Using Efron-Stein's inequality for Optimal Transport Central Limit Theorem and Fairness issues in Machine Learning

Jean-Michel Loubes¹

¹ Institut de Mathématiques de Toulouse

Résumé. Nous utilisons l'inégalité de Efron-Stein pour prouver des résultats sur des théorèmes de limite centrale pour la distance de Wasserstein. Nous montrerons quelques applications pour détecter un usage déloyal d'un algorithme d'apprentissage. Ces résultats sont tirés de del Barrio et Loubes (2019) et del Barrio et al. (2018).

Mots-clés. Efron-Stein inequality, Wasserstein distance, CLT, Fair Learning.

1 Main Result

With inferential goals in mind, the main object of interest is the transportation cost between two sets of random points or between an empirical and a reference measure. In the, by now classical, Kantorovich formulation, for probabilities P and Q on R^d a transportation plan is a joint probability, say π , on $R^d \times R^d$ with marginals P and Q. The associated transportation cost is

$$I[\pi] = \int_{R^d \times R^d} c(x, y) d\pi(x, y),$$

where c is some cost function, and the optimal transportation cost is the minimal value of $I[\pi]$ among all choices of transportation plans, π , between P and Q. The problem admits a much more general formulation, but for our present purposes it is enough to know that for the choice $c(x, y) = c_p(x, y) = ||x - y||^p$, $p \ge 1$, if we denote by $\mathcal{W}_p^p(P, Q)$ the corresponding optimal transportation cost, then \mathcal{W}_p defines a metric in the set $\mathcal{F}_p(\mathbb{R}^d)$ of probabilities on \mathbb{R}^d with finite p-th moment. We refer to Villani (2009) for general background on these facts.

Assume that P and Q are Borel probabilities on \mathbb{R}^d with finite second moments and

Q has a positive density in the interior of its convex support. (1)

Consider P_n the empirical distribution.

We show first a variance bound for $\mathcal{W}_2^2(P_n, Q)$.

Theorem 1.1 If Q has a density and P and Q have finite fourth moments then

$$\operatorname{Var}(\mathcal{W}_{2}^{2}(P_{n},Q)) \leq \frac{C(P,Q)}{n},$$

where $C(P,Q) = 8 \Big(E(||X_1 - X_2||^2 ||X_1||^2) + (E||X_1 - X_2||^4)^{1/2} \Big(\int_{\mathbb{R}^d} ||y||^4 dQ(y) \Big)^{1/2} \Big).$

The main ingredient in the proof is the Efron-Stein inequality for variances, namely, that if $Z = f(X_1, \ldots, X_n)$ with X_1, \ldots, X_n independent random variables, (X'_1, \ldots, X'_n) is an independent copy of (X_1, \ldots, X_n) and $Z_i = f(X_1, \ldots, X'_i, \ldots, X_n)$ then

$$\operatorname{Var}(Z) \le \sum_{i=1}^{n} E(Z - Z_i)_{+}^{2}.$$

We refer, for instance, to Boucheron, Lugosi et Massart (2013) for a proof. In the particular case when X_1, \ldots, X_n are i.i.d. and f is a symmetric function of x_1, \ldots, x_n all the values $E(Z - Z_i)^2_+$ are equal and the bound simplifies to

$$\operatorname{Var}(Z) \le nE(Z - Z')_{+}^{2} \tag{2}$$

with $Z' = f(X'_1, X_2, ..., X_n)$.

Theorem 1.1 provides a simple bound with explicit constants for the variance of $W_2^2(P_n, Q)$ and implies tightness of $\sqrt{n}(W_2^2(P_n, Q) - E(W_2^2(P_n, Q)))$ with the only requirement of finite fourth moments and a density for Q. Next, we present a different application of the Efron-Stein inequality that will result in an approximation bound from which a CLT can be concluded.

Theorem 1.2 Assume that P and Q satisfy (1) and have finite moments of order $4 + \delta$ for some $\delta > 0$. Write φ_0 for the optimal transportation potential from P to Q. If

$$R_n = \mathcal{W}_2^2(P_n, Q) - \int_{R^d} (\|x\|^2 - 2\varphi_0(x)) dP_n(x)$$

then

$$n\operatorname{Var}(R_n) \to 0$$

as $n \to \infty$.

We will show how this result is an important step for the proof of a CLT for Monge-Kantorovich a.k.a Wasserstein distance and the consequences in fair learning as developed in del Barrio et al. (2018).

Bibliographie

Villani, C. (2009). Optimal Transport: Old and New. Grundlehren der Mathematischen Wissenschaften [Fundamental Principles of Mathematical Sciences] 338. Springer, Berlin.
Boucheron, S., Lugosi, G. and Massart, P. (2013). Concentration Inequalities. A Nonasymptotic Theory of Independence. Oxford.
del Barrio, E. and Loubes, J-M. (2019). Central limit theorems for empirical transportation cost in general dimension, The Annals of Probability 47 (2), 926-951
E Del Barrio, F Gamboa, P Gordaliza, JM Loubes (2018). Obtaining fairness using

optimal transport theory. arXiv preprint arXiv:1806.03195