Low-parametric deep learning



Francesco Tudisco GSSI – Gran Sasso Science Institute

Data-driven paradigm

• Input

- set of observations $(x, y) \in \mathbb{R}^d \times \mathbb{R}^k$

– underlying model M

• Goal

infer the map $f: x \mapsto y$ in a **fast** and **reliable** way

Supervised classification

$$x =$$

Dog 0.43 Mop 0.56 Pole 0.01

Prediction & recommendation





Scientific computing

Field measurements

 $\begin{cases} \mathcal{F}(f(x);\theta) = a(x) \\ \mathcal{B}(f(x)) = b(x) \end{cases}$



Artificial Neural Network ansatz To compute $f: x \mapsto y$ we need an *ansatz* function space

Artificial Neural Networks (NNs) are the space of choice in most applications



Feed-forward Neural Nets

Parametric family of functions that we can write as

$$f(x; W) = \{\sigma_{\ell} \circ A_{\ell} \circ \cdots \circ \sigma_{1} \circ A_{1}\}(x)$$

where:

- $W = (W_1, ..., W_\ell)$ are the parameters
- $A_i : \mathbb{R}^{N_i} \to \mathbb{R}^{N_{i+1}}$ are simple mappings, for instance
 - MatMult $A_i(x) = W_i x$
 - Convolution $A_i(x) = W_i * x$
 - Dot product $A_i(x) = W_i x (W_i x)^T$
- σ_i are entrywise nonlinear activation functions

Each time we do inference on *x*

Memory and computation footprints This requires

N

N.

We need to compute f(x)

- To evaluate ℓ mappings $A_i \approx \sum_i N_i N_{i+1}$
- To evaluate ℓ activations $\sigma_i \approx \sum_i N_i$
- To store the weights $W_i \sim N_{i+1} \times N_i \approx \sum_i N_i N_{i+1}$

 $O(\ell N^2)$ $O(\ell N^2)$ Operations to make Variables to store

Examples

AlexNet ≈ 62 million parameters VGG16 ≈ 135 million parameters ResNet50 ≈ 23 million parameters M-BERT ≈ 4 billion parameters DALL-E ≈ 3.5 billion parameters

Prohibitive for **online learning** and for **limited-resource devices**, e.g. smartphones, satellites, drones, etc.



Not all weights are equally important



Model compression, aka **pruning** Constrain the parameter space

Most common examples:

- sparsity: $nnz(W_i) = O(N)$
- quantization: $(W_i)_{ij} \in \{-1, 0, +1\}$
- low-rank: $rank(W_i) = r_i \ll N$





Lottery ticket hypothesis

The lottery ticket hypothesis: finding sparse, trainable neural networks, J Frankle, M Carbin, ICLR 2019

Dense NNs contain subnetworks that, when trained in isolation and for the same number of epochs, can match the accuracy of the original full net.



Performance of different pruning strategies as the compression rate increases on LeNet5/MNIST $98\% \rightarrow 97\%$

In practice



Is this happening by chance?

A network of depth ℓ can be approximated arbitrarily well by pruning a network of depth 2ℓ

Proving the lottery ticket hypothesis, E Malach et al, ICML20 For any $\varepsilon > 0$, and any neural network f(x; W) of depth ℓ and width N, there exist:

- a larger network g of depth 2ℓ and width $poly(\ell, N, \varepsilon^{-1})$,
- and a subnetwork \tilde{f} with only $O(\ell N^2)$ parameters

such that $|\tilde{f} - f| \leq \varepsilon$ with probability at least $(1 - \varepsilon)$.

Problem(s): the result is not constructive and no proof of the existence of winning tickets

How about the training phase?

This approach completely **ignores training costs**

bass	<u> </u>		
Forward	Evaluate $f(W, x_b)$	Evaluate $\nabla_W L(W, x_b)$	Backward pa
			SSI

Training costs

Each iteration requires $O(2\ell N^2)$

Proposed approach

Manifold training



Low-rank matrix manifold

- $\mathcal{M} = \{W: \operatorname{rank}(W) = r \ll N\}$
- $W \in \mathcal{M}$ can be written as:



Memory requirement O(2Nr)Operation cost O(2Nr)

Manifold SGD



- 1. Computing $\mathcal{P}_{\mathcal{M}}(W \lambda \nabla_W L(W))$ can be expensive
- 2. How can we find "the right \mathcal{M} "?

DLRT algorithm's key ingredients: – efficient – adaptive 1. Compute $\mathcal{P}_{\mathcal{M}}(W - \lambda \nabla_W L(W))$ using only the three factors U, S, V

• $U \leftarrow U - \lambda \nabla_U L(U; S, V, x)$

$$V \leftarrow V - \lambda \nabla_V L(V; U, S, x)$$

- $S \leftarrow S \lambda \nabla_S L(S; U, V, x)$
- 2. Adjust the rank of *S* to match a desired compression rate

 $\vartheta = \frac{\# \text{ final parameters}}{\# \text{ initial parameters}}$



 $\vartheta = \text{compression rate}$

 $\widetilde{W}^{(n)} = U^{(n)}S^{(n)}V^{(n)}$ comp low-rank weights If the loss is smooth enough and its gradient is " ε -far" from the low-rank manifold, then

$$f(x; \widetilde{W}^{(n)}) - f(x; W^{(n)}) \le O(\varepsilon + \vartheta)$$

 $W^{(n)} =$ ideal full weights

Fully-connected on MNIST

Rank Evolution



Fully-connected on MNIST

Timings & Accuracy



LeNet5 on MNIST

method	mean test acc.	ranks	params	c.r.	params	c.r.
LeNet5	$\mathbf{98.8\%} \pm 0.06$	$\left[20, 50, 500, 10 ight]$	430500	0%	861000	0%
$\vartheta = 0.09$	$98.2\%\pm0.26$	$\left[10, 23, 62, 10\right]$	37445	$90.9\%\pm0.3$	532176	$35.5\%\pm1.8$
$\vartheta = 0.11$	$98.2\%\pm0.44$	$\left[10, 20, 48, 10\right]$	30278	$93.1\%\pm0.45$	412898	$53.3\%\pm3.5$
$\vartheta = 0.13$	$97.9\%\pm0.49$	$\left[9, 16, 37, 10\right]$	24542	$94.3\%\pm0.17$	316997	$63.2\%\pm1.1$
$arphi$ $\vartheta = 0.15$	$98.1\%\pm0.33$	$\left[9, 16, 28, 10\right]$	20033	$95.4\%\pm0.23$	251477	$71.4\%\pm1.83$
$\vartheta = 0.2$	$98.1\%\pm0.34$	$\left[8,8,15,10\right]$	13091	$96.9\% \pm 0.16$	135536	$83.4\%\pm1.21$
$\vartheta = 0.3$	$97.5\%\pm0.48$	$\left[4,6,8,10\right]$	9398	$97.9\%\pm0.08$	80792	$91.2\%\pm0.59$
$\vartheta = 0.4$	$96.0\%\pm0.94$	$\left[2,4,4,10\right]$	7250	$98.3\%\pm0.06$	47882	$94.4\%\pm0.3$
$\vartheta = 0.45$	$94.1\%\pm0.49$	$\left[2,2,3,10\right]$	6647	$\mathbf{98.4\%}\pm0.07$	35654	$\mathbf{95.4\%} \pm 0.4$
(SSL)(ft)	99.18%		110000	74.4%		< 0%
(NISP) (ft)	99.0%		100000	76.5%		< 0%
(GAL)	98.97%		30000	93.0%		< 0%
(LRNN)	98.67%	$\left[3,3,9,9\right]$	18075	95.8%		< 0%
(SVD prune)	94.0%	$\left[2,5,89,10\right]$	123646	71.2%		< 0%

Imagenet1K – 1million images & 1000 classes

Model	test acc. [%]	c.r. eval [%]	c.r. train [%]
ResNet-50	-0.56	54.1	14.2
VGG16	-2.19	86	78.4

Test accuracy w.r.t. full-rank baseline for $\vartheta = 0.1$



The bound does not depend on the singular values



Open Review

Low-rank lottery tickets @ NeurIPS 2022

Thank you!