

Private Causal Inference (Appendix)

Matt J. Kusner

Washington University in St. Louis
mkusner@wustl.edu

Yu Sun

Cornell University
ys646@cornell.edu

Karthik Sridharan

Cornell University
sridharan@cs.cornell.edu

Kilian Q. Weinberger

Cornell University
kqw4@cornell.edu

In this appendix we first reproduce the Propose, Test, Release algorithm of [1] which we use to privatize the Interquartile Range dependence score (IQR), described in Section 4.3 of the main paper. We then prove various claims in the paper.

Algorithm 1 IQR Propose-Test-Release [1]

- 1: **Input:** data $\mathbf{X} = \{x_1, \dots, x_m\}$, privacy $\epsilon, \delta > 0$
 - 2: $k = \lfloor \log \text{IQR}(\mathbf{X}) \rfloor$
 - 3: $B_1 = [e^k, e^{k+1})$
 - 4: $B_2 = [e^{k-0.5}, e^{k+0.5})$
 - 5: **for** $j = 1, 2$ **do**
 - 6: $A_j :=$ number of data-points to modify to move IQR(\mathbf{X}) out of interval B_j
 - 7: $R_j = A_j + z$, where $z \sim \text{Lap}(0, \frac{1}{\epsilon})$
 - 8: **if** $R_j > 1 + \log(1/\delta)$ **then**
 - 9: **return** $\log \text{IQR}(\mathbf{X}) + z$, where $z \sim \text{Lap}(0, \frac{1}{\epsilon})$
 - 10: **end if**
 - 11: **end for**
 - 12: **return** \perp
-

Proof of Theorem 2. Let $x_1 \sim \text{Lap}(\mu_1, \sigma)$ and $x_2 \sim \text{Lap}(\mu_2, \sigma)$ be two independent Laplace random variables with $\mu_1 < \mu_2$, then the probability of failure is $\Pr(x_1 > x_2)$. We would like to compute the probability of failure in closed form. We know that by independence, the joint probability is equal to the product of marginal probabilities. We also know that the Laplace cdf. is

$$F(x; \mu, \sigma) = \begin{cases} F_1(x; \mu, \sigma) = \frac{1}{2} \exp(\frac{x-\mu}{\sigma}) & \text{if } x \leq \mu \\ F_2(x; \mu, \sigma) = 1 - \frac{1}{2} \exp(-\frac{x-\mu}{\sigma}) & \text{if } x > \mu \end{cases}$$

where F_1 and F_2 are only defined on the specified domains.

There are six mutually exclusive and collective exhaustive ways for which a failure could happen:

- ① $\mu_1 < x_2 < x_1 < \mu_2$
- ② $x_2 < \mu_1 < \mu_2 < x_1$
- ③ $\mu_1 < x_2 < \mu_2 < x_1$

$$\textcircled{4} x_2 < \mu_1 < x_1 < \mu_2$$

$$\textcircled{5} \mu_1 < \mu_2 < x_2 < x_1$$

$$\textcircled{6} x_2 < x_1 < \mu_1 < \mu_2$$

By symmetry of the Laplace distribution, we know that $\Pr(\textcircled{3}) = \Pr(\textcircled{4})$ and $\Pr(\textcircled{5}) = \Pr(\textcircled{6})$. Thus we only need to calculate $\Pr(\textcircled{1}), \Pr(\textcircled{2}), \Pr(\textcircled{3}),$ and $\Pr(\textcircled{5})$.

$$\begin{aligned} \Pr(\textcircled{1}) &= \int_{\mu_1}^{\mu_2} \int_{x_2}^{\mu_2} p(x_1)p(x_2)dx_1dx_2 \\ &= \int_{\mu_1}^{\mu_2} [F_2(\mu_2; \mu_1, \sigma) - F_2(x_2; \mu_1, \sigma)]p(x_2)dx_2 \end{aligned}$$

Now consider the quantity being integrated, which is equal to

$$-\frac{1}{2} \exp(-\frac{\mu_2 - \mu_1}{\sigma})p(x_2) + \underbrace{\frac{1}{2} \exp(-\frac{x_2 - \mu_1}{\sigma})p(x_2)}_*$$

The right-hand term is,

$$\begin{aligned} * &= \frac{1}{2} \exp(-\frac{x_2 - \mu_1}{\sigma}) \frac{1}{2\sigma} \exp(-\frac{\mu_2 - x_2}{\sigma}) \quad \text{since } x_2 < \mu_2 \\ &= \frac{1}{4\sigma} \exp(-\frac{\mu_2 - \mu_1}{\sigma}) \end{aligned}$$

So,

$$\begin{aligned} \Pr(\textcircled{1}) &= \frac{\mu_2 - \mu_1}{4\sigma} \exp(-\frac{\mu_2 - \mu_1}{\sigma}) \\ &\quad - \frac{1}{2} \exp(-\frac{\mu_2 - \mu_1}{\sigma}) \int_{\mu_1}^{\mu_2} p(x_2) \\ &= \frac{\mu_2 - \mu_1}{4\sigma} \exp(-\frac{\mu_2 - \mu_1}{\sigma}) \\ &\quad - \frac{1}{2} \exp(-\frac{\mu_2 - \mu_1}{\sigma}) [\frac{1}{2} - F_1(\mu_1; \mu_2, \sigma)] \end{aligned}$$

Next, we have that

$$\begin{aligned} \Pr(\textcircled{2}) &= \Pr(x_1 > \mu_2) \Pr(x_2 < \mu_1) \\ &= (1 - F_2(\mu_2; \mu_1, \sigma)) F_1(\mu_1; \mu_2, \sigma) \\ &= \frac{1}{2} \exp(-\frac{\mu_2 - \mu_1}{\sigma}) F_1(\mu_1; \mu_2, \sigma) \end{aligned}$$

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{z_1+z_2-z_3}^{\infty} p(z_1 | a, s)p(z_2 | b, s)p(z_3 | c, s)p(z_4 | d, s)dz_4dz_3dz_2dz_1 \quad (1)$$

And similarly,

$$\begin{aligned} \Pr(\textcircled{3}) &= \Pr(x_1 > \mu_2) \Pr(\mu_1 < x_2 < \mu_2) \\ &= (1 - F_2(\mu_2; \mu_1, \sigma)) \left[\frac{1}{2} - F_1(\mu_1; \mu_2, \sigma) \right] \\ &= \frac{1}{2} \exp\left(-\frac{\mu_2 - \mu_1}{\sigma}\right) \left[\frac{1}{2} - F_1(\mu_1; \mu_2, \sigma) \right] \end{aligned}$$

and as stated, $\Pr(\textcircled{4})$ is the same. Moving on,

$$\begin{aligned} \Pr(\textcircled{5}) &= \int_{\mu_2}^{\infty} \int_{x_2}^{\infty} p(x_1)p(x_2)dx_1dx_2 \\ &= \int_{\mu_2}^{\infty} [1 - F_2(x_2; \mu_1, \sigma)]p(x_2)dx_2 \\ &= \int_{\mu_2}^{\infty} \frac{1}{2} \exp\left(-\frac{x_2 - \mu_1}{\sigma}\right) \frac{1}{2\sigma} \exp\left(-\frac{x_2 - \mu_2}{\sigma}\right) dx_2 \\ &= \frac{1}{4} \int_{\mu_2}^{\infty} \frac{1}{2(\sigma/2)} \exp\left(-\frac{x_2 - (\mu_1 + \mu_2)/2}{\sigma/2}\right) dx_2 \\ &= \frac{1}{4} (1 - F_2(\mu_2; \mu', \sigma')) \\ &= \frac{1}{8} \exp\left(-\frac{\mu_2 - (\mu_1 + \mu_2)/2}{\sigma/2}\right) = \frac{1}{8} \exp\left(-\frac{\mu_2 - \mu_1}{\sigma}\right) \end{aligned}$$

and as stated, $\Pr(\textcircled{6})$ is the same. So lastly,

$$\begin{aligned} \Pr(x_1 > x_2) &= 2 \Pr(\textcircled{5}) + 2 \Pr(\textcircled{3}) + \Pr(\textcircled{2}) + \Pr(\textcircled{1}) \\ &= \frac{1}{4} \exp\left(-\frac{\mu_2 - \mu_1}{\sigma}\right) \\ &\quad + \frac{1}{2} \exp\left(-\frac{\mu_2 - \mu_1}{\sigma}\right) \\ &\quad - \exp\left(-\frac{\mu_2 - \mu_1}{\sigma}\right) F_1(\mu_1; \mu_2, \sigma) \\ &\quad + \frac{1}{2} \exp\left(-\frac{\mu_2 - \mu_1}{\sigma}\right) F_1(\mu_1; \mu_2, \sigma) \\ &\quad + \frac{\mu_2 - \mu_1}{4\sigma} \exp\left(-\frac{\mu_2 - \mu_1}{\sigma}\right) \\ &\quad - \frac{1}{4} \exp\left(-\frac{\mu_2 - \mu_1}{\sigma}\right) \\ &\quad + \frac{1}{2} \exp\left(-\frac{\mu_2 - \mu_1}{\sigma}\right) F_1(\mu_1; \mu_2, \sigma) \\ &= \frac{\mu_2 - \mu_1 + 2\sigma}{4\sigma} \exp\left(-\frac{\mu_2 - \mu_1}{\sigma}\right) \end{aligned}$$

This completes the derivation. \blacksquare

Proof of Theorem 4. Given a set of Laplace random variables: $z_1 \sim \text{Lap}(a, s)$, $z_2 \sim \text{Lap}(b, s)$, $z_3 \sim \text{Lap}(c, s)$, $z_4 \sim \text{Lap}(d, s)$ such that $a + b < c + d$, we would like to compute the probability that $z_1 + z_2 < z_3 + z_4$. Let $p(x | \mu, \sigma)$ be the pdf of the Laplace distribution $\text{Lap}(\mu, \sigma)$. Computing the above probability

requires evaluating the expression in eq. (1). Similar to the above proof, we can compute this integral by enumerating all of the possible cases, which gives the stated result. \blacksquare

Proof of Theorem 5.

$$\begin{aligned} |r'_{i,Y} - \tilde{r}'_{i,Y}| &= |\mathbf{w}^\top \phi(X) - \tilde{\mathbf{w}}^\top \phi(X)| \\ &\leq \|\mathbf{w} - \tilde{\mathbf{w}}\|_{\mathcal{H}} \|\phi(X)\|_{\mathcal{H}} \leq \|\mathbf{w} - \tilde{\mathbf{w}}\|_{\mathcal{H}} \quad (2) \end{aligned}$$

In the above we used the fact that $\|\phi(X)\|_{\mathcal{H}} = \sqrt{K(x, x)} \leq 1$. On the other hand note that \mathbf{w} is the minimizer of regularized objective on data set $(x_1, y_1), \dots, (x_n, y_n)$ and $\tilde{\mathbf{w}}$ is the minimizer on set $(x_1, y_1), \dots, (x_{n-1}, y_{n-1}), (x'_n, y'_n)$ (we assume the last coordinate is the one that is changes w.l.o.g.). By strong convexity of the regularized objective we have,

$$\begin{aligned} \frac{\lambda}{2} \|\mathbf{w} - \tilde{\mathbf{w}}\|^2 &\leq \frac{\lambda}{2} \|\tilde{\mathbf{w}}\|_{\mathcal{H}}^2 + \frac{1}{n} \sum_{i=1}^n (\tilde{\mathbf{w}}^\top \phi(x_i) - y_i)^2 \\ &\quad - \frac{\lambda}{2} \|\mathbf{w}\|_{\mathcal{H}}^2 - \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^\top \phi(x_i) - y_i)^2 \\ &\leq \frac{\lambda}{2} \|\tilde{\mathbf{w}}\|_{\mathcal{H}}^2 + \frac{1}{n} \sum_{i=1}^{n-1} (\tilde{\mathbf{w}}^\top \phi(x_i) - y_i)^2 + (\tilde{\mathbf{w}}^\top \phi(\tilde{x}_n) - \tilde{y}_n)^2 \\ &\quad - \frac{\lambda}{2} \|\mathbf{w}\|_{\mathcal{H}}^2 - \frac{1}{n} \sum_{i=1}^{n-1} (\mathbf{w}^\top \phi(x_i) - y_i)^2 - (\mathbf{w}^\top \phi(\tilde{x}_n) - \tilde{y}_n)^2 \\ &\quad + \frac{1}{n} (\tilde{\mathbf{w}}^\top \phi(x_n) - y_n)^2 - (\mathbf{w}^\top \phi(x_n) - y_n)^2 \\ &\quad - \frac{1}{n} (\tilde{\mathbf{w}}^\top \phi(\tilde{x}_n) - \tilde{y}_n)^2 + (\mathbf{w}^\top \phi(\tilde{x}_n) - \tilde{y}_n)^2 \\ &\leq \frac{2}{n} \sup_{x, y \in [-1, 1]} \left((\tilde{\mathbf{w}}^\top \phi(x) - y)^2 - (\mathbf{w}^\top \phi(x) - y)^2 \right) \\ &\leq \frac{2}{n} \|\tilde{\mathbf{w}} - \mathbf{w}\| \times (\|\tilde{\mathbf{w}}\| + \|\mathbf{w}\| + 2) \end{aligned}$$

Now note that since $0 \in \mathcal{H}$ we can conclude that,

$$\|\mathbf{w}\| \leq \frac{1}{\sqrt{\lambda}}$$

(The above is got by plugging in the 0 in the regularized objective which yields a value of 1 and since loss is non-negative, we can conclude that the norm of the minimizer of the regularized objective is atmost $1/\sqrt{\lambda}$. Plugging this in yields:

$$\|\mathbf{w} - \tilde{\mathbf{w}}\| \leq \frac{8}{\lambda^{3/2}n}$$

Plugging this in Eq. 2 yields the theorem. \blacksquare

Proof of tighter HSIC bound. Let D be the original dataset and D' be the dataset with one column modified. We subscript $HSIC$ with l and k implicitly. This is the quantity of interest

$$\begin{aligned} & |H\hat{S}IC(D) - H\hat{S}IC(D')| \\ &= \frac{1}{(N-1)^2} |tr(K'HL'H) - tr(KHLH)| \end{aligned}$$

Pulling out the constant, we have

$$\begin{aligned} & (N-1)^2 |H\hat{S}IC(D) - H\hat{S}IC(D')| \\ &= |tr(K'HL'H) - tr(KHLH)| \\ &= |tr((K'HL' - KHL)H)| \quad \text{linearity of trace} \\ &= |tr(H(K'HL' - KHL))| \quad \text{cyclicity of trace} \end{aligned}$$

Let $\mathbf{1}$ be the square matrix of *ones*(N). We know that since $H = I - \frac{1}{N}\mathbf{1}$ by definition,

$$\begin{aligned} H(K'HL' - KHL) &= (K'HL' - KHL) - \\ & \quad \frac{1}{N}\mathbf{1}(K'HL' - KHL) \end{aligned}$$

so we have that

$$\begin{aligned} & (N-1)^2 |H\hat{S}IC(D) - H\hat{S}IC(D')| \\ &= |tr(K'HL' - KHL) - \frac{1}{N}tr(\mathbf{1}(K'HL' - KHL))| \end{aligned} \quad (4)$$

Next, we need three identities. Let $sum(A) = \sum_{i,j} A$, then

$$\begin{aligned} \textbf{Identity 1:} & \quad tr(\mathbf{1}A) = sum(A) \\ \textbf{Identity 2:} & \quad tr(\mathbf{1}A\mathbf{1}B) = sum(A)sum(B) \\ \textbf{Identity 3:} & \quad sum(AB) = sum(BA) \end{aligned}$$

Where Identity 3 holds only for symmetric matrices. Identity 3 is obvious since $AB = (BA)^T$ and $sum(C) = sum(C^T)$, while the first two can be proven by expanding out the matrices and using the row-column rule, or just trying random matrices on MATLAB until you believe that it works. I did both, they are sure to be correct.

And again from the definition of H , we know that

$$KHL = K(L - \frac{1}{N}\mathbf{1}L) = KL - \frac{1}{N}K\mathbf{1}L$$

so

$$K'HL' - KHL = (K'L' - \frac{1}{N}K'\mathbf{1}L') - (KL - \frac{1}{N}K\mathbf{1}L)$$

Now we continue our derivation of eq. (4)

$$\begin{aligned} & (N-1)^2 |H\hat{S}IC(D) - H\hat{S}IC(D')| = \\ & \underbrace{|tr(K'HL' - KHL)|}_{\star} - \underbrace{\frac{1}{N}sum(K'HL' - KHL)}_{\diamond} \end{aligned}$$

We can rewrite each term \star and \diamond using our traces identities as follows,

$$\begin{aligned} \star &= [tr(K'L') - \frac{1}{N}sum(L'K')] - [tr(KL) - \frac{1}{N}sum(LK)] \\ \diamond &= \frac{1}{N}[sum(K'L' - \frac{1}{N}K'\mathbf{1}L') - sum(KL - \frac{1}{N}K\mathbf{1}L)] \\ &= \frac{1}{N}[sum(K'L') - \frac{1}{N}sum(K'\mathbf{1}L')] \\ & \quad - \frac{1}{N}[sum(KL) - \frac{1}{N}sum(K\mathbf{1}L)] \end{aligned}$$

By identity 3, we see that the $sum(KL)$ and $sum(LK)$ as well as the $sum(K'L')$ and $sum(L'K')$ terms in \star and \diamond are identical. Thus we are left with

$$\begin{aligned} & (N-1)^2 |H\hat{S}IC(D) - H\hat{S}IC(D')| = \star - \diamond \\ &= |[tr(K'L') - tr(KL)] \\ & \quad - \frac{2}{N}[sum(K'L') - sum(KL)] \\ & \quad + \frac{1}{N}[\frac{1}{N}sum(K'\mathbf{1}L') - \frac{1}{N}sum(K\mathbf{1}L)]| \\ &= |[sum(K'.*L') - sum(K.*L)] \\ & \quad - \frac{2}{N}[sum(K'L') - sum(KL)] \\ & \quad + \frac{1}{N^2}[sum(K')sum(L') - sum(K)sum(L)]| \quad (5) \end{aligned}$$

where the last line comes from applying Identity 1 backwards so we have, for example, $tr(\mathbf{1}K\mathbf{1}L)$, then applying Identity 2. We use MATLAB[®] notation $.*$ for the element-wise product of two matrices.

We bound eq. (5) by the triangle inequality,

$$\begin{aligned} & (N-1)^2 |H\hat{S}IC(D) - H\hat{S}IC(D')| \\ & \leq |sum(K'.*L') - sum(K.*L)| \quad \textcircled{1} \\ & \quad + \frac{2}{N}|sum(K'L') - sum(KL)| \quad \textcircled{2} \\ & \quad + \frac{1}{N^2}|sum(K')sum(L') - sum(K)sum(L)| \quad \textcircled{3} \end{aligned}$$

$$\begin{aligned} & (N-1)^2 |H\hat{S}IC(D) - H\hat{S}IC(D')| \\ & \leq [\textcircled{1} + \textcircled{2} + \textcircled{3}] \leq \max_{K,K',L,L'} \textcircled{1} + \max_{K,K',L,L'} \textcircled{2} + \max_{K,K',L,L'} \textcircled{3} \end{aligned}$$

Recall that the kernels k and l are bounded by 1. And that the kernel pairs K, K' and L, L' differ in at most one row and column. Thus, for $\textcircled{1}$, it is clear that the maximum occurs when a row and column c (no matter what c is) is changed from all 0 to 1 in both L and K , so $\max_{K,K',L,L'} \textcircled{1} = 2N - 1$

For $\textcircled{3}$, the maximum occurs at exactly same the point as $\textcircled{1}$, and the value achieved is

$$\frac{1}{N^2}[N^4 - (N^2 - 2N + 1)^2] \leq 4N - 5$$

for $N > 3$. For ②, applying the row-column rule and reasoning on small matrices inductively suggest that the maximum is also achieved when we change one row and column of L and K from all 0 to 1, and is thus

$$\frac{2}{N}[N^2 + (N - 1)(2N - 1)] \leq 6N - 5$$

for $N \geq 2$. As the argmax of all three terms coincide, we have that

$$\max_{\mathcal{C}} [\textcircled{1} + \textcircled{2} + \textcircled{3}] = \max_{\mathcal{C}} \textcircled{1} + \max_{\mathcal{C}} \textcircled{2} + \max_{\mathcal{C}} \textcircled{3}$$

Therefore, we have derived that for all practical purposes, the overall bound is $\frac{12N-11}{N^2-1}$. ■

References

- [1] Dwork, Cynthia and Lei, Jing. Differential privacy and robust statistics. In *Proceedings of the forty-first annual ACM symposium on Theory of computing*, pp. 371–380. ACM, 2009.