

A Note on Reverse Pinsker Inequalities

Olivier Binette

Université du Québec à Montréal

Email: olivier.binette@gmail.com

Abstract—A simple method is shown to provide optimal variational bounds on f -divergences with possible constraints on relative information extremums. Known results are refined or proved to be optimal as particular cases.

Index Terms—Kullback-Leibler divergence, reverse Pinsker inequalities, f -divergences, range of values, upper bounds

I. INTRODUCTION

This note is concerned with optimal upper bounds on relative entropy and other f -divergences in terms of the total variation distance and relative information extremums. When taking relative entropy as the f -divergence, such upper variational bounds have been referred to as *reverse Pinsker inequalities* [1], [2]. They are used in the optimal quantization of probability measures [2] and have also appeared in Bayesian nonparametrics for controlling the prior probability of relative entropy neighbourhoods (see e.g. (A.2) in [3]).

Our main theorem demonstrates a simple method that yields optimal “reverse Pinsker inequalities” for any f -divergence. This refines or shows the optimality of previously best known inequalities while avoiding arguments that are tuned to particular cases. In particular, Simic [4] uses a global upper bound on the Jensen function to bound relative entropy by a function of relative information extremums. Corollary 2 below refines their inequality to best possible. More recently, three different bounds on relative entropy involving the total variation distance have been proposed in Theorem 23 of [1] in Theorem 7 of [5] and in Theorem 1 of [6]. Our results show that the inequalities of [1] and [5] are in fact optimal in related contexts. Another direct application of the method improves Theorem 34 in [1], which is an upper bound on Rényi’s divergence in terms of the variational distance and relative information maximum, while providing a simpler proof for this type of inequality. Vajda’s well-known “range of values theorem” (see [7]–[11]) is also recovered as an application.

The rest of the paper is organized as follows. Section II presents the definitions and main results. Examples with particular f -divergences are provided in section III and proofs are given in section IV.

II. MAIN RESULTS

Let (P, Q) be a pair of probability measures defined on a common measurable space. It is assumed throughout that $P \ll Q$. Given a convex function $f : [0, \infty) \rightarrow (-\infty, \infty]$ such that $f(1) = 0$, the f -divergence between P and Q is defined as

$$D_f(P||Q) = \mathbb{E}_Q \left[f \left(\frac{dP}{dQ} \right) \right]. \quad (1)$$

In particular, the relative entropy $D(P||Q)$ and the total variation distance $D_{TV}(P, Q) = \sup_A |P(A) - Q(A)|$ correspond to the cases $f(t) = t \log(t)$ and $f(t) = \frac{1}{2}|t - 1|$ respectively.

For fixed $\delta \geq 0$, $m \geq 0$ and $M \leq \infty$, we consider the set $\mathcal{A}(\delta, m, M)$ of all probability measure pairs (P, Q) defined on a common measurable space and respecting the conditions : $P \ll Q$,

$$\text{ess inf } \frac{dP}{dQ} = m, \quad \text{ess sup } \frac{dP}{dQ} = M \quad (2)$$

and

$$D_{TV}(P, Q) = \delta. \quad (3)$$

Here ess inf and ess sup represent the essential infimum and supremum taken with respect to Q .

The following theorem provides the best upper bound on the f -divergence over the class $\mathcal{A}(\delta, m, M)$ determined by (2) and (3).

Theorem 1. *If $\delta \geq 0$, $m \geq 0$ and $M < \infty$ are such that $\mathcal{A}(\delta, m, M) \neq \emptyset$, then*

$$\sup_{(P, Q) \in \mathcal{A}(\delta, m, M)} D_f(P||Q) = \delta \left(\frac{f(m)}{1-m} + \frac{f(M)}{M-1} \right). \quad (4)$$

Remark 1. In the case where $m = 1$ or $M = 1$, any $(P, Q) \in \mathcal{A}(\delta, m, M)$ must be such that $\delta = D_{TV}(P, Q) = 0$. The right hand side of (4) is then understood as being equal to 0.

We can obtain from Theorem 1 tight bounds for more general families of distributions. Consider for instance

$$\mathcal{B}(m, M) = \bigcup_{\delta \geq 0} \mathcal{A}(\delta, m, M) \quad (5)$$

and

$$\mathcal{C}(\delta) = \bigcup_{\substack{m \in [0, 1] \\ M \in [1, \infty]}} \mathcal{A}(\delta, m, M). \quad (6)$$

Using the first family, Corollary 2 below provides the range of D_f as a function of relative information bounds.

Corollary 2. *If $m \geq 0$ and $M < \infty$ are such that $\mathcal{B}(m, M) \neq \emptyset$, then*

$$\sup_{(P, Q) \in \mathcal{B}(m, M)} D_f(P||Q) = \frac{(M-1)f(m) + (1-m)f(M)}{M-m}. \quad (7)$$

Using the second family (6), we re-obtain Theorem 4 of [1] (see also Lemma 11.1 in [12]). Taking the union over possible values of δ also yields Vajda’s well-known “range of values theorem” (see [7]–[11]).

Corollary 3. *If $0 \leq \delta \leq 1$, then*

$$\sup_{(P,Q) \in \mathcal{C}(\delta)} D_f(P\|Q) = \delta \left(f(0) + \lim_{M \rightarrow \infty} \frac{f(M)}{M} \right). \quad (8)$$

Remark 2. Theorem 1 generalizes Theorem 23 in [1] with $f(t) = t \log(t)$ for the relative entropy: the upper bounds obtained are the same in this case. The proofs also share similarities. A decomposition equivalent to (12) is used in [1] and their proof is concluded by using the monotonicity of the function $t \mapsto t \log(t)/(1-t)$, continuously extended at 0 and 1.

III. EXAMPLES

This section lists applications to particular f -divergences and follows the standard definitions of [1]. The bounds obtained are compared to similar inequalities recently shown in the literature.

A. Relative entropy (Kullback-Leibler divergence)

The relative entropy corresponds to $f(t) = t \log(t)$ in (1) and is denoted $D(P\|Q)$. The results are more neatly stated in this case as functions of $a = \text{ess inf } \frac{dQ}{dP} = M^{-1}$ and $b = \text{ess sup } \frac{dQ}{dP} = m^{-1}$, assuming both quantities are well defined. Theorem 1 then shows

$$\sup_{(P,Q) \in \mathcal{A}(\delta, m, M)} D(P\|Q) = \delta \left(\frac{\log(a)}{a-1} + \frac{\log(b)}{1-b} \right). \quad (9)$$

In particular, the resulting upper bound on $D(P\|Q)$ is Theorem 23 of [1]. Letting $b \rightarrow \infty$ gives the related Theorem 7 of [5] and the inequality presented therein is consequently optimal over $\bigcup_{0 \leq m \leq 1} \mathcal{A}(\delta, m, M)$.

Also, Corollary 2 yields

$$\sup_{(P,Q) \in \mathcal{B}(m, M)} D(P\|Q) = \frac{(a-1) \log(b) + (1-b) \log(a)}{b-a}.$$

For comparison, Theorem I of [4] (which also appears as Theorem I in [13] and is related to results in [14], [15]) provides the weaker upper bound

$$\frac{a \log(b) - b \log(a)}{b-a} + \log \left(\frac{b-a}{\log(b) - \log(a)} \right) - 1$$

on $D(P\|Q)$ over $(P, Q) \in \mathcal{B}(m, M)$ as an application of their ‘‘best possible global bound’’ for the Jensen functional.

B. Hellinger divergence of order α

Let $\alpha \in (0, 1) \cup (1, \infty)$ and $f(t) = (t^\alpha - 1)/(\alpha - 1)$. The corresponding divergence is denoted $\mathcal{H}_\alpha(P\|Q)$. Theorem 1 shows in this case

$$\sup_{(P,Q) \in \mathcal{A}(\delta, m, M)} \mathcal{H}_\alpha(P\|Q) = \frac{\delta}{1-\alpha} \left(\frac{1-m^\alpha}{1-m} - \frac{M^\alpha-1}{M-1} \right).$$

When $\alpha = 2$, $\mathcal{H}_\alpha = D_{\chi^2}$ is the χ^2 divergence and the above can be rewritten as

$$\sup_{(P,Q) \in \mathcal{A}(\delta, m, M)} D_{\chi^2}(P\|Q) = \delta(M-m).$$

For comparison, Example 6 of Theorem 5 in [1] is the weaker inequality

$$D_{\chi^2}(P\|Q) \leq 2\delta \max\{M-1, 1-m\}.$$

C. Rényi's divergence

Also related is Rényi's α -divergence, defined as

$$D_\alpha(P\|Q) = \frac{1}{\alpha-1} \log(1 + (\alpha-1)\mathcal{H}_\alpha(P\|Q))$$

and which is a monotonous transform of \mathcal{H}_α . Correspondingly we obtain

$$D_\alpha(P\|Q) \leq \frac{1}{\alpha-1} \log \left(1 + \delta \left(\frac{M^\alpha-1}{M-1} - \frac{1-m^\alpha}{1-m} \right) \right).$$

Taking $m = 0$ recovers Theorem 34 of [1]. Their inequality, which is also appears in Theorem 3 of [16] for $\alpha > 2$, is improved when $m > 0$.

IV. PROOFS

The starting point of our analysis is the following simple known application of convexity.

Lemma 4. *Let κ be a random variable with values in a bounded interval $I = [a, b]$, let $\varphi : I \rightarrow (-\infty, \infty]$ be a convex function and let $\bar{\alpha} = (b - \mathbb{E}[\kappa])/(b - a)$. Then*

$$\mathbb{E}[\varphi(\kappa)] \leq \bar{\alpha}\varphi(a) + (1-\bar{\alpha})\varphi(b). \quad (10)$$

Proof. Let α be the non-negative random variable such that $\kappa = \alpha a + (1-\alpha)b$. Then $\mathbb{E}[\alpha] = \bar{\alpha}$ and by convexity of φ we find

$$\begin{aligned} \mathbb{E}[\varphi(\kappa)] &\leq \mathbb{E}[\alpha\varphi(a) + (1-\alpha)\varphi(b)] \\ &= \bar{\alpha}\varphi(a) + (1-\bar{\alpha})\varphi(b). \end{aligned}$$

□

As a particular case, we obtain a bound on the total variation distance that is of use in the proof of Theorem 1.

Corollary 5. *If $m \geq 0$, $M < \infty$ and $(P, Q) \in \mathcal{B}(m, M)$, then*

$$D_{TV}(P, Q) \leq \frac{(M-1)(1-m)}{M-m}. \quad (11)$$

Proof. Lemma 4, applied with $\kappa = \frac{dP}{dQ}$, $\varphi(x) = |x-1|$, $a = m$ and $b = M$, shows that

$$\begin{aligned} \frac{1}{2} \mathbb{E}_Q \left[\left| \frac{dP}{dQ} - 1 \right| \right] &\leq \frac{1}{2} \left\{ \frac{M-1}{M-m} |m-1| + \frac{1-m}{M-m} |M-1| \right\} \\ &= \frac{(M-1)(1-m)}{M-m}. \end{aligned}$$

□

We now proceed with the proof of Theorem 1.

Proof of Theorem 1. Let $(P, Q) \in \mathcal{A}(\delta, m, M)$. If $A = \{x \mid \frac{dP}{dQ}(x) \leq 1\}$, then $\delta = Q(A) - P(A)$ and we may write

$$\begin{aligned} D_f(P\|Q) &= Q(A) \mathbb{E}_Q \left[f \left(\frac{dP}{dQ} \right) \Big| A \right] \\ &\quad + Q(A^c) \mathbb{E}_Q \left[f \left(\frac{dP}{dQ} \right) \Big| A^c \right]. \end{aligned} \quad (12)$$

To bound the first term on the right-hand side of (12), note that $\mathbb{E}_Q \left[\frac{dP}{dQ} \Big| A \right] = \frac{P(A)}{Q(A)}$ and that $x \in A$ implies $m \leq \frac{dP}{dQ}(x) \leq 1$.

An application of Lemma 4, using the fact that $f(1) = 0$, therefore yields

$$\begin{aligned} \mathbb{E}_Q \left[f \left(\frac{dP}{dQ} \right) \middle| A \right] &\leq \frac{1 - \frac{P(A)}{Q(A)}}{1 - m} f(m) \\ &= \frac{\delta f(m)}{Q(A)(1 - m)}. \end{aligned} \quad (13)$$

The second term is similarly bounded as to obtain

$$\begin{aligned} \mathbb{E}_Q \left[f \left(\frac{dP}{dQ} \right) \middle| A^c \right] &\leq \left(1 - \frac{M - \frac{P(A^c)}{Q(A^c)}}{M - 1} \right) f(M) \\ &= \frac{\delta f(M)}{Q(A^c)(M - 1)}. \end{aligned} \quad (14)$$

Together with (12), the inequalities (13) and (14) show that

$$D_f(P||Q) \leq \delta \left(\frac{f(m)}{1 - m} + \frac{f(M)}{M - 1} \right) \quad (15)$$

whenever $(P, Q) \in \mathcal{A}(\delta, m, M)$.

We now show that the supremum of (4) equals this bound. If $\delta = 0$, then the upper bound given by (4) is zero and the supremum trivially attains this bound. If $\delta = 1$, then $\mathcal{A}(\delta, m, M) = \emptyset$ and the statement of Theorem 1 is trivially satisfied. We can therefore assume $0 < \delta < 1$. Let $q = \frac{M-1}{M-m}$, $p = mq$, $t = \delta(M - m)[(M - 1)(1 - m)]^{-1}$ and consider the pair of discrete measures

$$\begin{cases} P = (tp, t(1 - p), 1 - t), \\ Q = (tq, t(1 - q), 1 - t). \end{cases} \quad (16)$$

Corollary 5 ensures $0 < t < 1$ and thus P and Q are probability measures. It is also straightforward to verify that $(P, Q) \in \mathcal{A}(\delta, m, M)$ with $t(q - p) = \delta$, $p/q = m$ and $(1 - p)/(1 - q) = M$. Some algebraic manipulations then show

$$\begin{aligned} D_f(P||Q) &= tqf \left(\frac{p}{q} \right) + t(1 - q)f \left(\frac{1 - p}{1 - q} \right) \\ &= \delta \left(\frac{f(m)}{1 - m} + \frac{f(M)}{M - 1} \right). \end{aligned}$$

□

Proof of Corollary 2. Combining Corollary 5 with equation (4) of Theorem 1 yields the upper bound

$$D_f(P||Q) \leq \frac{(M - 1)f(m) + (1 - m)f(M)}{M - m}$$

for every $(P, Q) \in \mathcal{B}(m, M)$. To see that the supremum over $\mathcal{B}(m, M)$ equals this bound, it suffices to let $\delta \rightarrow (M - 1)(1 - m)/(M - m)$ in (4). □

Proof of Corollary 3. Some care has to be taken when considering the elements of $\mathcal{A}(\delta, 0, \infty)$. To see that the right-hand side of (8) also upper bounds the elements of this set, we again use the decomposition (12). The first term is treated as in (13).

For the second term, let $\frac{dP}{dQ} \wedge K = \min\{\frac{dP}{dQ}, K\}$. By Fatou's lemma and Lemma 4, using that $f(1) = 0$,

$$\begin{aligned} \mathbb{E}_Q \left[f \left(\frac{dP}{dQ} \right) \middle| A^c \right] &\leq \liminf_{K \rightarrow \infty} \mathbb{E}_Q \left[f \left(\frac{dP}{dQ} \wedge K \right) \middle| A^c \right] \\ &\leq \liminf_{K \rightarrow \infty} \frac{\mathbb{E}_Q \left[\frac{dP}{dQ} \wedge K \middle| A^c \right] - 1}{K - 1} f(K). \end{aligned}$$

By the monotone convergence theorem,

$$\lim_{K \rightarrow \infty} \mathbb{E}_Q \left[\frac{dP}{dQ} \wedge K \middle| A^c \right] = \frac{P(A^c)}{Q(A^c)}$$

and hence

$$\mathbb{E}_Q \left[f \left(\frac{dP}{dQ} \right) \middle| A^c \right] \leq \frac{\delta}{Q(A^c)} \lim_{M \rightarrow \infty} \frac{f(M)}{M - 1}.$$

We note that $\lim_{M \rightarrow \infty} \frac{f(M)}{M - 1}$ exists by convexity of f and can be infinite. The required upper bound on $D_f(P||Q)$ is then obtained as in the proof of Theorem 1.

To see that the supremum equals this bound, it suffices to let $M \rightarrow \infty$ in Theorem 1. □

ACKNOWLEDGMENT

The author would like to thank Alexis Langlois-Rémillard and Jean-François Coeurjolly for helpful comments and suggestions.

REFERENCES

- [1] I. Sason and S. Verdú, “ f -divergence inequalities,” *IEEE Transactions on Information Theory*, vol. 62, no. 11, pp. 5973–6006, 2016.
- [2] G. Böcherer and B. C. Geiger, “Optimal quantization for distribution synthesis,” *IEEE Transactions on Information Theory*, vol. 62, no. 11, pp. 6162–6172, 2016.
- [3] O. Binette and S. Guillote, “Bayesian Nonparametrics for Directional Statistics,” *arXiv e-prints*, 2018, arXiv:1807.00305v2.
- [4] S. Simic, “Best possible global bounds for jensen's inequality,” *Applied Mathematics and Computation*, vol. 215, no. 6, pp. 2224 – 2228, 2009.
- [5] S. Verdú, “Total variation distance and the distribution of relative information,” in *2014 Information Theory and Applications Workshop (ITA)*, 2014, pp. 1–3.
- [6] I. Sason, “On Reverse Pinsker Inequalities,” *ArXiv e-prints*, 2015, arXiv:1503.07118.
- [7] I. Vajda, “On the f -divergence and singularity of probability measures,” *Periodica Mathematica Hungarica*, vol. 2, no. 1, pp. 223–234, 1972.
- [8] F. Liese and I. Vajda, “On divergences and informations in statistics and information theory,” *IEEE Transactions on Information Theory*, vol. 52, no. 10, pp. 4394–4412, 2006.
- [9] I. Vajda, “On metric divergences of probability measures,” *Kybernetika*, vol. 45, no. 6, pp. 885–900, 2009.
- [10] P. Kumar and L. Hunter, “On an information divergence measure and information inequalities,” *Carpathian Journal of Mathematics*, vol. 20, no. 1, pp. 51–66, 2004.
- [11] P. Kumar and S. Chhina, “A symmetric information divergence measure of the csiszár's f -divergence class and its bounds,” *Computers & Mathematics with Applications*, vol. 49, no. 4, pp. 575 – 588, 2005.
- [12] A. Basu, H. Shioya, and C. Park, *Statistical inference: the minimum distance approach*. CRC Press, 2011.
- [13] S. Simic, “Sharp global bounds for jensen's inequality,” *Rocky Mountain J. Math.*, vol. 41, no. 6, pp. 2021–2031, 2011.
- [14] —, “On certain new inequalities in information theory,” *Acta Mathematica Hungarica*, vol. 124, no. 4, pp. 353–361, 2009.
- [15] —, “Jensen's inequality and new entropy bounds,” *Applied Mathematics Letters*, vol. 22, no. 8, pp. 1262 – 1265, 2009.
- [16] I. Sason and S. Verdú, “Upper bounds on the relative entropy and rényi divergence as a function of total variation distance for finite alphabets,” in *2015 IEEE Information Theory Workshop - Fall (ITW)*, 2015, pp. 214–218.

- [17] I. Csiszár, "On topological properties of f-divergences," *Studia Scientiarum Mathematicarum Hungarica*, pp. 329–339, 1967.
- [18] I. J. Taneja and P. Kumar, "Relative information of type s, csiszár's f-divergence, and information inequalities," *Information Sciences*, vol. 166, no. 1, pp. 105 – 125, 2004.
- [19] S. M. Ali and S. D. Silvey, "A general class of coefficients of divergence of one distribution from another," *Journal of the Royal Statistical Society. Series B (Methodological)*, vol. 28, pp. 131–142, 1966.