Trade-offs in Large-Scale Distributed Tuplewise Estimation and Learning

Robin Vogel^{1,2} Aurélien Bellet³ Stephan Clémençon¹ Ons Jelassi¹ Guillaume Papa¹

¹ Telecom Paris, ² IDEMIA, ³ Inria

ECML-PKDD 2019

Outline

Introduction

Overview of U-statistics

Estimators for distributed data

Learning with distributed U-statistics

Large scale distributed data processing

Very large datasets are common nowadays in ML. \rightarrow Distribute the data in **partitions** over several machines.

Cluster computing frameworks:



- \cdot Abstract network and communication aspects of distribution. \oplus
- \cdot Restrict the types of operations efficiently achieved. igodot
- E.g. Apache Spark [Meng et al., 2016], Petuum [Xing et al., 2015], ...

Most ML techniques optimize **standard means** $\hat{L} = \sum_i \ell(x_i)/n$, Those are **separable across partitions**.

 \rightarrow One can **efficiently estimate** those.

Very common statistics - e.g. U-statistics - are not. \rightarrow Estimation can be **slow or inaccurate**.

Our contribution

Problem:

When a statistic is **not separable across partitions**, frameworks may be **unsuited to computing accurate estimators.**

Contribution: Quantified analysis for U-statistics,

- 1. Efficient estimators of U-statistics in a distributed setting.
- 2. The analysis of their accuracy-vs-time tradeoff.
- 3. Learning experiment with those as gradient estimators.

Plan:

- Properties of U-statistics, see [Hoeffding, 1948],
- Contributions 1-3.

Outline

Introduction

Overview of U-statistics

Estimators for distributed data

Learning with distributed U-statistics

U-statistics: definition and examples

Introduce *K* independent i.i.d. samples $\mathcal{D}_k = \{X_1^{(k)}, \ldots, X_{n_k}^{(k)}\} \subset \mathcal{X}_k$, and the **kernel** $h : \mathcal{X}_1^{d_1} \times \cdots \times \mathcal{X}_K^{d_K} \to \mathbb{R}$ with $h \in L^2$.

The **generalized** *U*-statistic of degrees (d_1, \ldots, d_K) is defined as:

$$U_{\mathbf{n}}(h) = \frac{1}{\prod_{k=1}^{K} \binom{n_k}{d_k}} \sum_{l_1} \dots \sum_{l_K} h(\mathbf{X}_{l_1}^{(1)}, \mathbf{X}_{l_2}^{(2)}, \dots, \mathbf{X}_{l_K}^{(K)}),$$

where \sum_{l_k} is the sum over all $\binom{n_k}{d_k}$ subsets of d_k elements of \mathcal{D}_k .

Examples:

- a sample variance $h(x_1, x_2) = (x_1 x_2)^2$,
- Kendall's tau $h((x_1, y_1), (x_2, y_2)) = \mathbb{I}\{(x_1 x_2) \cdot (y_1 y_2) > 0\},\$
- clustering, metric learning and ranking criterions.

Here: study quantitatively a common case in a distributed setting.

Two-sample U-statistics

Introduce, with $m \ll n$:

- the **abundant positive** i.i.d. sample $\mathcal{D}_n = \{X_1, \ldots, X_n\} \subset \mathcal{X}$,
- the **scarce negative** i.i.d. sample $Q_m = \{Z_1, \ldots, Z_m\} \subset Z$,
- a kernel $h : \mathcal{X} \times \mathcal{Z} \mapsto \mathbb{R}$.

Objective: Estimate $U(h) = \mathbb{E}[h(X_1, Z_1)]$.

With $\mathbf{n} = (n, m)$, the **U-statistic** $U_{\mathbf{n}}$, where:

$$U_{\mathbf{n}} := \frac{1}{nm} \sum_{i=1}^{n} \sum_{j=1}^{m} h(X_i, Z_j),$$

is the **unbiased** estimator of U(h) with **lowest variance**.

 U_n is an average of nm elements \rightarrow What if $n = 10^9$, $m = 10^3$? One answer is **incomplete** *U*-statistics [Clémençon et al., 2016],

$$\widetilde{U}_B := \frac{1}{B} \sum_{(i,j)\in \mathcal{D}_B} h(X_i, Z_j),$$

where \mathcal{D}_B is a set of *B* elements sampled WR in $\{(i,j)\}_{i \in [\![n]\!], j \in [\![m]\!]}$.

Properties of two-sample *U*-statistics With the **Hajèk projections** $h_1(x) = \mathbb{E}[h(x, Z_1)], h_2(z) = \mathbb{E}[h(X_1, z)],$

and
$$h_0(x,z) = h(x,z) - h_1(x) - h_2(z) + U(h)$$
,

then $U_{\mathbf{n}} = T_1 + T_2 + W_0 - U$ with

$$T_1 = \frac{1}{n} \sum_{i=1}^n h_1(X_i) \text{ and } T_2 = \frac{1}{m} \sum_{j=1}^m h_2(Z_j) \text{ and } W_0 = \frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m h_0(X_i, Z_j),$$

which is called the **second Hoeffding decomposition**.

A U-statistic $U_{\mathbf{n}}$ is called **degenerate** when $h_1 = U$ and $h_2 = U$ a.s.. Introducing $\sigma_1^2 = \operatorname{Var}(h_1(X)), \sigma_2^2 = \operatorname{Var}(h_2(Z)), \sigma_0^2 = \operatorname{Var}(h_0(X_1, Z_1))$:

$$\operatorname{Var}(U_{\mathbf{n}}) = \frac{\sigma_1^2}{n} + \frac{\sigma_2^2}{m} + \frac{\sigma_0^2}{nm}$$

Examples with $X_1, Z_1 \sim \mathcal{U}[-1, 1]$: h(x, z) = x + z gives $\sigma_0^2 = 0$. | $h(x, z) = x \cdot z$ gives $\sigma_1^2, \sigma_2^2 = 0$.

Outline

Introduction

Overview of U-statistics

Estimators for distributed data

Learning with distributed U-statistics

Distributed data

We distribute the data on N workers.



One can distribute the instances:

- ▶ with SWOR (SWR): sampling without (with) replacement,
- ▶ **proportionally**: each cluster contains n/N instances from \mathcal{D}_n , and m/M from \mathcal{Q}_m .

Here: Focus on proportional SWOR (see paper for others).

Naive estimators (1/2)

Averaging full *U*-statistics from each cluster gives $U_{n,N}$.



Proposed statistics: U_{n,N}

Naive estimators (2/2)

Averaging *B* pairs SWR from each cluster gives $\widetilde{U}_{\mathbf{n},N,B}$.



Proposed statistics: $U_{\mathbf{n},N}$, $\tilde{U}_{\mathbf{n},N,B}$,

Estimators with redistribution (1/2)

Averaging $U_{n,N}$ on T redistributions of the data gives $\hat{U}_{n,N,T}$.



Proposed estimators: $U_{\mathbf{n},N}$, $\widetilde{U}_{\mathbf{n},N,B}$, $\widehat{U}_{\mathbf{n},N,T}$. With σ_t , π_t random permutations at time t.

Estimators with redistribution (2/2)



Proposed estimators: $U_{\mathbf{n},N}$, $\widetilde{U}_{\mathbf{n},N,B}$, $\widehat{U}_{\mathbf{n},N,T}$, and $\widetilde{U}_{\mathbf{n},N,B,T}$. With σ_t , π_t random permutations at time *t*.

All of the estimators are **unbiased**.

Variance expressions

We have: Variances in closed form.

For the **naive estimators**:

$$\operatorname{Var}(\boldsymbol{U}_{\mathbf{n},N}) = \operatorname{Var}(\boldsymbol{U}_{\mathbf{n}}) + (N-1)\frac{\sigma_0^2}{nm},$$
$$\operatorname{Var}(\widetilde{\boldsymbol{U}}_{\mathbf{n},N,B}) = \left(1 - \frac{1}{B}\right)\operatorname{Var}(\boldsymbol{U}_{\mathbf{n},N}) + \frac{\sigma^2}{NB}$$

 \rightarrow Term in $N\sigma_0^2/nm$.

For the estimators with redistribution:

$$\operatorname{Var}(\widehat{U}_{\mathbf{n},N,T}) = \operatorname{Var}(U_{\mathbf{n}}) + (N-1)\frac{\sigma_{0}^{2}}{nmT},$$
$$\operatorname{Var}(\widetilde{U}_{\mathbf{n},N,B,T}) = \operatorname{Var}(\widehat{U}_{\mathbf{n},N,T}) - \frac{1}{TB}\operatorname{Var}(U_{\mathbf{n},N}) + \frac{\sigma^{2}}{NTB}$$

 \rightarrow Term in $N\sigma_0^2/Tnm$.



Conclusions:

- · Redistribution is useful when $\sigma_2^2 \ll \sigma_0^2$,
- For any same # of pairs, $\hat{U}_{\mathbf{n},N,B,T}$ is worse than $\hat{U}_{\mathbf{n},N,T}$.

Outline

Introduction

Overview of *U***-statistics**

Estimators for distributed data

Learning with distributed U-statistics

Optimize the convexified AUC

Objective: optimize the AUC using a its convex upper-bound:

$$U_n(h_{AUC}) := \frac{1}{nm} \sum_{i,j} \mathbb{I}\{w^\top (X_i - Z_j) > 0\},$$

$$\leq \frac{1}{nm} \sum_{i,j} \left[1 + w^\top (X_i - Z_j)\right]_+ =: U_n(h_{conv}).$$

Optimizer: Momentum BGD with LR 10⁻², mom. 0.9 for 5×10^3 iter. **Loss:** $U_{\mathbf{n},N,B}(h_{\text{conv}}) + \lambda ||w||_2^2$, where we redistribute each n_r iter.

Parameters: $N = 100, B = 500, \lambda = 0.01, n_r \in \{1, 5, 25, +\infty\}.$

Dataset: Shuttle (outlier dataset), n = 45,000, m = 3,500, 20% as test set, train monitored on a fixed set of 450*K* pairs, see [Rayana, 2016].

Results



- \cdot Not redistributing ightarrow end performance very noisy \bigcirc
- \cdot Redistributing ightarrow better performance (with high. prob.) igodot

Danke, please come see us at the poster session !



References



Clémençon, S., Colin, I., and Bellet, A. (2016).

Scaling-up Empirical Risk Minimization: Optimization of Incomplete *U*-statistics. *Journal of Machine Learning Research*, 17(76):1–36.

Hoeffding, W. (1948).

A class of statistics with asymptotically normal distribution. *The Annals of Mathematical Statistics*, 19:293–325.



Rayana, S. (2016). Odds library.

Xing, E. P., Ho, Q., Dai, W., Kim, J. K., Wei, J., Lee, S., Zheng, X., Xie, P., Kumar, A., and Yu, Y. (2015). Petuum: A New Platform for Distributed Machine Learning on Big Data.

IEEE Transactions on Big Data, 1(2):49–67.