

Distributed Learning of Generalized Low Rank Models

— scientific program at École Polytechnique, April to June, 2018 —

Hoi-To Wai, Eric Moulines

March 3, 2018

Generalized low rank model (GLRM) is a powerful tool for describing real data such as health records, social networks, and gene expressions, etc. [1]. The model proposes easy-to-infer, yet general latent structure in the data, and therefore inferring these models allows us to *making sense out of the big-data*. In fact, the associated inference problems already cover popular and successful statistical tools [2] such as matrix completion, robust PCA, dictionary learning, etc.. For example, the low rank latent structure is efficient for clustering of data [3], discovering community structure [4], etc.

From a computational perspective, handling high-dimensional data (potentially low rank) has always been a prime issue as we seek ways to tame with the growing complexity. In particular, the saturation of Moore's Law in hardware design has made it increasingly difficult to scale up the computation power of *individual machines*. To our rescue, recent work have developed *distributed* and *parallelized* methods [5,6] that aim at distributing the computation load to *multiple machines*. This is relevant especially due to the now popular parallel computation architectures (such as GPUs) and internet-of-things hardwares, where machines with computation power are connected to each other.

The proposed scientific program is aimed at developing distributed learning strategies for the GLRMs inference. Specifically, in the model that we plan to study, we assume that the array of data, \mathbf{Y} , is observed through an exponential family distribution parameterized by an array of latent variable \mathbf{X} . Moreover, \mathbf{X} admits a decomposition such that $\mathbf{X} = \mathbf{\Theta} + \mathbf{S}$, where $\mathbf{\Theta}$ is a low rank matrix, and \mathbf{S} is a matrix that admits a low-dimensional representation such as sparseness, sparseness in a transformed domain, etc.. A particular feature in our model is that discrete data such as counts can be handled naturally. Note that these are the common data types in applications such as health records, social networks and gene expressions.

Concretely, we plan to study the distributed GLRM inference methods developed from a number of approaches including the alternating direction method of multiplier (ADMM) [7], the consensus optimization method [5] and its projection-free variant [8], etc. We also exploit the structure in ways that the data are distributed among machines or agents. As a side benefit, the method developed should also be *privacy preserving* such that the data shared among machines are anonymized, since a part of our goals is to analyze the sensitive health record data.

The visitor, Dr. Hoi-To Wai, has worked with the host, Prof. Eric Moulines, in summer 2015 and their collaboration has led to the aforementioned projection-free consensus optimization method [8].

REFERENCES

- [1] M. Udell *et al.* "Generalized low rank models." Foundations and Trends in Machine Learning, 2016.
- [2] E. Candès, X. Li, Y. Ma, and J. Wright. "Robust principal component analysis?." Journal of the ACM, 2011.
- [3] A. Ng, M. Jordan, and Y. Weiss. "On spectral clustering: Analysis and an algorithm." In NIPS 2002.
- [4] H.-T. Wai *et al.*, "Community Detection from Low Rank Excitations of a Graph Filter", to appear in ICASSP 2018.
- [5] A. Nedić, and A. Ozdaglar. "Distributed subgradient methods for multi-agent optimization." IEEE TAC, 2009.
- [6] B. Recht *et al.* "Hogwild: A lock-free approach to parallelizing stochastic gradient descent." In NIPS 2011.
- [7] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein. "Distributed optimization and statistical learning via the alternating direction method of multipliers." Foundations and Trends in Machine Learning, 2011.
- [8] H.-T. Wai, J. Lafond, A. Scaglione and E. Moulines, "Decentralized Frank-Wolfe Algorithm for Convex and Non-convex Optimization," IEEE TAC, 2017.