# Model compression as constrained optimization, with application to neural nets

**Miguel Á. Carreira-Perpiñán** and **Yerlan Idelbayev**

Electrical Engineering and Computer Science

University of California, Merced

`http://eecs.ucmerced.edu`

# Why neural net compression?

❖ Deep learning very successful in some practical applications
computer vision, speech, NLP, etc.

❖ Very good classification accuracy if training large, deep neural nets on large datasets (with several GPUs over several days. . . ).

❖ It is of interest to deploy high-accuracy deep nets on limited-computation devices (mobile phones, IoT, etc.), but this requires compressing the deep net to meet the device's limitations in memory, runtime, energy, bandwidth, etc.

❖ Surprisingly high compression rates often possible because deep nets are vastly over-parameterized in practice.
It seems this makes it easier to learn a high-accuracy net.

❖ Still, lossy compression will degrade the accuracy, so we want to compress optimally.

# Related work on neural net compression

Many papers in the last 2–3 years, although neural net compression was already studied in the 1980–90s. Some common approaches:

- ❖ Direct compression: train a reference net, compress its weights using standard 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 backprop training.

  - ✦ remove neuron/weight values with low magnitude on the fly
  - ✦ binarize or round weight values in the backprop forward pass

  Often no convergence guarantees.

- ❖ Other ad-hoc algorithms proposed for specific compression techniques.

- ❖ Multiple compression techniques can be combined.

# 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 <span style="color:yellow">lowest loss possible for a given compression technique and compression rate</span>.

Then, we would like an algorithm that:

- ❖ is generic, applicable to many compression techniques
  so the user can easily try different types of compression

- ❖ has optimality guarantees inasmuch as possible

- ❖ is easy to integrate in existing deep learning toolboxes

- ❖ is efficient in training.

# An illustrative example: low-rank compression

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

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

We want $\mathbf{W} = \mathbf{U}\mathbf{V}^T$ with $\mathbf{U}_{d \times r}$, $\mathbf{V}_{D \times r}$ and $r < \min(d, D)$. Hence the problem is (RRR, reduced-rank regression):

$$\min_{\mathbf{U}, \mathbf{V}} L(\mathbf{U}, \mathbf{V}) = \sum_{n=1}^{N} \|\mathbf{y}_n - \mathbf{U}\mathbf{V}^T\mathbf{x}_n\|^2 \quad \left(\begin{smallmatrix}\text{and orthog. constr.} \\ \text{on } \mathbf{U} \text{ or } \mathbf{V}\end{smallmatrix}\right)$$

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

- ❖ using gradient-based optimization
- ❖ using alternating optimization over $\mathbf{U}$ and $\mathbf{V}$
- ❖ in fact, there is a closed-form solution as an eigenproblem.

But we want a more generally applicable algorithm. Instead, introduce auxiliary variables $\mathbf{W} = \mathbf{U}\mathbf{V}^T$ and define as a constrained problem:

$$\min_{\mathbf{W}, \mathbf{U}, \mathbf{V}} L(\mathbf{W}) = \sum_{n=1}^{N} \|\mathbf{y}_n - \mathbf{W}\mathbf{x}_n\|^2 \text{ s.t. } \mathbf{W} = \mathbf{U}\mathbf{V}^T \quad \left(\begin{smallmatrix}\text{and orthog. constr.} \\ \text{on } \mathbf{U} \text{ or } \mathbf{V}\end{smallmatrix}\right)$$

This separates the loss part from the compression part.

# Model compression as constrained optimization

General formulation:

uncompressed weights
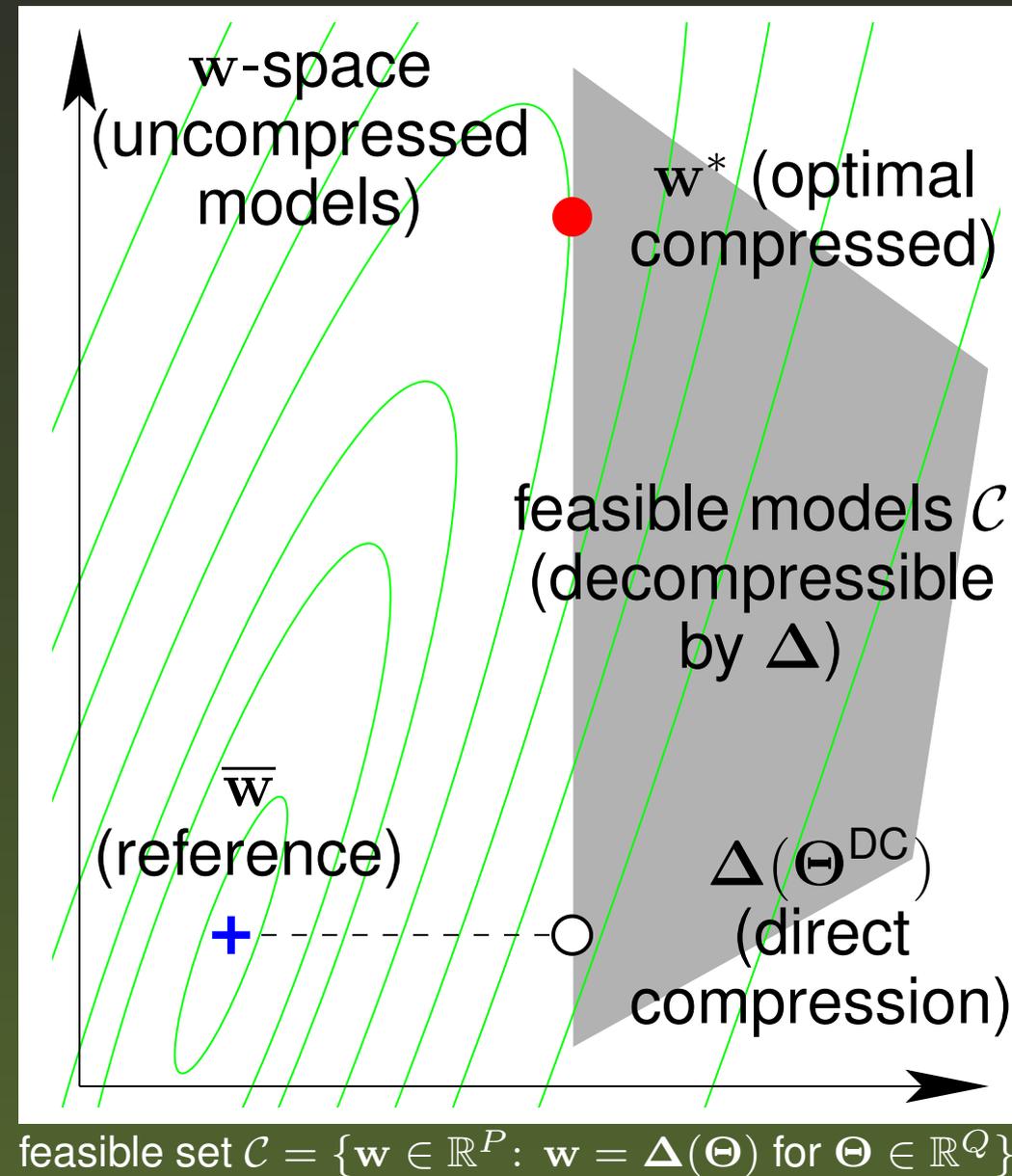
low-dim. params

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

task loss

decompression mapping

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 $\Delta(\Theta)$.

$\mathbf{w}$-space (uncompressed models)

$\mathbf{w}^*$ (optimal compressed)

feasible models $\mathcal{C}$ (decompressible by $\Delta$)

$\overline{\mathbf{w}}$ (reference)

$\Delta(\Theta^{\text{DC}})$ (direct compression)

feasible set $\mathcal{C} = \{\mathbf{w} \in \mathbb{R}^P : \mathbf{w} = \Delta(\Theta) \text{ for } \Theta \in \mathbb{R}^Q\}$

# The "learning-compression" (LC) algorithm

Use a penalty method (e.g. quadratic penalty): as $\mu \to \infty$, minimize:

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

using alternating optimization over $\mathbf{w}$ and $\boldsymbol{\Theta}$:

- ❖ L step: $\min_{\mathbf{w}} L(\mathbf{w}) + \frac{\mu}{2}\|\mathbf{w} - \boldsymbol{\Delta}(\boldsymbol{\Theta})\|^2$
  independent of the compression technique

- ❖ C step: $\min_{\boldsymbol{\Theta}} \|\mathbf{w} - \boldsymbol{\Delta}(\boldsymbol{\Theta})\|^2$
  independent of the loss and dataset.

This algorithm converges to a local stationary point of the constrained problem under standard assumptions as $\mu \to \infty$

- ❖ smooth loss $L(\mathbf{w})$ and decompression mapping $\boldsymbol{\Delta}(\boldsymbol{\Theta})$
- ❖ sufficiently accurate optimizations on L and C steps.

Generic: one algorithm, many compression techniques.

# LC algorithm: solving the L step

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

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

❖ 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 $\mu \to \infty$.

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 $\mathbf{w}$:

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

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

❖ Equivalent to orthogonal projection $\boldsymbol{\Theta} = \boldsymbol{\Pi}(\mathbf{w})$ on the feasible set.
Find the closest compressed model to $\mathbf{w}$.

❖ The solution depends on the choice of the decompression mapping $\boldsymbol{\Delta}$, and is known for many common compression techniques:

✦ low rank: SVD

✦ binarization: binarize weights

✦ quantization in general: $k$-means

✦ pruning: zero all but top weights

✦ etc.

So the C step simply requires calling the corresponding compression subroutine as a black box.

# Example C step: quantization with adaptive codebook

The C step solution is derived from the form of the decompression mapping $\Delta\colon \Theta \to \mathbf{w} \in \mathbb{R}^P$.

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

❖ Compressed net: codebook ($K$ floats) + $\lceil \log_2 K \rceil$ bits per weight.

❖ The decompression mapping $\Delta$ is a table lookup $w_i = c_{\kappa(i)}$ where $\kappa\colon \{1, \dots, P\} \to \{1, \dots, K\}$ is a discrete mapping that assigns each weight to one codebook entry.

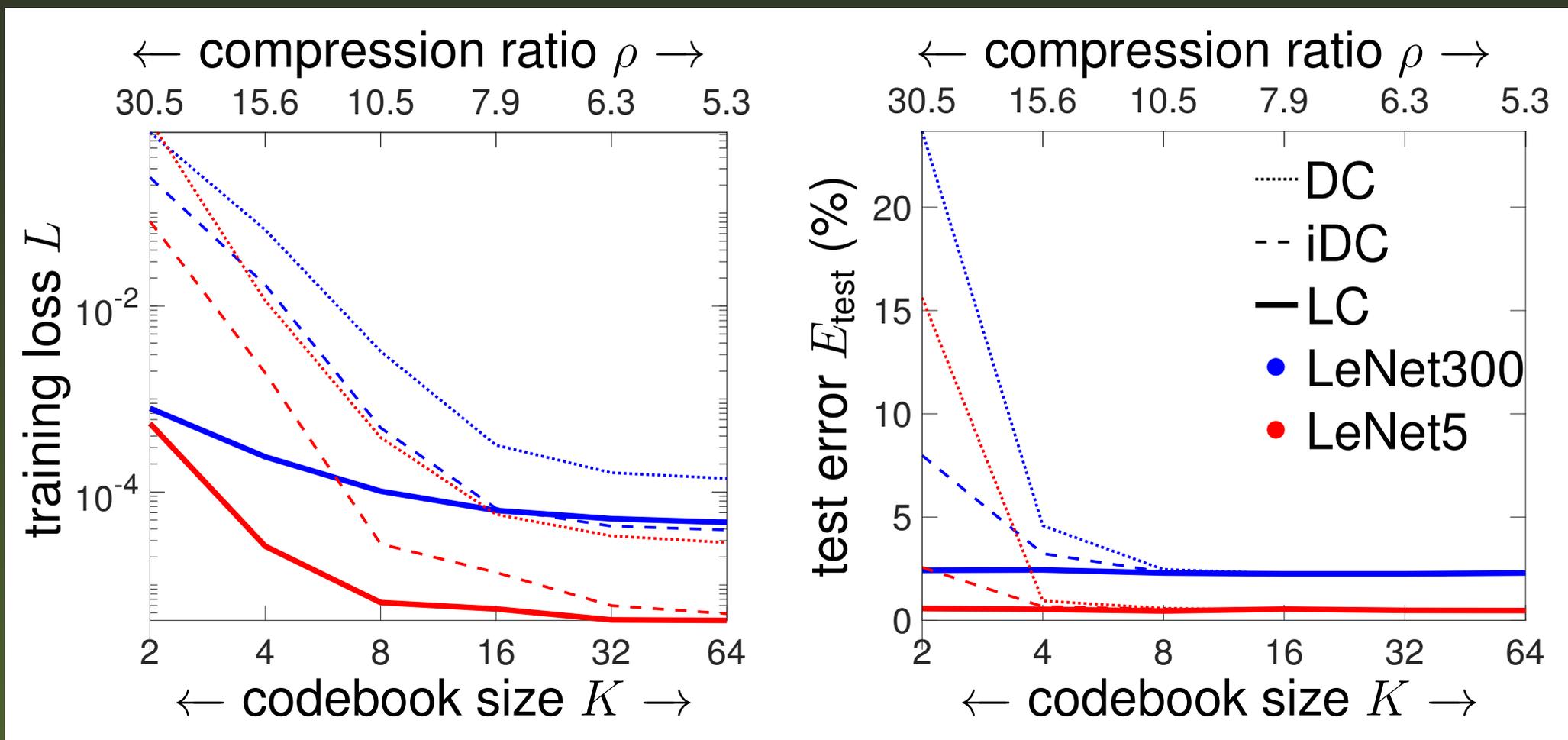❖ Rewriting $\kappa()$ using binary assignment variables $\mathbf{Z} \in \{0, 1\}^{P \times K}$:

$$\min_{\mathcal{C}, \kappa} \sum_{p=1}^{P} \left\| w_i - c_{\kappa(i)} \right\|^2 \Leftrightarrow \min_{\mathcal{C}, \mathbf{Z}} \sum_{i,k=1}^{P,K} z_{ik} \| w_i - c_k \|^2 \text{ s.t. } \begin{cases} \sum_{k=1}^{K} z_{ik} = 1 \\ i = 1, \dots, P \end{cases}$$

❖ So $\Theta = \{\mathcal{C}, \mathbf{Z}\}$, and the decompression mapping is $w_i = \sum_{k=1}^{K} z_{ik} c_k$.

❖ This is the squared distortion problem. It is NP-complete.

❖ Approximate solution: $k$-means.

# Experimental results: quantization

## 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:



VGG-like net (14M parameters, 12 layers) on CIFAR10 using $K = 2$ (compression ratio $\times 32$). Error: reference 13.15%, compressed 13.05%.

# 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:
  - ✦ 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 task 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.

- ❖ Guaranteed to converge to local optimum under std assumptions.

- ❖ Easy to implement in deep learning toolboxes.

- ❖ Very effective in practice.

We are developing this framework for number of compression types.