for selected values of $\lambda~(\times 10^{-4})$

Low-rank with scheme selection: experiments on CIFAR10 (selected ranks & schemes)

The final selected ranks and reshaping schemes for some of the compressed VGG16 on CIFAR10																			
	test error	FI OPs	narameters	The selected scheme and rank over layers															
		I LOI 3	parameters	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
$\lambda = 2.0 \times 10^{-5}$	5.90%	156M	4.8M	${\mathcal S}_1$ 16	${\mathcal S}_2 \ {\mathbf 32}$	${\mathcal S}_2 \ 71$	${\mathcal S}_2 onumber 97$	<i>S</i> 3 116	${\mathcal S}_2$ 238	${\mathcal S}_2$ 263	S ₃ 254	${\mathcal S}_2$ 292	${\mathcal S}_2$ 172	${\mathcal S}_2$ 122	${\mathcal S}_2 onumber 99$	${\mathcal S}_2$ 105	31	11	9
$\lambda = 7.5 \times 10^{-5}$	5.97%	78M	2.5M	$rac{{\mathcal S}_1}{{ extsf{15}}}$	${\mathcal S}_2$ 20	$rac{{\mathcal S}_2}{43}$	${\mathcal S}_2 \ 53$	$\frac{\mathcal{S}_2}{113}$	$\frac{\mathcal{S}_2}{110}$	$\frac{\mathcal{S}_2}{116}$	$rac{{\mathcal S}_2}{239}$	$\frac{\mathcal{S}_2}{124}$	$rac{{\cal S}_2}{79}$	${\mathcal S}_2 \ {\bf 72}$	$\overline{\mathcal{S}_2}$ 74	$rac{{\mathcal S}_2}{89}$	23	10	9

Quantization, pruning, and low-rank: Comparison on MNIST



Quantization, pruning, and low-rank: Comparison on CIFAR10



Why not to train smaller model in the first place?

We certainly can train smaller model (with fewer weights or layers) from scratch using SGD. However, it requires many trials to find best small network architecture.

Another option: if model (loss L and model function f) and decompression mapping are both differentiable, we can directly optimize:

 $\min_{\boldsymbol{\varTheta}} L(\boldsymbol{\Delta}(\boldsymbol{\varTheta}))$

- Unfortunately, most of the compressions of interest are non differentiable
- Many tailored algorithms were proposed to bypass the non-differentiability of specific compression forms (BinaryConnect, STE)
- There is still the problem of selecting the best form of Δ (say, what should be the ranks per layer?)

Generally, training a large model and then compressing it allows to peek into upper bound of the network performance, and adjust expectations wrt smaller, compressed model.

Convergence of the LC algorithm

The framework includes a large variety of compression problems that can be continuous (low-rank with given rank, ℓ_1 pruning) or combinatorial (rank selection, quantizaiton, ℓ_0 pruning)

Convergence of the LC alternation to a local minimizer can be established for:

- convex, differentiable cases (quite generally)
- nonconvex, fully differential cases (more restrictive results)

However, most interesting cases are those with nonconvex losses and non-differentiable constraints (quantization, ℓ_0 pruning). For these problems:

- convergence results are of general interest and is an open problem
- have flavor of *k*-means: algorithm stops at a point of no further improvement
- empirically, algorithm performs very well