

Supplementary Material to Paper "Local Differential Privacy for Sampling"

Abstract

This is the Supplementary Material to Paper "Local Differential Privacy for Sampling".
Notation "main file" indicates reference to the submitted draft.

Appendix: table of contents

| | |
|----------------------------------|-------|
| Proofs and formal results | Pg 2 |
| Proof of Lemma 2 | Pg 2 |
| Proof of Theorem 4 | Pg 2 |
| Proof of Theorem 5 | Pg 3 |
| Proof of Theorem 6 | Pg 6 |
| Proof of Theorem 7 | Pg 7 |
| Additional formal results | Pg 9 |
| Additional experiments | Pg 12 |

1 Proofs and formal results

1.1 Proof of Lemma 2

For any $x \in \mathcal{X}$ and $P, P' \in \mathcal{D}(\mathcal{X})$, we have $\Pr[A(P) = x] \in \mathcal{M}$ and $\Pr[A(P') = x] \in \mathcal{M}$ by the fact that $A(P)$ samples from densities that lie in the mollifier \mathcal{M} . By definition of ε -mollifiers, the density ratio between any two densities in the ε -mollifiers is bounded by $\exp(\varepsilon)$, meaning we have

$$\frac{\Pr[A(P) = x]}{\Pr[A(P') = x]} \leq \exp(\varepsilon), \quad (1)$$

and thus A is an ε -private sampler.

1.2 Proof of Theorem 4

The proof follows from two Lemma which we state and prove.

Lemma 1 *For any $T \in \mathbb{N}_*$, we have that*

$$\sum_{t=1}^T \theta_t(\varepsilon) = \sum_{t=1}^T \left(\frac{\varepsilon}{\varepsilon + 4 \log(2)} \right)^t < \frac{\varepsilon}{4 \log(2)}. \quad (2)$$

Proof Since $(\varepsilon/(\varepsilon + 4 \log(2))) < 1$ for any ε and noting that $\theta_t(\varepsilon) = (\varepsilon/(\varepsilon + 4 \log(2)))\theta_{t-1}(\varepsilon)$, we can conclude that $\theta_t(\varepsilon)$ is a geometric sequence. For any geometric series with ratio r , we have that

$$\sum_{t=1}^T r^t = r \left(\frac{1 - r^T}{1 - r} \right) \quad (3)$$

$$= \frac{r}{1 - r} - \frac{r^{T+1}}{1 - r} \quad (4)$$

$$< \frac{r}{1 - r} \quad (5)$$

Indeed, $\frac{r}{1-r}$ is the limit of the geometric series above when $T \rightarrow \infty$. In our case, we let $r = (\varepsilon/(\varepsilon + 4 \log(2)))$ to show that

$$\frac{r}{1 - r} = \frac{\frac{\varepsilon}{\varepsilon + 4 \log(2)}}{1 - \frac{\varepsilon}{\varepsilon + 4 \log(2)}} = \frac{\frac{\varepsilon}{\varepsilon + 4 \log(2)}}{\frac{4 \log(2)}{\varepsilon + 4 \log(2)}} = \frac{\varepsilon}{4 \log(2)}, \quad (6)$$

which concludes the proof. ■

Lemma 2 *For any $\varepsilon > 0$ and $T \in \mathbb{N}_*$, let $\theta(\varepsilon) = (\theta_1(\varepsilon), \dots, \theta_T(\varepsilon))$ denote the parameters and $c = (c_1, \dots, c_t)$ denote the sufficient statistics returned by Algorithm 1, then we have*

$$-\frac{\varepsilon}{2} \leq \langle \theta(\varepsilon), c \rangle - \varphi(\theta(\varepsilon)) \leq \frac{\varepsilon}{2}. \quad (7)$$

Proof Since the algorithm returns classifiers such that $c_t(x) \in [-\log 2, \log 2]$ for all $1 \leq t \leq T$, we have from Lemma 1,

$$\sum_{t=1}^T \theta_t(\varepsilon) c_t \leq \log(2) \sum_{t=1}^T \theta_t(\varepsilon) < \log(2) \frac{\varepsilon}{4 \log(2)} = \frac{\varepsilon}{4}, \quad (8)$$

and similarly,

$$\sum_{t=1}^T \theta_t(\varepsilon) c_t \geq -\log(2) \sum_{t=1}^T \theta_t(\varepsilon) > -\log(2) \frac{\varepsilon}{4 \log(2)} = -\frac{\varepsilon}{4}. \quad (9)$$

Thus we have

$$-\frac{\varepsilon}{4} \leq \langle \theta(\varepsilon), c \rangle \leq \frac{\varepsilon}{4}. \quad (10)$$

By taking exponential, integrand (w.r.t Q_0) and logarithm of 10, we get

$$\log \int_x \exp\left(-\frac{\varepsilon}{4}\right) dQ_0 \leq \log \int_x \exp(\langle \theta(\varepsilon), c \rangle) dQ_0 \leq \log \int_x \exp\left(\frac{\varepsilon}{4}\right) dQ_0 \quad (11)$$

$$-\frac{\varepsilon}{4} \leq \varphi(\theta(\varepsilon)) \leq \frac{\varepsilon}{4} \quad (12)$$

Since $\langle \theta(\varepsilon), c \rangle \in [-\varepsilon/4, \varepsilon/4]$ and $\varphi(\theta(\varepsilon)) \in [-\varepsilon/4, \varepsilon/4]$, the proof concludes by considering highest and lowest values. ■

The proof of Theorem 4 now follows from taking the exp of all quantities in (7), which makes appear Q_T in the middle and conditions for membership to \mathcal{M}_ε in the bounds.

1.3 Proof of Theorem 5

We begin by first deriving the KL drop expression. At each iteration, we learn a classifier c_t , fix some step size $\theta > 0$ and multiply Q_{t-1} by $\exp(\theta \cdot c_t)$ and renormalize to get a new distribution which we will denote by $Q_t(\theta)$ to make the dependence of θ explicit.

Lemma 3 For any $\theta > 0$, let $\varphi(\theta) = \log \int_x \exp(\theta \cdot c_t) dQ_{t-1}$. The drop in KL is

$$DROP(\theta) := KL(P, Q_{t-1}) - KL(P, Q_t(\theta)) = \theta \cdot \int_x c_t dP - \varphi(\theta) \quad (13)$$

Proof Note that $Q_t(\theta)$ is indeed a one dimensional exponential family with natural parameter θ , sufficient statistic c_t , log-partition function $\varphi(\theta)$ and base measure Q_{t-1} . We can write out the KL

divergence as

$$\text{KL}(P, Q_{t-1}) - \text{KL}(P, Q_t(\theta)) = \int_x \log \left(\frac{P}{Q_{t-1}} \right) dP - \int_x \log \left(\frac{P}{\exp(\theta \cdot c_t - \varphi(\theta)) Q_{t-1}} \right) dP \quad (14)$$

$$= \int_x \log \left(\frac{\exp(\theta \cdot c_t - \varphi(\theta)) Q_{t-1}}{Q_{t-1}} \right) dP \quad (15)$$

$$= \int_x \theta \cdot c_t - \varphi(\theta) dP \quad (16)$$

$$= \theta \cdot \int_x c_t dP - \varphi(\theta) \quad (17)$$

■

It is not hard to see that the drop is indeed a concave function of θ , suggesting that there exists an optimal step size at each iteration. We split our analysis by considering two cases and begin when $\gamma_Q^t < 1/3$. Since $\theta > 0$, we can lowerbound the first term of the KL drop using WLA. The trickier part however, is bounding $\varphi(\theta)$ which we make use of Hoeffding's lemma.

Lemma 4 (Hoeffding's Lemma) *Let X be a random variable with distribution Q , with $a \leq X \leq b$ such that $\mathbb{E}_Q[X] = 0$, then for all $\lambda > 0$, we have*

$$\mathbb{E}_Q[\exp(\lambda \cdot X)] \leq \exp \left(\frac{\lambda^2(b-a)^2}{8} \right) \quad (18)$$

Lemma 5 *For any classifier c_t satisfying Assumption 3 (WLA), we have*

$$\mathbb{E}_{Q_{t-1}}[\exp(\theta_t(\varepsilon) \cdot c_t)] \leq \exp \left(\theta_t^2(\varepsilon) \cdot \frac{(c_t^*)^2}{2} - \theta_t(\varepsilon) \cdot \gamma_Q^t \cdot c_t^* \right) \quad (19)$$

Proof Let $X = c_t - \mathbb{E}_{Q_{t-1}}[c_t]$, $b = c_t^*$, $a = -c_t^*$ and $\lambda = \theta_t(\varepsilon)$ and noticing that

$$\mathbb{E}_{Q_{t-1}}[\lambda \cdot X] = \mathbb{E}_{Q_{t-1}}[c_t - \mathbb{E}_{Q_{t-1}}[c_t]] = \mathbb{E}_{Q_{t-1}}[c_t] - \mathbb{E}_{Q_{t-1}}[c_t] = 0, \quad (20)$$

allows us to apply Lemma 4. By first realizing that

$$\exp(\lambda \cdot X) = \exp(\theta_t(\varepsilon) \cdot c_t) \cdot \exp(\theta_t(\varepsilon) \cdot \mathbb{E}_{Q_{t-1}}[-c_t]), \quad (21)$$

We get that

$$\mathbb{E}_{Q_{t-1}}[\exp(\theta_t(\varepsilon) \cdot c_t)] \cdot \exp(\theta_t(\varepsilon) \cdot \mathbb{E}_{Q_{t-1}}[-c_t]) \leq \exp \left(\theta_t^2(\varepsilon) \cdot \frac{(c_t^*)^2}{2} \right). \quad (22)$$

Re-arranging and using the WLA inequality yields

$$\mathbb{E}_{Q_{t-1}}[\exp(\theta_t(\varepsilon) \cdot c_t)] \leq \exp \left(\theta_t^2(\varepsilon) \cdot \frac{(c_t^*)^2}{2} - \theta_t(\varepsilon) \cdot \mathbb{E}_{Q_{t-1}}[-c_t] \right) \quad (23)$$

$$\leq \exp \left(\theta_t^2(\varepsilon) \cdot \frac{(c_t^*)^2}{2} - \theta_t(\varepsilon) \cdot \gamma_Q^t \cdot c_t^* \right) \quad (24)$$

Applying Lemma 5 and Lemma 3 (writing $Q_t = Q_t(\varepsilon)$) together gives us

$$\text{KL}(P, Q_t) = \text{KL}(P, Q_{t-1}) - \text{DROP}(\theta_t(\varepsilon)) \quad (25)$$

$$= \text{KL}(P, Q_{t-1}) - \theta_t(\varepsilon) \cdot \int_{\mathcal{X}} c_t dP + \log \mathbb{E}_{Q_{t-1}}[\exp(\theta_t(\varepsilon) \cdot c_t)] \quad (26)$$

$$\leq \text{KL}(P, Q_{t-1}) - c_t^* \cdot \theta_t(\varepsilon) \cdot \left(\frac{1}{c_t^*} \int_{\mathcal{X}} c_t dP \right) + \left(\theta_t^2(\varepsilon) \cdot \frac{(c_t^*)^2}{2} - \theta_t(\varepsilon) \cdot \gamma_Q^t \cdot c_t^* \right) \quad (27)$$

$$\leq \text{KL}(P, Q_{t-1}) - c_t^* \theta_t(\varepsilon) \left(\gamma_P^t + \gamma_Q^t - \frac{c_t^* \cdot \theta_t(\varepsilon)}{2} \right) \quad (28)$$

Now we move to the case of $\gamma_Q^t \geq 1/3$.

Lemma 6 *For any classifier c_t returned by Algorithm 1, we have that*

$$\mathbb{E}_{Q_{t-1}}[\exp(c_t)] \leq \exp(-\Gamma(\gamma_Q^t)) \quad (29)$$

where $\Gamma(z) = \log(4/(5-3z))$.

Proof Consider the straight line between $(-\log 2, 1/2)$ and $(\log 2, 2)$ given by $y = 5/4 + (3/(4 \cdot \log 2))x$, which by convexity is greater than $y = \exp(x)$ on the interval $[-\log 2, \log 2]$. To this end, we define the function

$$f(x) = \begin{cases} \frac{5}{4} + \frac{3}{4 \cdot \log 2} \cdot x, & \text{if } x \in [-\log 2, \log 2] \\ 0, & \text{otherwise} \end{cases} \quad (30)$$

Since $c_t(x) \in [-\log 2, \log 2]$ for all $x \in \mathcal{X}$, we have that $f(c_t(x)) \geq \exp(c_t(x))$ for all $x \in \mathcal{X}$. Taking $\mathbb{E}_{Q_{t-1}}[\cdot]$ over both sides and using linearity of expectation gives

$$\mathbb{E}_{Q_{t-1}}[\exp(c_t(x))] \leq \mathbb{E}_{Q_{t-1}}[f(c_t(x))] \quad (31)$$

$$= \frac{5}{4} + \frac{3}{4 \log 2} (\mathbb{E}_{Q_{t-1}}[c_t(x)]) \quad (32)$$

$$= \frac{5}{4} - \frac{3}{4} \left(\frac{1}{\log 2} \mathbb{E}_{Q_{t-1}}[-c_t(x)] \right) \quad (33)$$

$$< \frac{5}{4} - \frac{3}{4} \gamma_Q^t \quad (34)$$

$$= \exp \left(-\log \left(\frac{5 - 3\gamma_Q^t}{4} \right)^{-1} \right) \quad (35)$$

$$= \exp \left(-\log \left(\frac{4}{5 - 3\gamma_Q^t} \right) \right) \quad (36)$$

$$= \exp(-\Gamma(\gamma_Q^t)), \quad (37)$$

as claimed. ■

Now we use Lemma 3 and Jensen's inequality since $\theta_t(\varepsilon) < 1$ so that

$$\text{KL}(P, Q_t) = \text{KL}(P, Q_{t-1}) - \text{DROP}(\theta) \quad (38)$$

$$= \text{KL}(P, Q_{t-1}) - \theta_t(\varepsilon) \cdot \int_x c_t dP + \log \mathbb{E}_{Q_{t-1}}[\exp(\theta_t \cdot c_t)] \quad (39)$$

$$\leq \text{KL}(P, Q_{t-1}) - \theta_t(\varepsilon) \cdot \mathbb{E}_P[c_t] + \theta_t \cdot \log \mathbb{E}_{Q_{t-1}}[\exp(c_t)] \quad (40)$$

$$\leq \text{KL}(P, Q_{t-1}) - \theta_t(\varepsilon) (\mathbb{E}_P[c_t] - \log \mathbb{E}_{Q_{t-1}}[\exp(c_t)]) \quad (41)$$

$$= \text{KL}(P, Q_{t-1}) - \theta_t(\varepsilon) \left(c_t^* \left(\frac{1}{c_t^*} \mathbb{E}_P[c_t] \right) - \log \mathbb{E}_{Q_{t-1}}[\exp(c_t)] \right) \quad (42)$$

$$< \text{KL}(P, Q_{t-1}) - \theta_t(\varepsilon) (c_t^* \gamma_P^t - \log(\exp(-\Gamma(\gamma_Q^t)))) \quad (43)$$

$$= \text{KL}(P, Q_{t-1}) - \theta_t(\varepsilon) (c_t^* \gamma_P^t + \Gamma(\gamma_Q^t)). \quad (44)$$

1.4 Proof of Theorem 6

We first note that for any $Q \in \mathcal{M}_\varepsilon$,

$$\text{KL}(P, Q) = \int_x \log \left(\frac{P}{Q} \right) dP \quad (45)$$

$$= \int_x \log \left(\frac{P}{Q_0} \frac{Q_0}{Q} \right) dP \quad (46)$$

$$= \int_x \log \left(\frac{P}{Q_0} \right) dP - \int_x \log \left(\frac{Q_0}{Q} \right) dP \quad (47)$$

$$\geq \text{KL}(P, Q_0) - \int_x \frac{\varepsilon}{2} dP \quad (48)$$

$$\geq \text{KL}(P, Q_0) - \frac{\varepsilon}{2}, \quad (49)$$

which completes the proof of the upperbound To show (13), we have that

$$\text{KL}(P, Q_t) \leq \text{KL}(P, Q_{T-1}) - \theta_t(\varepsilon) \cdot \Lambda_t \quad (50)$$

$$\leq \text{KL}(P, Q_0) - \sum_{t=1}^{T-1} \theta_t(\varepsilon) \cdot \Lambda_t \quad (51)$$

$$= \text{KL}(P, Q_0) - \sum_{t=1}^{T-1} \theta_t(\varepsilon) \cdot (c_t^* \gamma_P^t + \Gamma(\gamma_Q^t)) \quad (52)$$

$$\leq \text{KL}(P, Q_0) - \sum_{t=1}^{T-1} \theta_t(\varepsilon) \cdot (\log 2 \cdot \gamma_P + \Gamma(\gamma_Q)) \quad (53)$$

$$\leq \text{KL}(P, Q_0) - (\log 2 \cdot \gamma_P + \log 2 \cdot \gamma_Q) \cdot \sum_{t=1}^{T-1} \theta_t(\varepsilon) \quad (54)$$

$$\leq \text{KL}(P, Q_0) - (\log 2 \cdot \gamma_P + \log 2 \cdot \gamma_Q) \cdot \sum_{t=1}^{T-1} \theta_t(\varepsilon) \quad (55)$$

$$= \text{KL}(P, Q_0) - \log 2 \cdot (\gamma_P + \gamma_Q) \cdot \theta_1(\varepsilon) \cdot \left(\frac{1 - \theta_t(\varepsilon)}{1 - \theta_1(\varepsilon)} \right) \quad (56)$$

$$= \text{KL}(P, Q_0) - \varepsilon \cdot \left(\frac{\gamma_P + \gamma_Q}{4} \right) \cdot (1 - \theta_t(\varepsilon)), \quad (57)$$

where we used the fact that $\Gamma(x) \geq \log 2 \cdot x$ and explicit geometric summation expression.

1.5 Proof of Theorems 7

We start by a general Lemma.

Lemma 7 *For any region of the support B , we have that*

$$\int_B dQ_t \geq \int_B dP - \int_B \log \left(\frac{P}{Q_t} \right) dP \quad (58)$$

Proof By first noting that for any region B ,

$$\int_B (dP - dQ_t) = \int_B \left(1 - \frac{dQ_t}{dP} \right) dP \quad (59)$$

we then use the inequality $1 - x \leq \log(1/x)$ to get

$$\int_B (dP - dQ_t) = \int_B \left(1 - \frac{dP}{dQ_t} \right) dP \leq \int_B \log \left(\frac{dP}{dQ_t} \right) dP = \int_B \log \left(\frac{P}{Q_t} \right) dP \quad (60)$$

Re-arranging the above inequality gives us the bound. ■

Lemma 7 allows us to understand the relationship between two distributions P and Q_t in terms regions they capture. The general goal is to show that for a given region B (which includes the

highly dense mode regions), the amount of mass captured by the model $\int_B dQ_t$, is lower bounded by the target mass $\int_B dP$, and some small quantity. The inequality in Lemma 7 comments on this precisely with the small difference being a term that looks familiar to the KL-divergence - rather one that is bound to the specific region B . Though, this term can be understood to be small since by Theorem 5, we know that the global KL decreases, we give further refinements to show the importance of privacy parameters ε . We show that the term $\int_B \log(P/Q_t)dP$ can be decomposed in different ways, leading to our two Theorems to prove.

Lemma 8

$$\int_B \log\left(\frac{P}{Q_t}\right) dP \leq \int_B \log\left(\frac{P}{Q_0}\right) dP - \Delta + \frac{\varepsilon}{2} \left(1 - \int_B dP\right). \quad (61)$$

where $\Delta = KL(P, Q_0) - KL(P, Q_t)$

Proof We decompose the space \mathcal{X} into B and the complement B^c to get

$$\int_B \log\left(\frac{P}{Q_t}\right) dP = \int_{\mathcal{X}} \log\left(\frac{P}{Q_t}\right) dP - \int_{B^c} \log\left(\frac{P}{Q_t}\right) dP \quad (62)$$

$$= \mathbf{KL}(P, Q_t) - \int_{B^c} \log\left(\frac{P}{Q_t}\right) dP \quad (63)$$

$$\leq \mathbf{KL}(P, Q_0) - \Delta - \int_{B^c} \log\left(\frac{P}{Q_t}\right) dP, \quad (64)$$

where we used Theorem 5, and letting $\theta = \theta(\varepsilon)$ for brevity, we also have

$$\int_{B^c} \log\left(\frac{P}{Q_t}\right) dP = \int_{B^c} \log\left(\frac{P}{Q_0 \exp(\langle \theta, c \rangle - \varphi(\theta))}\right) dP \quad (65)$$

$$= \int_{B^c} \log\left(\frac{P}{Q_0}\right) dP - \int_{B^c} \exp(\langle \theta, c \rangle - \varphi(\theta)) dP \quad (66)$$

$$\geq \int_{B^c} \log\left(\frac{P}{Q_0}\right) dP - \int_{B^c} \frac{\varepsilon}{2} dP \quad (67)$$

$$= \int_{B^c} \log\left(\frac{P}{Q_0}\right) dP - \frac{\varepsilon}{2} \left(1 - \int_B dP\right) \quad (68)$$

Combining these inequalities together gives us:

$$\int_B \log\left(\frac{P}{Q_t}\right) dP \leq \mathbf{KL}(P, Q_0) - \Delta - \left(\int_{B^c} \log\left(\frac{P}{Q_0}\right) dP - \frac{\varepsilon}{2} \left(1 - \int_B dP\right)\right) \quad (69)$$

$$= \int_{\mathcal{X}} \log\left(\frac{P}{Q_0}\right) dP - \int_{B^c} \log\left(\frac{P}{Q_0}\right) dP - \Delta + \frac{\varepsilon}{2} \left(1 - \int_B dP\right) \quad (70)$$

$$= \int_B \log\left(\frac{P}{Q_0}\right) dP - \Delta + \frac{\varepsilon}{2} \left(1 - \int_B dP\right) \quad (71)$$

■

We are now in a position to prove Theorem 7. Using Lemma 8 into the inequality in Lemma 7 yields

$$\int_B dQ_t \geq \int_B dP - \left(\int_B \log \left(\frac{P}{Q_0} \right) dP - \Delta + \frac{\varepsilon}{2} \left(1 - \int_B dP \right) \right) \quad (72)$$

$$= \left(1 + \frac{\varepsilon}{2} \right) \int_B dP - \frac{\varepsilon}{2} - \int_B \log \left(\frac{P}{Q_0} \right) + \Delta. \quad (73)$$

Reorganising and using the Theorem's notations, we get

$$\mathfrak{M}(B, Q) \geq \mathfrak{M}(B, P) - KL(P, Q_0; B) + \frac{\varepsilon}{2} \cdot J(P, Q; B, \varepsilon), \quad (74)$$

where we recall that $J(P, Q; B, \varepsilon) \doteq \mathfrak{M}(B, P) + \frac{2\Delta(Q)}{\varepsilon} - 1$. Theorem 6 says that we have in the high boosting regime $2\Delta(Q_T)/\varepsilon \geq (\gamma_P + \gamma_Q)/2 - \theta_T(\varepsilon) \cdot (\gamma_P + \gamma_Q)/2$. Letting $\bar{\gamma} \doteq (\gamma_P + \gamma_Q)/2$ and $K \doteq 4 \log 2$, we have from MBDE in the high boosting regime:

$$\begin{aligned} \frac{2\Delta(Q)}{\varepsilon} &\geq \bar{\gamma} \cdot \left(1 - \left(\frac{1}{1 + \frac{K}{\varepsilon}} \right)^T \right) \\ &\geq \bar{\gamma} \cdot \left(1 - \frac{1}{1 + \frac{TK}{\varepsilon}} \right) \\ &= \bar{\gamma} \cdot \frac{TK}{TK + \varepsilon}. \end{aligned} \quad (75)$$

To have $J(P, Q; B, \varepsilon) \geq -(2/\varepsilon) \cdot \alpha \mathfrak{M}(B, P)$, it is thus sufficient that

$$\begin{aligned} \mathfrak{M}(B, P) &\geq \frac{1}{1 + \frac{2\alpha}{\varepsilon}} \cdot \left(1 - \bar{\gamma} \cdot \frac{TK}{TK + \varepsilon} \right) \\ &= \varepsilon \cdot \frac{\varepsilon + (1 - \bar{\gamma})TK}{(\varepsilon + 2\alpha)(\varepsilon + TK)}. \end{aligned} \quad (76)$$

In this case, we check that we have from (74)

$$\mathfrak{M}(B, Q) \geq (1 - \alpha)\mathfrak{M}(B, P) - KL(P, Q_0; B), \quad (77)$$

as claimed.

1.6 Additional formal results

One might ask what such a strong model of privacy allows to keep from the accuracy standpoint in general. Perhaps paradoxically at first sight, it is not hard to show that privacy can bring approximation guarantees on learning: *if we learn Q_ε within an ε -mollifier \mathcal{M} (hence, we get ε -privacy for sampling from Q_ε), then each time some Q_ε in \mathcal{M} accurately fits P , we are guaranteed that the one we learn also accurately fits P — albeit eventually more moderately —. We let $Q_\varepsilon(\cdot; \cdot)$ denote the density learned, where \cdot is the dataset argument.*

Lemma 9 Suppose $\exists \varepsilon$ -mollifier \mathcal{M} s.t. $Q_\varepsilon \in \mathcal{M}$, then $(\exists P, D', \delta : \text{KL}(P, Q_\varepsilon(; P')) \leq \delta) \Rightarrow (\forall D, \text{KL}(P, Q_\varepsilon(; P)) \leq \delta + \varepsilon)$.

Proof The proof is straightforward; we give it for completeness: for any dataset D , we have

$$\text{KL}(P, Q_\varepsilon(; P)) = \int_x \log \left(\frac{P}{Q_\varepsilon(; P)} \right) dP \quad (78)$$

$$= \int_x \log \left(\frac{P}{Q_\varepsilon(; P')} \right) dP + \int_x \log \left(\frac{Q_\varepsilon(; P)}{Q_\varepsilon(; P')} \right) dP \quad (79)$$

$$\leq \int_x \log \left(\frac{P}{Q_\varepsilon(; P')} \right) dP + \varepsilon \cdot \int_x dP \quad (80)$$

$$= \text{KL}(P, Q_\varepsilon(; P')) + \varepsilon \quad (81)$$

$$\leq \delta + \varepsilon, \quad (82)$$

from which we derive the statement of Lemma 9 assuming \mathcal{A} is ε -IP (the inequalities follow from the Lemma's assumption). ■

In the jargon of (computational) information geometry [1], we can summarize Lemma 9 as saying that if there exists an eligible¹ density in a small KL-ball relatively to P , we are guaranteed to find a density also in a small KL-ball relatively to P . This result is obviously good when the premises hold true, but it does not tell the full story when they do not. In fact, when there exists an eligible density outside a big KL-ball relatively to P , it is not hard to show using the same arguments as for the Lemma that we *cannot* find a good one, and this is not a feature of MBDE: this would hold regardless of the algorithm. This limitation is intrinsic to the likelihood ratio constraint of differential privacy, as the following Lemma shows. In the context of ε -DP, we assume that all input datasets have the same size, say m .

Lemma 10 Let \mathcal{A} denote an algorithm learning an ε -differentially private density. Denote $D \sim P$ an input of the algorithm and $\mathcal{Q}_\varepsilon(D)$ the set of all densities that can be the output of \mathcal{A} on input D , taking in considerations all internal randomisations of \mathcal{A} . Suppose there exists an input D' for which one of these densities is far from the target: $\exists D', \exists Q \in \mathcal{Q}_\varepsilon(D') : \text{KL}(P, Q(; D')) \geq \Delta$ for some "big" $\Delta > 0$. Then the output Q of \mathcal{A} obtained from **any** input $D \sim P$ satisfies: $\text{KL}(P, Q(; D)) \geq \Delta - m\varepsilon$.

Proof Denote D the actual input of \mathcal{A} . There exists a sequence \mathcal{D} of datasets of the same size, whose length is at most m , which transforms D into D' by repeatedly changing one observation in the current dataset: call it $\mathcal{D} = \{D, D_1, D_2, \dots, D_k, D'\}$, with $k \leq m - 1$. Denote $Q(; D'')$ any

¹Within the chosen ε -mollifier.

element of $\mathcal{Q}_\varepsilon(D'')$ for $D'' \in \mathcal{D}$. Since \mathcal{A} is ε -differentially private, we have:

$$\Delta \leq \mathbf{KL}(P, Q(; D')) \tag{83}$$

$$= \int_x \log \left(\frac{P}{Q(; D')} \right) dP \tag{84}$$

$$= \int_x \log \left(\frac{P}{Q(; D)} \right) dP + \int_x \log \left(\frac{Q(; D)}{Q(; D_1)} \right) dP + \sum_{j=1}^{k-1} \int_x \log \left(\frac{Q(; D_j)}{Q(; D_{j+1})} \right) dP + \int_x \log \left(\frac{Q(; D_k)}{Q(; D')} \right) dP \tag{85}$$

$$= \mathbf{KL}(P, Q(; D)) + \int_x \log \left(\frac{Q(; D)}{Q(; D_1)} \right) dP + \sum_{j=1}^{k-1} \int_x \log \left(\frac{Q(; D_j)}{Q(; D_{j+1})} \right) dP + \int_x \log \left(\frac{Q(; D_k)}{Q(; D')} \right) dP \tag{86}$$

$$\leq \mathbf{KL}(P, Q(; D)) + m\varepsilon, \tag{87}$$

from which we derive the statement of Lemma 10. ■

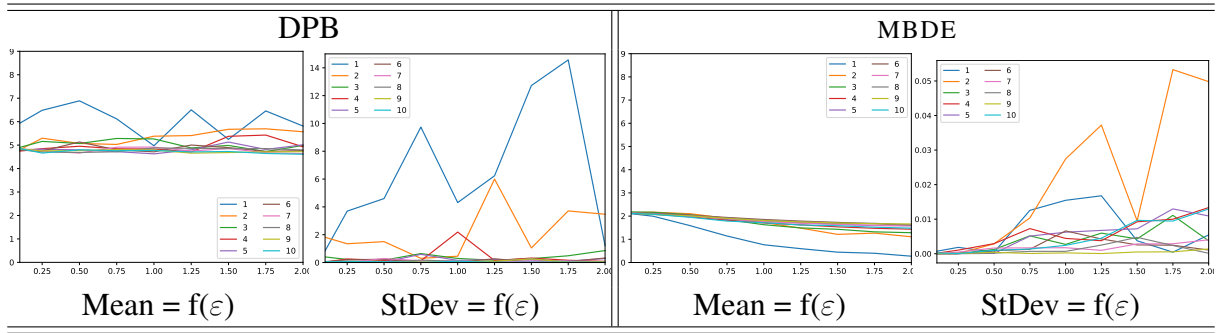


Figure 1: NLL metrics (mean and standard deviation) on the 1D random Gaussian problem for DPB (left pane) and MBDE (right pane), for a varying number of $m = 1, \dots, 10$ random Gaussians. The lower the better on each metric. Remark the different scales for StDev (see text).

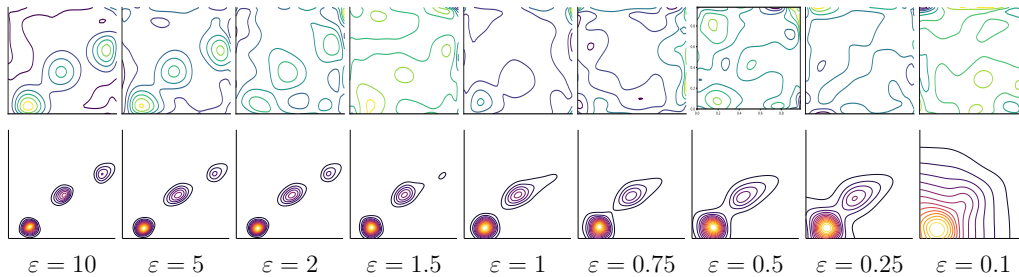


Figure 2: Randomly placed Gaussian convergence comparison for DPB (upper) against MBDE (lower).

2 Additional experiments

We provide here additional results to the main file. Figure 1 provides NLL values for the random 1D Gaussian problem. Figure 2 displays that picking Q_0 a standard Gaussian does not prevent to obtain good results — and beat DPB — when sampling random Gaussians.

References

- [1] J.-D. Boissonnat, F. Nielsen, and R. Nock. Bregman voronoi diagrams. *DCG*, 44(2):281–307, 2010.