

Supplementary Materials for Efficient Optimization Methods for Extreme Similarity Learning with Nonlinear Embeddings

Bowen Yuan
National Taiwan University
f03944049@csie.ntu.edu.tw

Pengrui Quan
University of California, Los Angeles
prquan@g.ucla.edu

Yu-Sheng Li
National Taiwan University
r07922087@csie.ntu.edu.tw

Chih-Jen Lin
National Taiwan University
cjlin@csie.ntu.edu.tw

1 MORE IMPLEMENTATION DETAILS

In evaluating $L(\theta)$, the required $P, Q, \hat{P}_c, \hat{Q}_c, P_c, Q_c$, and Frobenius inner products can be directly computed by the forward and other built-in functions. However, when we use GPUs for large data sets, we face two difficulties. First, as the memory consumption of a forward process is proportional to the number of data, calculating the whole P and Q at once may be infeasible for GPUs. Second, for later computations, we need to cache P and Q in $O((m+n)k)$ space, which may also be infeasible.

For the first difficulty, fortunately, as the computation of each row of P and Q is independent of others, we follow [2] to have a mini-batch setting. Specifically, by splitting \mathbb{U} and \mathbb{V} into multiple subsets, we sequentially calculate rows of P or Q corresponding to indices in each subset. For $\tilde{P}_c, \tilde{Q}_c, \hat{P}_c, \hat{Q}_c, P_c$, and Q_c , each of which is the sum of m or n matrices of size $k \times k$, we can calculate the partial results on each subset and accumulate them for the final output.

To overcome the second difficulty, we apply a heterogeneous computing scheme by storing P and Q in CPU, which has a larger memory capacity than GPU. After finishing the above computation in GPU on each subset, we move the resulting partial P and Q to the CPU memory. Some subsequent sparse operations such as $L^+(\theta)$ in (10) are conducted in CPU. As the advantage of GPU over CPU is more significant on dense operations, our setting of having sparse operations on CPU is appropriate. On the contrary, for $\tilde{P}_c, \tilde{Q}_c, \hat{P}_c, \hat{Q}_c, P_c$, and Q_c , we store these $k \times k$ matrices in the GPU memory. Then we compute $L^-(\theta)$ involving dense operations of Frobenius inner products on GPU.

To compute $\nabla L(\theta)$ in (27), similar to $L^+(\theta)$, XQ and $X^T P$ are computed on CPUs. Then by applying the mini-batch setting, for each subset, we feed the corresponding part of $XQ, X^T P, A, B, P, Q, \tilde{P}$ and \tilde{Q} into GPUs. With the cached $\hat{P}_c, \hat{Q}_c, P_c$, and Q_c , we first conduct dense operations to compute $XQ + \omega(APQ_c - A\tilde{P}\hat{Q}_c)$ and $X^T P + \omega(BQP_c - B\tilde{Q}\hat{P}_c)$ on GPU. Then we call the automatic differentiation function (e.g., `tf.gradients` in TensorFlow) to finish the transposed Jacobian-vector products.

To compute Gd in (40), as we mentioned, the Jacobian-vector products can be implemented with the forward-mode automatic differentiation. However, for compatibility with older versions of libraries where the forward-mode automatic differentiation is not supported, we apply the trick provided in [1] to compute W and H with the reverse-mode automatic differentiation. Similar to P

and Q , we apply the mini-batch setting and cache W and H in the CPU memory. Then with the cached W, H, P , and Q , we compute ZQ and $Z^T P$ on CPUs. Next, similar to $\nabla L(\theta)$, by the mini-batch setting, for each subset, we feed the required partial matrices to compute $AWQ_c + APH_c$ and $BHP_c + BQW_c$ in GPU. Finally, we call the automatic differentiation function to finish the remaining transposed Jacobian-vector products.

REFERENCES

- [1] Jamie Townsend. A new trick for calculating Jacobian vector products, 2017.
- [2] Chien-Chih Wang, Kent-Loong Tan, Chun-Ting Chen, Yu-Hsiang Lin, S. Sathya Keerthi, Dhruv Mahajan, Sellamanickam Sundararajan, and Chih-Jen Lin. Distributed Newton methods for deep learning. *Neural Comput.*, 30:1673–1724, 2018.