Notes for Week 2/8/19 - 2/15/19

Alden Green

February 14, 2019

Consider $\mathbf{X} = (X_1, \dots, X_n) \sim \mathbb{P}$ and $\mathbf{Y} = (Y_1, \dots, Y_m) \sim \mathbb{Q}$ independently generated random variables. We will assume \mathbb{P} and \mathbb{Q} are supported on $\mathcal{D} \subset \mathbb{R}^d$ and are absolutely continuous with density functions f and g, respectively. Our statistical goal is hypothesis testing: that is, we wish to construct a test function ϕ which differentiates between

$$\mathbb{H}_0: f = g \text{ and } \mathbb{H}_1: f \neq g$$

For a given function class \mathcal{H} , some $\epsilon > 0$, and test function ϕ a Borel measurable function of the data with range $\{0,1\}$, we evaluate the quality of the test using worst-case risk

$$R_{\epsilon}^{m,n}(\phi;\mathcal{H}) = \sup_{f \in \mathcal{H}} \mathbb{E}_{f,f}^{(m,n)}(\phi) + \sup_{\substack{f,g \in \mathcal{H} \\ \delta(f,g) \ge \epsilon}} \mathbb{E}_{f,g}^{m,n}(1-\phi)$$

Total variation test. Let $\mathbf{Z} = (X_1, \dots, X_n, Y_1, \dots, Y_m)$, and let t = n + m. As in [1], define the K-NN graph $G_K = (V, E_K)$ to have vertex set $V = \{1, \dots, t\}$ and edge set E_K which contains the pair (i,j) if and only if x_i is among the K-nearest neighbors (with respect to Euclidean distance) of x_j , or vice versa. Let D_K denote the incidence matrix of G_K .

Define the kNN-total variation test statistic to be

$$T_{TV} = \sup_{\substack{\theta \in \mathbb{R}^t: \\ \|D_k \theta\|_1 \le C_{n,k}}} \left(\sum_{i=1}^n \theta_i - \sum_{j=n+1}^t \theta_j \right)$$
 (1)

Hereafter, take $\mathcal{D} = [0, 1]^d$, and consider

$$\mathcal{H}_{lip}(L) = \left\{ f : [0,1]^d \to \mathbb{R}^+ : \int_{\mathcal{D}} f = 1, f \text{ L-piecewise Lipschitz, bounded above and below} \right\}$$

where

Definition 0.1 (Piecewise Lipschitz). A function f is L-piecewise lipschitz over $[0,1]^d$ if there exists a set $\mathcal{S} \subset [0,1]^d$ such that

- (a) $\nu(S) = 0$
- (b) There exist $C_{\mathcal{S}}$, ϵ_0 such that $\mu((\mathcal{S}_{\epsilon} \cup (\partial \mathcal{D})_{\epsilon}) \cap [0,1]^d) \leq C_{\mathcal{S}}\epsilon$ for all $0 < \epsilon \leq \epsilon_0$.
- (c) For any z, z' in the same connected component of $[0, 1]^d \setminus (\mathcal{S}_{\epsilon} \cup (\partial \mathcal{D})_{\epsilon})$,

$$|g(z) - g(z')|_2 \le L ||z - z'||_2$$

and

Definition 0.2 (Bounded above and below). A function $f: \mathcal{D} \to \mathbb{R}$ is bounded above and below if there exists p_{\min}, p_{\max} such that

$$0 < p_{\min} < f(x) < p_{\max} < \infty \tag{} \forall x \in \mathcal{D})$$

Conjecture 1. For $\tau = ????$ and $K \approx \log^{1+2r}(n)$ for some $r \geq 0$, the test $\phi_{TV} = \{T_{TV} \geq \tau\}$ has worst-case risk

$$R_{\epsilon}^{(m,m)}(\mathcal{H}_{lip}(L)) \leq c_1/a^2$$

whenever $\epsilon \ge c_2 a \log^{\alpha} mm^{-1/d}$ where $\alpha = 3r + 5/2 + (2r + 1)/d$ and c_1 and c_2 are constants which depend only on (d, L).

Proof. Write

$$\sum_{i=1}^{m} \theta_i - \sum_{j=m+1}^{2m} \theta_j = \langle \theta, 1_X - 1_Y \rangle$$

where

$$1_X = \begin{cases} 1, & i = 1, \dots, m \\ 0, & \text{otherwise} \end{cases} \quad 1_Y = \begin{cases} 1, & i = j + 1, \dots, 2m \\ 0, & \text{otherwise} \end{cases}$$

Let

$$\widehat{\theta} \in \operatorname*{argmax}_{\theta \in \mathbb{R}^{2m}} \left\{ \langle \theta, 1_X - 1_Y \rangle : \|D_K \theta\|_1 \le C_{n,k} \right\}$$

satisfy $T_{TV} = \langle \widehat{\theta}, 1_X - 1_Y \rangle$. For $i = 1, \dots, 2m$, introduce θ^* defined by

$$(\theta^*)_i = \left(\frac{f(z_i) - g(z_i)}{f(z_i) + g(z_i)}\right)$$

We have

$$T_{TV} = \langle \widehat{\theta} - \theta^*, 1_X - 1_Y \rangle + \langle \theta^*, 1_X - 1_Y \rangle$$

We bound the first term, and begin by employing a basic inequality argument. Let

$$\widetilde{\theta} = \frac{\theta^{\star}}{\|D\theta^{\star}\|_{1}}.$$

Since $\widetilde{\theta}$ is feasible and $\widehat{\theta}$ optimal for (1), we have

$$\langle \widehat{\theta}, 1_X - 1_Y \rangle \ge \langle \widetilde{\theta}, 1_X - 1_Y \rangle$$

which, letting $w = (1_X - 1_Y) - \theta^*$, means

$$\langle \widehat{\theta}, \theta^* \rangle \ge \langle \widetilde{\theta}, \theta^* \rangle - \langle \widehat{\theta} - \widetilde{\theta}, w \rangle$$
 (2)

As a result, we have

$$\begin{split} \left| \langle \widehat{\theta} - \theta^{\star}, 1_{X} - 1_{Y} \rangle \right| &\leq \left| \langle \widehat{\theta} - \theta^{\star}, \theta^{\star} \rangle \right| + \left| \langle \widehat{\theta} - \theta^{\star}, w \rangle \right| \\ &\leq \left| \langle \widetilde{\theta} - \theta^{\star}, \theta^{\star} \rangle \right| + \left| \langle \widehat{\theta} - \widetilde{\theta}, w \rangle \right| + \left| \langle \widehat{\theta} - \theta^{\star}, w \rangle \right| \\ &\leq \left| \langle \widetilde{\theta} - \theta^{\star}, \theta^{\star} \rangle \right| + 2 \left| \langle \widehat{\theta} - \widetilde{\theta}, w \rangle \right| + \left| \langle \widetilde{\theta} - \theta^{\star}, w \rangle \right| \end{split}$$

where (i) follows from (2).

REFERENCES

[1] Oscar Hernan Madrid Padilla, James Sharpnack, Yanzhen Chen, and Daniela M Witten. Adaptive non-parametric regression with the k-nn fused lasso. $arXiv\ preprint\ arXiv:1807.11641$, 2018.