Team#26 project
**"Adaptive Mixture of Low-Rank Factorizations for Compact Neural Modeling" (ICLR 2019)**

Yuriy Gabuev, Van Khachatryan, Stanislav Tsepa
Denis Zuenko, Aleksandr Rubashevskii

Skoltech
20 December, 2018

## Outline

- Key concepts

- Experiment reproduction

- Conclusion

- References

# Problem

- Modern NNs have large weight matrices, most are not suitable for mobile deployment

- Low-rank factorization is the popular instrument to reduce matrix size

- Large weight matrix W can be represented as a product of two small rank-d matrices

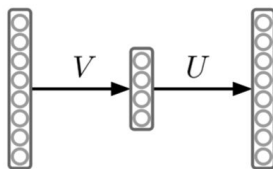$$W = UV^\top \qquad U \in \mathbb{R}^{m \times d}, V \in \mathbb{R}^{n \times d}$$

- Large complexity decrease: $O(d(m+n)) \ll O(mn)$

- Problem - loss of information when projecting onto low-dimensional space (Linear bottleneck)
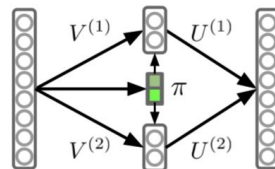
# Adaptive mixture of low-rank factorizations

- Make decompositions ~~great again!~~ data-dependent

$$W(h) = \sum_{k=1}^{K} \pi_k(h) U^{(k)} \big( V^{(k)} \big)^{\top} \qquad \pi(\cdot) \; : \; \mathbb{R}^n \; \to \; \mathbb{R}^K$$

- Replace large matvec with adaptive mixture of low-rank matvecs



(a) regular low-rank          (b) adaptive low-rank

- Strictly speaking, **not a decomposition**, but a new learnable module

- $\pi(\cdot)$ is a small non-linear data-dependent function, e.g. $\pi(h) = \sigma(P(\mathrm{pool}(h)), \; P \in \mathbf{R}^{K \times n_{\mathrm{pooled}}}$
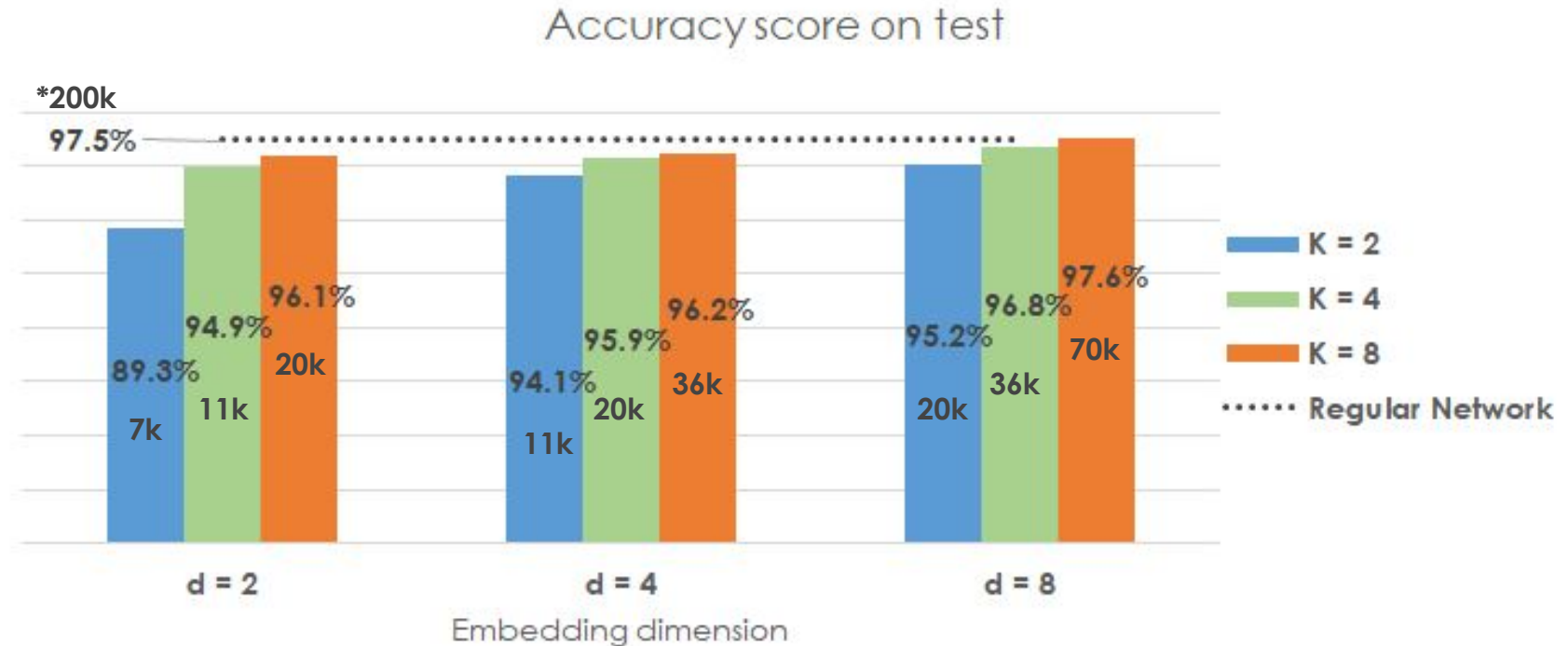
# Experiment reproduction

MLP (multi-layer perceptron)

- Digit recognition on MNIST dataset
- Simple one-layer MLP of 300 hidden units, input and output sizes are 784 and 10, respectively
- Rank-d matrices with d = 2, 4, 8
- K = 2, 4, 8
- Computed mixed weights with x (784 x 1) reduced to x (28 x 1)
- Accuracy of adaptive versions of low-rank factorization is 89-97.5%, depending on (K, d)

# Experiment reproduction

MLP (multi-layer perceptron): Results



Accuracy score on test

*200k

97.5%

89.3% — 7k, 94.9% — 11k, 96.1% — 20k (d = 2)

94.1% — 11k, 95.9% — 20k, 96.2% — 36k (d = 4)

95.2% — 20k, 96.8% — 36k, 97.6% — 70k (d = 8)

K = 2
K = 4
K = 8
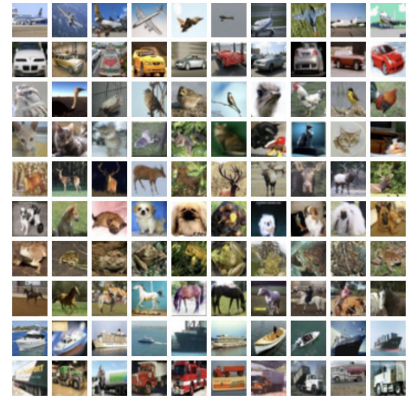Regular Network

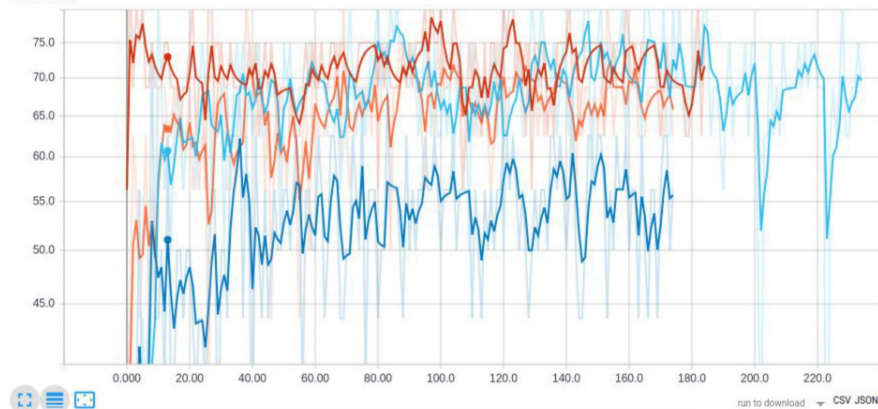Embedding dimension

*number of parameters

# Experiment reproduction

Convolutional Neural Networks

- CIFAR-10

- MobileNet, pointwise convolutions replaced with adaptive low-rank approximation

- Reducing of parametres -> reducing of accuracy

- The best accuracy for MobileNet is 96.875

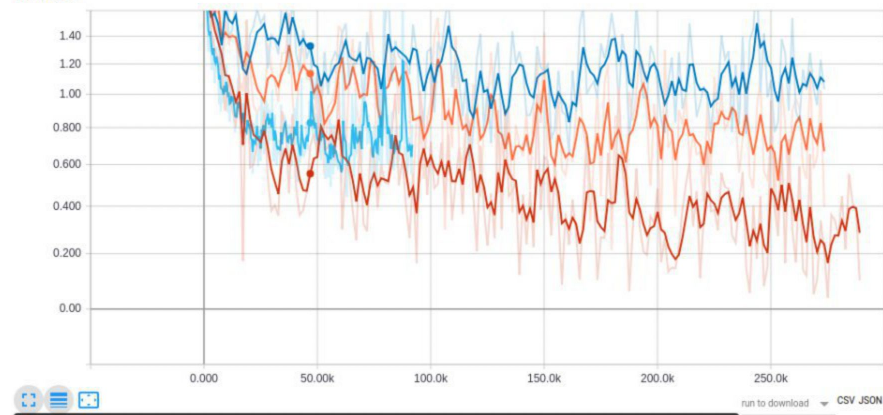- The best accuracy for MobileNet-low-rank is 90.625

Test/Acc_1

| Name | Smoothed | Value | Step | Time | Relative |
|---|---|---|---|---|---|
| MobileNet_1545271992 | 72.96 | 75.00 | 13.00 | Thu Dec 20, 05:38:21 | 21m 53s |
| MobileNet_LowRank_-d=2, K=8, pi_size=8- 1545271981 | 51.03 | 62.50 | 13.00 | Thu Dec 20, 05:39:27 | 23m 2s |
| MobileNet_LowRank_-d=4, K=4, pi_size=8-1545291806 | 60.72 | 62.50 | 13.00 | Thu Dec 20, 10:57:13 | 11m 50s |
| MobileNet_LowRank_-d=8, K=2, pi_size=8-1545271971 | 63.36 | 62.50 | 13.00 | Thu Dec 20, 05:39:19 | 23m 22s |

Train/Loss

| Name | Smoothed | Value | Step | Time | Relative |
|---|---|---|---|---|---|
| MobileNet_1545271992 | 0.5542 | 0.6986 | 46.89k | Thu Dec 20, 06:04:54 | 48m 30s |
| MobileNet_LowRank_-d=2, K=8, pi_size=8- 1545271981 | 1.327 | 1.328 | 46.89k | Thu Dec 20, 06:07:22 | 51m 11s |
| MobileNet_LowRank_-d=4, K=4, pi_size=8-1545291806 | 0.8276 | 0.8154 | 46.92k | Thu Dec 20, 12:33:00 | 1h 47m 41s |
| MobileNet_LowRank_-d=8, K=2, pi_size=8- 1545271971 | 1.135 | 1.096 | 46.89k | Thu Dec 20, 06:07:11 | 51m 31s |

| | best acc Top-1 | Params |
|---|---|---|
| MobileNet-CIFAR10 | 96.875 | 121.61k |
| MobileNet-CIFAR10 low-rank, (2,8) | 81.25 | 23.1k |
| MobileNet-CIFAR10 low-rank, (8,2) | 90.625 | 46.52k |
| MobileNet-CIFAR10 low-rank, (4,4) | 81.25 | 46.7k |

# Conclusion

- adaptive low-rank factorization is an original method proved to have results much better than regular low-rank decomposition
- it achieves up to 60-80% compression (dependent on the model) without significant decrease of accuracy
- in contrast to regular low-rank methods it learns non-linear low-rank manifolds due to learnable

$$\pi(\cdot)$$

# References

[1] T. Chen, J. Lin, T. Lin, C. Wang, D. Zhou, S. Han, Adaptive Mixture of Low-Rank Factorizations for Compact, 2018
[2] François Chollet. Xception: Deep learning with depth-wise separable convolutions. *arXiv preprint*, 2016.

# Code

https://github.com/zuenko/ALRF