# On $f$-divergences between Cauchy distributions[*]

Frank Nielsen

Sony Computer Science Laboratories Inc.

E-mail: `Frank.Nielsen@acm.org`

Kazuki Okamura

Department of Mathematics, Faculty of Science, Shizuoka University

E-mail: `okamura.kazuki@shizuoka.ac.jp`

**Abstract**

We prove that the $f$-divergences between univariate Cauchy distributions are all symmetric, and can be expressed as strictly increasing scalar functions of the symmetric chi-squared divergence. We report the corresponding scalar functions for the total variation distance, the Kullback-Leibler divergence, the squared Hellinger divergence, and the Jensen-Shannon divergence among others. Next, we give conditions to expand the $f$-divergences as converging infinite series of higher-order power chi divergences, and illustrate the criterion for converging Taylor series expressing the $f$-divergences between Cauchy distributions. We then show that the symmetric property of $f$-divergences holds for multivariate location-scale families with prescribed matrix scales provided that the standard density is even which includes the cases of the multivariate normal and Cauchy families. However, the $f$-divergences between multivariate Cauchy densities with different scale matrices are shown asymmetric. Finally, we present several metrizations of $f$-divergences between univariate Cauchy distributions and further report geometric embedding properties of the Kullback-Leibler divergence.

Keywords: Univariate and multivariate location-scale families; Cauchy distributions; Circular Cauchy distributions; Wrapped Cauchy distributions; Log-Cauchy distributions; Complex analysis; Maximal invariant; Information geometry; Divergence; Hilbert embeddings; Elliptic integrals.

# Contents

# 1  Introduction

Let $\mathbb{R}$, $\mathbb{R}_+$ and $\mathbb{R}_{++}$ be the sets of real numbers, non-negative real numbers, and positive real numbers, respectively. The probability density function of a Cauchy distribution (also called a Lorentzian distribution [24] in physics) is

$$p_{l,s}(x) := \frac{1}{\pi s \left(1 + \left(\frac{x-l}{s}\right)^2\right)} = \frac{s}{\pi(s^2 + (x-l)^2)},$$

where $l \in \mathbb{R}$ denotes the location parameter and $s \in \mathbb{R}_{++}$ the scale parameter of the Cauchy distribution, and $x \in \mathbb{R}$. The space of Cauchy distributions form a location-scale family

$$\mathcal{C} = \left\{ p_{l,s}(x) := \frac{1}{s} p\left(\frac{x-l}{s}\right) \ : \ (l,s) \in \mathbb{R} \times \mathbb{R}_{++} \right\},$$

with standard density

$$p(x) := \frac{1}{\pi(1 + x^2)}. \tag{1}$$

To measure the dissimilarity between two continuous probability distributions $P$ and $Q$, we consider the class of statistical $f$-divergences [14, 58] between their corresponding probability densities functions $p(x)$ and $q(x)$ assumed to be strictly positive on $\mathbb{R}$:

$$I_f(p : q) := \int_{\mathbb{R}} p(x) f\left(\frac{q(x)}{p(x)}\right) \mathrm{d}x,$$

where $f(u)$ is a convex function on $(0, \infty)$, strictly convex at $u = 1$ (to ensure reflexivity $I_f(p : q) = 0$ iff $p = q$), and satisfying $f(1) = 0$ (to ensure positive-definiteness $I_f(p, q) \geq 0$ since by Jensen's

inequality we have $I_f(p : q) \geq f(1) = 0$). The Kullback-Leibler divergence (KLD also called relative entropy) is an $f$-divergence obtained for $f_{\mathrm{KL}}(u) = -\log u$. In general, the $f$-divergences are oriented dissimilarities: $I_f(p : q) \neq I_f(q : p)$ (eg., the KLD). The reverse $f$-divergence $I_f(q : p)$ can be obtained as a forward $f$-divergence for the conjugate function $f^*(u) := uf\left(\frac{1}{u}\right)$ (convex with $f^*(1) = 0$): $I_f(q : p) = I_{f^*}(p : q)$. We have $I_f = I_g$ when there exists $\lambda \in \mathbb{R}$ such that $f(u) = g(u) + \lambda(u - 1)$. Thus an $f$-divergence is symmetric when there exists a real $\lambda$ such that $f(u) = uf\left(\frac{1}{u}\right) + \lambda(u - 1)$, and $f$-divergences can always be symmetrized by taking the generator $s_f(u) = \frac{1}{2}(f(u) + uf\left(\frac{1}{u}\right))$. In general, calculating the definite integrals of $f$-divergences is non trivial: For example, the formula for the KLD between Cauchy densities was only recently obtained [11]:

$$
\begin{aligned}
D_{\mathrm{KL}}(p_{l_1,s_1} : p_{l_2,s_2}) &:= I_{f_{\mathrm{KL}}}(p : q) = \int p_{l_1,s_1}(x) \log \frac{p_{l_1,s_1}(x)}{p_{l_2,s_2}(x)} \mathrm{d}x \\
&= \log\left(\frac{(s_1 + s_2)^2 + (l_1 - l_2)^2}{4s_1 s_2}\right).
\end{aligned}
$$

Let $\lambda = (\lambda_1 = l, \lambda_2 = s)$. Then we can rewrite the KLD formula as

$$
D_{\mathrm{KL}}(p_{\lambda_1} : p_{\lambda_2}) = \log\left(1 + \frac{1}{2}\chi(\lambda_1, \lambda_2)\right), \tag{2}
$$

where

$$
\chi(\lambda, \lambda') := \frac{(\lambda_1 - \lambda'_1)^2 + (\lambda_2 - \lambda'_2)^2}{2\lambda_2 \lambda'_2} = \frac{\|\lambda - \lambda'\|^2}{2\lambda_2 \lambda'_2}.
$$

See (3) for complex representations.

We observe that the KLD between Cauchy distributions is symmetric: $D_{\mathrm{KL}}(p_{l_1,s_1} : p_{l_2,s_2}) = D_{\mathrm{KL}}(p_{l_2,s_2} : p_{l_1,s_1})$. Let

$$
D_\chi^N(p : q) := \int \frac{(p(x) - q(x))^2}{q(x)} \mathrm{d}x \text{ and } D_\chi^P(p : q) := \int \frac{(p(x) - q(x))^2}{p(x)} \mathrm{d}x
$$

denote the Neyman and Pearson chi-squared divergences between densities $p(x)$ and $q(x)$. These divergences are $f$-divergences [58] for the generators $f_\chi^P(u) = (u - 1)^2$ and $f_\chi^N(u) = \frac{1}{u}(u - 1)^2$, respectively. The $\chi^2$-divergences between Cauchy densities are symmetric [54]:

$$
D_\chi(p_{\lambda_1} : p_{\lambda_2}) := D_\chi^N(p_{\lambda_1} : p_{\lambda_2}) = D_\chi^P(p_{\lambda_1} : p_{\lambda_2}) = \chi(\lambda_1, \lambda_2),
$$

hence the naming of the function $\chi(\cdot, \cdot)$. Notice that we have

$$
\chi(p_{\lambda_1} : p_{\lambda_2}) = \rho(\lambda_1)\rho(\lambda_2)\frac{1}{2}D_E^2(\lambda_1, \lambda_2),
$$

where $D_E(\lambda_1, \lambda_2) := \sqrt{(\lambda_2 - \lambda_1)^\top(\lambda_2 - \lambda_1)}$ and $(\lambda_2 - \lambda_1)^\top$ denotes the transpose of the vector $(\lambda_2 - \lambda_1)$. That is, the function $\chi$ is a conformal half squared Euclidean divergence [62, 59] with conformal factor $\rho(\lambda) := \frac{1}{\lambda_2}$. When the Neyman and Pearson chi-squared divergences are not symmetric, we define the chi-squared symmetric divergence as

$$
D_\chi(p : q) = D_\chi^N(p : q) + D_\chi^P(p : q) = \int \frac{(p(x) + q(x))(p(x) - q(x))^2}{p(x)q(x)} \mathrm{d}x.
$$

3

In this work, we first prove in §2 that all $f$-divergences between univariate Cauchy distributions are symmetric (Theorem 1) and can be expressed as a strictly increasing scalar function of the chi-squared divergence (Theorem 2). We illustrate this result by reporting the corresponding functions for the total variation distance, the Kullback-Leibler divergence, the LeCam-Vincze divergence, the squared Hellinger divergence, and the Jensen-Shannon divergence. Further results for the $f$-divergences between the circular Cauchy, wrapped Cauchy and log-Cauchy distributions based on the invariance properties of the $f$-divergences are presented in §3. We report conditions to expand the $f$-divergences as infinite series of higher-order chi divergences and instantiate the results for the Cauchy distributions in §5. In §4, we then show that the symmetric property of $f$-divergence holds for multivariate location-scale families including the normal and Cauchy families with prescribed matrix scales provided that the standard density is even, but does not hold for general case of different matrix scales. We consider metrizations of the square roots of the KLD and the Bhattacharyya divergences in §6. Finally in §7 we investigate geometric properties of these metrics.

In the appendix, we first recall the information geometry of the Cauchy family in §A, explain the relationship of Cauchy distributions with the Möbius and Boole transformations in §B, give alternative simpler proofs of the Kullback-Leibler divergence (§C) and chi-squared divergence (§D) between Cauchy distributions, report a closed-form formula for the total variation distance between densities of a location-scale family in §E. In §F, we also recall the complete elliptic integrals, which are used in the proof of the metrization of the square root of the Bhattacharyya divergence. We discuss isometric embedding into a Hilbert space of the square root of the KLD in §G. We finally give a code snippet for calculating some converging truncated Taylor series of $f$-divergences between Cauchy distributions in §H.

## 2 Symmetric property of the $f$-divergences between univariate Cauchy distributions

Consider the location-scale non-abelian group LS(2) which can be represented as a matrix group [55]. A group element $g_{l,s}$ is represented by a matrix element $M_{l,s} = \begin{bmatrix} s & l \\ 0 & 1 \end{bmatrix}$ for $(l, s) \in \mathbb{R} \times \mathbb{R}_{++}$. The group operation $g_{l_{12},s_{12}} = g_{l_1,s_1} \times g_{l_2,s_2}$ corresponds to a matrix multiplication $M_{l_{12},s_{12}} = M_{l_1,s_1} \times M_{l_2,s_2}$ (with the group identity element $g_{0,1}$ being the matrix identity). A location-scale family is defined by the action of the location-group on a standard density $p(x) = p_{0,1}(x)$. That is, density $p_{l,s}(x) = g_{l,s}.p(x)$ where '.' denotes the action. We have the following invariance for the $f$-divergences between any two densities of a location-scale family [55] (including the Cauchy family):

$$I_f(g.p_{l_1,s_1} : g.p_{l_2,s_2}) = I_f(p_{l_1,s_1} : p_{l_2,s_2}), \forall g \in \mathrm{LS}(2).$$

Thus we have

$$I_f(p_{l_1,s_1} : p_{l_2,s_2}) = I_f\left(p : p_{\frac{l_2-l_1}{s_1},\frac{s_2}{s_1}}\right) = I_f\left(p_{\frac{l_1-l_2}{s_2},\frac{s_1}{s_2}} : p\right).$$

Therefore, we may always consider the calculation of the $f$-divergence between the standard density and another density of the location-scale family. For example, we check that

$$\chi((l_1, s_1), (l_2, s_2)) = \chi\left((0, 1), \left(\frac{l_2 - l_1}{s_1}, \frac{s_2}{s_1}\right)\right)$$

4

since $\chi((0,1),(l,s)) = \frac{(s-1)^2 + l^2}{2s}$. If we assume that the standard density $p$ is such that $E_p[X] = \int xp(x)\mathrm{d}x = 0$ and $E_p[X^2] = \int x^2 p(x)\mathrm{d}x = 1$ (hence unit variance), then the random variable $Y = \mu + \sigma X$ has mean $E[Y] = \mu$ and standard deviation $\sigma(Y) = \sqrt{E[(Y-\mu)^2]} = \sigma$. However, the expectation and variance of Cauchy distributions are not defined, hence we preferred $(l,s)$ parameterization over the $(\mu, \sigma^2)$ parameterization, where $l$ denotes the median and $s$ the probable error for the Cauchy location-scale family [46].

## 2.1 $f$-divergences between densities of a location family

Let us first prove that $f$-divergences between densities of a location family with *even* standard density are symmetric:

**Proposition 1** *Let $\mathcal{L}_p = \{p(x-l) \ : \ l \in \mathbb{R}\}$ denote a location family with even standard density (i.e., $p(-x) = p(x)$) on the support $\mathcal{X} = \mathbb{R}$. Then all $f$-divergences between two densities $p_{l_1}$ and $p_{l_2}$ of $\mathcal{L}$ are symmetric: $I_f(p_{l_1} : p_{l_2}) = I_f(p_{l_2} : p_{l_1})$.*

Proof.    Consider the change of variable $l_1 - x = y - l_2$ (so that $x - l_2 = l_1 - y$) with $\mathrm{d}x = -\mathrm{d}y$ and let us use the property that $p(z - l_1) = p(l_1 - z)$ since $p$ is an even standard density. We have:

$$
\begin{aligned}
I_f(p_{l_1} : p_{l_2}) \ &:= \ \int_{-\infty}^{+\infty} p(x - l_1) f\left(\frac{p(x - l_2)}{p(x - l_1)}\right) \mathrm{d}x, \\
&= \ \int_{+\infty}^{-\infty} p(l_1 - x) f\left(\frac{p(x - l_2)}{p(l_1 - x)}\right) (-\mathrm{d}y), \\
&= \ \int_{-\infty}^{+\infty} p(y - l_2) f\left(\frac{p(x - l_2)}{p(y - l_2)}\right) \mathrm{d}y, \\
&= \ \int_{-\infty}^{+\infty} p(y - l_2) f\left(\frac{p(l_1 - y)}{p(y - l_2)}\right) \mathrm{d}y, \\
&= \ \int_{-\infty}^{+\infty} p(y - l_2) f\left(\frac{p(y - l_1)}{p(y - l_2)}\right) \mathrm{d}y, \\
&=: \ I_f(p_{l_2} : p_{l_1}).
\end{aligned}
$$

QED.

Thus $f$-divergences between location Cauchy densities are symmetric since $p(x) = p(-x)$ for the standard Cauchy density of Eq. 1.

## 2.2 $f$-divergences between Cauchy distributions are symmetric

Let $\|\lambda\| = \sqrt{\lambda_1^2 + \lambda_2^2}$ denote the Euclidean norm of a 2D vector $\lambda = (\lambda_1, \lambda_2) \in \mathbb{R}^2$. We state the main theorem:

**Theorem 1** *All $f$-divergences between univariate Cauchy distributions $p_\lambda$ and $p_{\lambda'}$ with $\lambda = (l,s)$ and $\lambda' = (l', s')$ are symmetric and can be expressed as*

$$
I_f(p_\lambda : p_{\lambda'}) = h_f\left(\chi(\lambda, \lambda')\right)
$$

5

*where*

$$\chi(\lambda, \lambda') := \frac{\|\lambda - \lambda'\|^2}{2\lambda_2 \lambda_2'}$$

*and $h_f : \mathbb{R}_+ \to \mathbb{R}_+$ is a function (with $h_f(0) = 0$).*

The proof does not yield explicit closed-form formula for the $f$-divergences as it can be in general difficult to calculate in closed forms, and relies on McCullagh's complex parametrization [46] $p_\theta$ of the parameter of the Cauchy density $p_{l,s}$ with $\theta = l + is$:

$$p_\theta(x) = \frac{|\text{Im}(\theta)|}{\pi|x - \theta|^2},$$

since $|x - (l + is)|^2 = ((x - l) + is)((x - l) - is) = (x - l)^2 + s^2$. The parameter space $\theta$ is the complex plane $\mathbb{C}$ where we identify $\bar{\theta}$ with $\theta$, and the Cauchy distributions are degenerated to Dirac distributions $\delta_l(x)$ whenever $s = 0$.

We make use of the special linear group $\text{SL}(2, \mathbb{R})$ for $\theta$ the complex parameter:

$$\text{SL}(2, \mathbb{R}) := \left\{ \begin{bmatrix} a & b \\ c & d \end{bmatrix} \; : \; a, b, c, d \in \mathbb{R}, ad - bc = 1 \right\}.$$

Let $A.\theta := \frac{a\theta + b}{c\theta + d}$ (real linear fractional transformations) be the action of $A = \begin{bmatrix} a & b \\ c & d \end{bmatrix} \in \text{SL}(2, \mathbb{R})$. McCullagh proved that if $X \sim \text{Cauchy}(\theta)$ then $A.X \sim \text{Cauchy}(A.\theta)$, where $\theta \in \mathbb{C}$ is identified with $\bar{\theta}$ (hence $\lambda(\theta) = (\text{Re}(\theta), |\text{Im}(\theta)|)$). For example, if $X \sim \text{Cauchy}(is)$ then $\frac{1}{X} \sim \text{Cauchy}(\frac{1}{is}) = \text{Cauchy}(-\frac{i}{s}) \equiv \text{Cauchy}(\frac{1}{s})$. Using the $\lambda = (l, s)$ parameterization, we have

$$l_A = \frac{(al + b)(cl + d) + acs^2}{(cl + d)^2 + c^2 s^2},$$

$$s_A = \left| \frac{(ad - bc)s}{(cl + d)^2 + c^2 s^2} \right|.$$

We can also define an action of $\text{SL}(2, \mathbb{R})$ to the real line $\mathbb{R}$ by $x \mapsto \frac{ax+b}{cx+d}$, $x \in \mathbb{R}$, where we interpret $-\frac{d}{c} \mapsto \frac{a}{c}$ if $c \neq 0$. We remark that $d \neq 0$ if $c = 0$. This map is bijective between $\mathbb{R}$. We have the following invariance:

**Lemma 1 (Invariance of Cauchy $f$-divergence under $\text{SL}(2, \mathbb{R})$)** *For any $A \in \text{SL}(2, \mathbb{R})$ and $\theta_1, \theta_2 \in \mathbb{H}$, we have*

$$I_f(p_{A.\theta_1} : p_{A.\theta_2}) = I_f(p_{\theta_1} : p_{\theta_2}).$$

Proof.    We prove the invariance by the change of variable in the integral. Let $D(\theta_1 : \theta_2) := I_f(p_{\theta_1} : p_{\theta_2})$. We have

$$D(A.\theta_1 : A.\theta_2) = \int_{\mathbb{R}} \frac{\text{Im}(A.\theta_1)}{\pi|x - A.\theta_1|^2} f\left( \frac{\text{Im}(A.\theta_2)|x - A.\theta_1|^2}{\text{Im}(A.\theta_1)|x - A.\theta_2|^2} \right) dx.$$

Since $A \in \text{SL}(2, \mathbb{R})$, we have

$$\text{Im}(A.\theta_i) = \frac{\text{Im}(\theta_i)}{|c\theta_i + d|^2}, \quad i \in \{1, 2\}.$$

6

If $x = A.y$ then $\mathrm{d}x = \frac{\mathrm{d}y}{|cy+d|^2}$, and

$$|A.y - A.\theta_i|^2 = \frac{|y - \theta_i|^2}{|cy+d|^2 \ |c\theta_i + d|^2}, \quad i \in \{1, 2\}.$$

Hence we get:

$$\int_{\mathbb{R}} f\left(\frac{\mathrm{Im}(A.\theta_2)|x - A.\theta_1|^2}{\mathrm{Im}(A.\theta_1)|x - A.\theta_2|^2}\right) \frac{\mathrm{Im}(A.\theta_2)}{\pi|x - A.\theta_1|^2}\mathrm{d}x = \int_{\mathbb{R}} f\left(\frac{\mathrm{Im}(\theta_2)|y - \theta_1|^2}{\mathrm{Im}(\theta_1)|y - \theta_2|^2}\right) \frac{\mathrm{Im}(\theta_2)}{\pi|y - \theta_2|^2}\mathrm{d}y,$$
$$= I_f\left(p_{\theta_1} : p_{\theta_2}\right).$$

QED.

Let us notice that the Cauchy family is the only univariate location-scale family that is also closed by inversion [34]: That is, if $X \sim \mathrm{Cauchy}(l, s)$ then $\frac{1}{X} \sim \mathrm{Cauchy}(l', s')$. Therefore our results are specific to the Cauchy family and not to any other location-scale family. However the characterization by [34] yields some applications. See Appendix B for details.

We now prove Theorem 1 using the notion of maximal invariants of Eaton [18] (Chapter 2) that will be discussed in §2.5.

Let us rewrite the function $\chi$ with complex arguments as:

$$\chi(z, w) := \frac{|z - w|^2}{2\,\mathrm{Im}(z)\mathrm{Im}(w)}, \quad z, w \in \mathbb{C}. \tag{3}$$

**Proposition 2 (McCullagh [46])** *The function $\chi$ defined in Eq. 3 is a maximal invariant for the action of the special linear group $\mathrm{SL}(2, \mathbb{R})$ to $\mathbb{H} \times \mathbb{H}$ defined by*

$$A.(z, w) := \left(\frac{az + b}{cz + d}, \frac{aw + b}{cw + d}\right), \quad A = \begin{bmatrix} a & b \\ c & d \end{bmatrix} \in \mathrm{SL}(2, \mathbb{R}), \ z, w \in \mathbb{H}.$$

*That is, we have*

$$\chi(A.z, A.w) = \chi(z, w), \quad A \in \mathrm{SL}(2, \mathbb{R}), \ z, w \in \mathbb{H},$$

*and it holds that for every $z, w, z', w' \in \mathbb{H}$ satisfying that $\chi(z', w') = \chi(z, w)$, there exists $A \in \mathrm{SL}(2, \mathbb{R})$ such that $(A.z, A.w) = (z', w')$.*

By Lemma 1 and Theorem 2.3 of [18], there exists a unique function $h_f : [0, \infty) \to [0, \infty)$ such that $h_f(\chi(z, w)) = D(z, w)$ for all $z, w \in \mathbb{H}$.

**Theorem 2** *The $f$-divergence between two univariate Cauchy densities is symmetric and expressed as a function of the chi-squared divergence:*

$$I_f(p_{\theta_1} : p_{\theta_2}) = I_f(p_{\theta_2} : p_{\theta_1}) = h_f(\chi(\theta_1, \theta_2)), \quad \theta_1, \theta_2 \in \mathbb{H}. \tag{4}$$

Therefore we have proven that the $f$-divergences between univariate Cauchy densities are all symmetric. Note that we have $h_f = h_{f^*}$. In general, the $f$-divergences between two Cauchy mixtures $m(x) = \sum_{i=1}^{k} w_i p_{l_i, s_i}(x)$ and $m'(x) = \sum_{i=1}^{k'} w_i' p_{l_i', s_i'}(x)$ are asymmetric (i.e., $I_f(m : m') \neq I_f(m' : m)$) except when $k = k' = 1$.

Similarly, we proved in Proposition 1 that all $f$-divergences between two densities of a location family with even standard density are symmetric. These $f$-divergences $I_f[p_{l_1} : p_{l_2}]$ can be expressed as a function $k_f$ of the the absolute value $|l_1 - l_2|$:

$$I_f(p_{l_1} : p_{l_2}) = I_f(p_{l_1-l_2} : p) = I_f(p : p_{l_2-l_1}) = k_f(|l_1 - l_2|).$$

Since the Cauchy standard density is even, we have for a prescribed scale subfamily $I_f(p_{l_1,s} : p_{l_2,s}) = h_f(\chi((l_1, s), (l_2, s))) = k_{f,s}(|l_1 - l_2|)$. By the definition of $\chi$, it follows that we have $k_{f,s}(u) = h_f\left(\frac{u}{2s^2}\right)$.

**Remark 1** *It has been shown that Amari's dual $\pm\alpha$-connections [3] $^\alpha\Gamma$ all coincide with the Levi-Civita metric connection [48]. That is, the $\alpha$-geometry coincides with the Fisher-Rao geometry for the Cauchy family [54], for all $\alpha \in \mathbb{R}$ (see Appendix A). Moreover, Eguchi [19, 20] showed how to build an information-geometric dualistic structure $(M, {}^D g, {}^D\nabla, {}^D\nabla^*)$ from any arbitrary smooth divergence $D$, consisting of a pair of torsion-free affine connections $({}^D\nabla, {}^D\nabla^*)$ coupled to the metric tensor $^D g$ so that we have $^D\nabla^* = {}^{D^*}\nabla$ where $D^*(p : q) := D(q, p)$ denotes the reverse divergence. When the divergence $D$ is a $f$-divergence, it can be shown that the induced connections are $\alpha$-connections, $^D\nabla = {}^\alpha\nabla$ and $^D\nabla^* = {}^{-\alpha}\nabla$ with $\alpha = 3 + 2\frac{f'''(1)}{f''(1)}$, and the metric tensor $^D g = \frac{1}{f''(1)}{}^F g$ is proportional to the Fisher information metric tensor $^F g$ [53]. (Notice that Amari defined* standard *$f$-divergences in [3] by fixing their their scalings so that $f''(1) = 1$.) Since $f$-divergences are symmetric for Cauchy distributions, we have $^{I_f}\nabla = {}^{I_{f^*}}\nabla = {}^{I_f}\nabla^*$.*

**Remark 2** *Of course, not all statistical divergences between Cauchy densities are symmetric. For example, consider the statistical $q$-divergence [3] for a scalar $q \in [1, 3)$:*

$$D_q(p : r) := \frac{1}{(1-q)Z_q(p)}\left(1 - \int p^q(x)r^{1-q}(x)\mathrm{d}\mu(x)\right),$$

*where $Z_q(p) := \int p^q(x)\mathrm{d}\mu(x)$. Then the 2-divergence between two Cauchy densities $p_{\lambda_1}$ and $p_{\lambda_2}$ (with $\lambda_i = (l_i, s_i)$) is available in closed-form (as a corresponding Bregman divergence [3]):*

$$D_2(p_{\lambda_1} : p_{\lambda_2}) = \frac{\pi}{s_2}\|\lambda_1 - \lambda_2\|^2.$$

*Thus $D_2(p_{\lambda_1} : p_{\lambda_2}) \neq D_2(p_{\lambda_2} : p_{\lambda_1})$ when $s_1 \neq s_2$.*

Note that since $I_f(p_{\theta_2} : p_{\theta_1}) = h_f(\chi(\theta_1, \theta_1))$, Lemma 1 can *a posteriori* be checked for the chi-squared divergence: For any $A \in \mathrm{SL}(2, \mathbb{R})$ and $\theta \in \mathbb{H}$, we have

$$\chi(p_{A.\theta_1} : p_{A.\theta_2}) = \chi(p_{\theta_1} : p_{\theta_2}),$$

and therefore for any $f$-divergence, since we have $I_f(p_{A.\theta_1} : p_{A.\theta_2}) = I_f(p_{\theta_1} : p_{\theta_2})$ since

$$I_f(p_{A.\theta_2} : p_{A.\theta_1}) = h_f(\chi(A.\theta_1, A.\theta_1)) = h_f(\chi(\theta_1, \theta_1)) = I_f(p_{\theta_2} : p_{\theta_1}).$$

To prove that $\chi(p_{A.\theta_1} : p_{A.\theta_2}) = \chi(p_{\theta_1} : p_{\theta_2})$, let us first recall that $\text{Im}(A.\theta) = \frac{\text{Im}(\theta)}{|c\theta+d|^2}$ and $|A.\theta_1 - A.\theta_2|^2 = \frac{|\theta_1-\theta_2|^2}{|c\theta_1+d|^2 \, |c\theta_2+d|^2}$. Thus we have

$$
\begin{aligned}
\chi(A.\theta_1, A.\theta_2) &= \frac{|A.\theta_1 - A.\theta_2|^2}{2\,\text{Im}(A.\theta_1)\text{Im}(A.\theta_2)}, \\
&= \frac{|\theta_1 - \theta_2|^2|c\theta_1 + d|^2 \, |c\theta_2 + d|^2}{|c\theta_1 + d|^2 \, |c\theta_2 + d|^2 \, 2\,\text{Im}(\theta_1)\text{Im}(\theta_2)}, \\
&= \frac{|\theta_1 - \theta_2|^2}{2\,\text{Im}(\theta_1)\text{Im}(\theta_2)} = \chi(\theta_1, \theta_2).
\end{aligned}
$$

Alternatively, we may also define a bivariate function $g_f(l, s)$ so that using the action of the location-scale group, we have:

$$
h_f(\chi(\theta_1, \theta_2)) = g_f\left(\frac{l_1 - l_2}{s_2}, \frac{s_1}{s_2}\right),
$$

where $\theta_1 = l_1 + is_1$ and $\theta_2 = l_2 + is_2$. When the function $h_f$ is not explicitly known, we may estimate the $f$-divergences using Monte Carlo importance samplings [55].

## 2.3 Strictly increasing function $h_f$

We have proven that $I_f(p_{\theta_1} : p_{\theta_2}) = I_f(p_{\theta_2} : p_{\theta_1}) = h_f(\chi(\theta_1, \theta_2))$. Let us prove now that the function $h_f$ is is a strictly increasing function.

**Theorem 3** *Let $f : (0, \infty) \to \mathbb{R}$ be a convex function such that $f(1) = 0$ and $f \in C^1((0, 1)) \cap C^1((1, \infty))$ and $f'(x) < f'(y)$ for every $x < 1 < y$. Let $D_f(\lambda : \lambda')$ be the $f$-divergence between $p_\lambda$ and $p_{\lambda'}$, specifically,*

$$
D_f(\lambda : \lambda') = \int_{\mathbb{R}} p_\lambda(x) f\left(\frac{p_{\lambda'}(x)}{p_\lambda(x)}\right) dx.
$$

*Let $\chi$ be McCullagh's maximal invariant. Let $h_f : (0, \infty) \to [0, \infty)$ be the function such that*

$$
h_f(\chi(\lambda, \lambda')) = D_f(\lambda : \lambda'), \ \ \lambda, \lambda' \in \mathbb{H}.
$$

*Then, $h_f$ is a strictly increasing function.*

The assumption of $f$ is complicated as we would like to cover the important case of the TV distance.

Proof.   Let $u \geq 0$. Let $\lambda = i$ and $\lambda' = u + i$. Then, $\chi(i, u + i) = \frac{u^2}{2}$ and hence,

$$
h_f\left(\frac{u^2}{2}\right) = D_f(p_i : p_{u+i}).
$$

Hence it suffices to show that $F_1(u) := D_f(p_i : p_{u+i})$ is a strictly increasing function. We see that

$$
F_1(u) = \int_{\mathbb{R}} \frac{1}{\pi(x^2 + 1)} f\left(\frac{x^2 + 1}{(x - u)^2 + 1}\right) dx.
$$

Then,

9

**Lemma 2**

$$F_1'(u) = \int_{\mathbb{R}} \frac{2(x-u)}{\pi((x-u)^2+1)} f'\left(\frac{x^2+1}{(x-u)^2+1}\right) dx, \; u > 0,$$

*where we let $f'(1) = 0$.*

By the change-of-variable formula,

$$F_1'(u) = \frac{2}{\pi} \int_{\mathbb{R}} \frac{x}{x^2+1} f'\left(\frac{(x+u)^2+1}{x^2+1}\right) dx, \; u > 0,$$

We also see that

$$\int_{\mathbb{R}} \frac{x}{x^2+1} f'\left(\frac{(x+u)^2+1}{x^2+1}\right) dx$$

$$= \int_0^\infty \frac{x}{x^2+1} f'\left(\frac{(x+u)^2+1}{x^2+1}\right) dx + \int_{-\infty}^0 \frac{x}{x^2+1} f'\left(\frac{(x+u)^2+1}{x^2+1}\right) dx$$

$$= \int_0^\infty \frac{x}{x^2+1} \left(f'\left(\frac{(x+u)^2+1}{x^2+1}\right) - f'\left(\frac{(x-u)^2+1}{x^2+1}\right)\right) dx.$$

Since $f$ is convex and $x, u > 0$, it holds that

$$f'\left(\frac{(x+u)^2+1}{x^2+1}\right) \geq f'\left(\frac{(x-u)^2+1}{x^2+1}\right),$$

for every $x > 0$ except $x = u/2$. By the assumption,

$$f'\left(\frac{(x+u)^2+1}{x^2+1}\right) > f'\left(\frac{(x-u)^2+1}{x^2+1}\right), \quad x \in \left(\frac{99}{100}u, \frac{101}{100}u\right).$$

Hence,

$$\int_0^\infty \frac{x}{x^2+1} \left(f'\left(\frac{(x+u)^2+1}{x^2+1}\right) - f'\left(\frac{(x-u)^2+1}{x^2+1}\right)\right) dx > 0.$$

QED.

Proof. [Proof of Lemma 2] We show this assertion for $u = u_0 > 0$.

**Lemma 3** *For every $c > 1$,*

$$R_{f,c} := \sup_{1/c < a < b < c} \left|\frac{f(b)-f(a)}{b-a}\right| < +\infty.$$

Proof. We first remark that

$$\max_{x \in [1/c,c]} |f'(x)| \leq \max\left\{|f'(1/c)|, |f'(c)|\right\},$$

since $f$ is convex.

Assume that $a < b \leq 1$ or $1 \leq a < b$. Then, by the mean-value theorem,

$$\left|\frac{f(b)-f(a)}{b-a}\right| = |f'(\xi)| \leq \max\left\{|f'(1/c)|, |f'(c)|\right\}.$$

10

Finally we assume that $a < 1 < b$. Then,

$$\left| \frac{f(b) - f(a)}{b - a} \right| \le \left| \frac{f(1) - f(a)}{1 - a} \right| + \left| \frac{f(b) - f(1)}{b - 1} \right| \le 2 \max \left\{ |f'(1/c)|, |f'(c)| \right\}.$$

QED.

For each fixed $u > 0$,

$$\frac{1}{1 + u + u^2} \le \frac{x^2 + 1}{(x - u)^2 + 1} \le 1 + u + u^2.$$

Hence, for some $c_0 > 1$,

$$\frac{1}{c_0} < \inf_{x \in \mathbb{R}, u \in \left( \frac{99}{100} u_0, \frac{101}{100} u_0 \right)} \frac{x^2 + 1}{(x - u)^2 + 1} \le \sup_{x \in \mathbb{R}, u \in \left( \frac{99}{100} u_0, \frac{101}{100} u_0 \right)} \frac{x^2 + 1}{(x - u)^2 + 1} < c_0.$$

Assume that $0 < |h| < u_0/100$. Then, by Lemma 3,

$$\frac{1}{x^2 + 1} \left| \frac{1}{h} \left( f \left( \frac{x^2 + 1}{(x - u_0 - h)^2 + 1} \right) - f \left( \frac{x^2 + 1}{(x - u_0)^2 + 1} \right) \right) \right|$$

$$\le \frac{R_{f,c_0}}{x^2 + 1} \left| \frac{1}{h} \left( \frac{x^2 + 1}{(x - u_0 - h)^2 + 1} - \frac{x^2 + 1}{(x - u_0)^2 + 1} \right) \right|$$

$$= R_{f,c_0} \frac{2|x - u_0 - h| + u_0/100}{((x - u_0 - h)^2 + 1)((x - u_0)^2 + 1)} \le R_{f,c_0} \frac{1 + u_0/100}{(x - u_0)^2 + 1}.$$

We see that for $x \ne u_0/2$,

$$\frac{1}{x^2 + 1} \lim_{h \to 0} \frac{1}{h} \left( f \left( \frac{x^2 + 1}{(x - u_0 - h)^2 + 1} \right) - f \left( \frac{x^2 + 1}{(x - u_0)^2 + 1} \right) \right)$$

$$= \frac{2(x - u_0)}{(x - u_0)^2 + 1} f' \left( \frac{x^2 + 1}{(x - u_0)^2 + 1} \right).$$

Now the lemma follows from the dominated convergence theorem. QED.

**Remark 3** *It follows that the Chebyshev center [8] $p_{\lambda^*}$ of a set of $n$ Cauchy distributions $p_{\lambda_1}, \ldots, p_{\lambda_n}$ with respect to any $f$-divergence does not depend on the generator $f$ with $\lambda^* = \arg\min_\lambda \max_i I_f(p_{\lambda_i} : p_\lambda)$ since*

$$\begin{aligned} \arg\min_\lambda \max_i I_f(p_{\lambda_i} : p_\lambda) &= \arg\min_\lambda \max_i h_f(\chi(\lambda_i, \lambda)), \\ &= \arg\min_\lambda \max_i \chi(\lambda_i, \lambda), \\ &= \arg\min_\lambda \max_i \frac{\|\lambda_i - \lambda\|^2}{\lambda_i}. \end{aligned}$$

*Similarly, the Cauchy Voronoi diagrams with respect to $f$-divergences all coincide [54].*

**Remark 4** *It is interesting to consider whether if the symmetry of $f$-divergence between a location-scale family on $\mathbb{R}$ holds for every $f$, then, the family is limited to Cauchy or not. If this is true, then, it implies the characterization of the Cauchy distribution by [34] and [17]. See Proposition 14 in Appendix.*

11

## 2.4 Some illustrating examples

### 2.4.1 The Kullback-Leibler divergence

It was proven in [11] that

$$D_{\mathrm{KL}}(p_{l_1,s_1} : p_{l_2,s_2}) \;=\; \log\left(\frac{(s_1+s_2)^2 + (l_1-l_2)^2}{4s_1 s_2}\right).$$

Thus we have

$$h_{\mathrm{KL}}(u) = \log\left(1 + \frac{1}{2}u\right).$$

This plays an important role in establishing an equivalence criterion for two infinite products of Cauchy measures. See [64].

### 2.4.2 LeCam-Vincze triangular divergence

Let us consider another illustrating example: The LeCam-Vincze triangular divergence [38, 81] defined by

$$D_{\mathrm{LCV}}(p : q) := \int \frac{(p(x) - q(x))^2}{p(x) + q(x)}\,\mathrm{d}x.$$

This divergence is a symmetric $f$-divergence obtained for the generator $f_{\mathrm{LCV}}(u) = \frac{(u-1)^2}{1+u}$. The triangular divergence is a bounded divergence since $f(0) = f^*(0) = 1 < \infty$, and its square root $\sqrt{D_{\mathrm{LCV}}(p : q)}$ yields a metric distance. The LeCam triangular divergence between a Cauchy standard density $p_{0,1}$ and a Cauchy density $p_{l,s}$ is

$$D_{\mathrm{LCV}}(p_{0,1} : p_{l,s}) = 2 - 4\sqrt{\frac{s}{l^2 + s^2 + 2s + 1}} \leq 2.$$

Since $\chi(p_{0,1} : p_{l,s}) = \frac{l^2 + (s-1)^2}{2s}$, we can express the triangular divergence using the $\chi$-squared divergence as

$$D_{\mathrm{LCV}}(p_{l_1,s_1} : p_{l_2,s_2}) = 2 - 4\sqrt{\frac{1}{2(\chi(p_{l_1,s_1}, p_{l_2,s_2}) + 2)}}.$$

Thus we have the function:

$$h_{f_{\mathrm{LCV}}}(u) = 2 - 4\sqrt{\frac{1}{2(u+2)}}.$$

### 2.4.3 Total variation distance

The total variation distance (TVD) is a metric $f$-divergence obtained for the generator $f_{\mathrm{TV}}(u) = \frac{1}{2}|u-1|$:

$$D_{\mathrm{TV}}(p : q) = I_{f_{\mathrm{TV}}}(p : q) = \frac{1}{2}\int_{\mathbb{R}} |p(x) - q(x)|\,\mathrm{d}x.$$

Consider the TVD between two Cauchy densities $p_{l_1,s_1}$ and $p_{l_2,s_2}$: $D_{\mathrm{TV}}(p_{l_1,s_1}, p_{l_2,s_2})$.

- When $s_2 = s_1 = s$, we have one root $r$ for $p_{l_1,s}(x) = p_{l_2,s}(x)$ since the Cauchy standard density $p(x)$ is even: $r = \frac{l_1+l_2}{2}$. Assume without loss of generality that $l_1 < l_2$. Then we have

$$D_{\text{TV}}(p_{l_1,s} : p_{l_2,s}) = \frac{1}{2}\left( \int_{-\infty}^{\frac{l_1+l_2}{2}} (p_{l_1,s}(x) - p_{l_2,s}(x))\mathrm{d}x + \int_{\frac{l_1+l_2}{2}}^{\infty} (p_{l_2,s}(x) - p_{l_1,s}(x))\mathrm{d}x \right),$$

$$= \frac{2}{\pi}\arctan\left( \frac{|l_2 - l_1|}{2s} \right) \le 1.$$

Notice that we have $\lim_{x\to\infty} \arctan(x) = \frac{\pi}{2}$. We can express $D_{\text{TV}}(p_{l_1,s} : p_{l_2,s})$ using $\chi(p_{l_1,s}, p_{l_2,s}) = \frac{(l_2-l_1)^2}{2s^2}$:

$$D_{\text{TV}}(p_{l_1,s} : p_{l_2,s}) = \frac{2}{\pi}\arctan\left( \sqrt{\frac{\chi(p_{l_1,s}, p_{l_2,s})}{2}} \right).$$

See also Appendix E for the total variation between two densities of a location family.

- We calculate the two roots $r_1$ and $r_2$ of $p_{l_1,s_1}(x) = p_{l_2,s_2}(x)$ when $s_2 \ne s_1$:

$$r_1 = \frac{\sqrt{s_1\,s_2{}^3 - 2s_1{}^2\,s_2{}^2 + \left(s_1{}^3 + \left(l_2{}^2 - 2l_1\,l_2 + l_1{}^2\right)s_1\right)s_2} + l_1 s_2 - l_2 s_1}{s_2 - s_1},$$

$$r_2 = \frac{\sqrt{s_1\,s_2{}^3 - 2s_1{}^2\,s_2{}^2 + \left(s_1{}^3 + \left(l_2{}^2 - 2l_1\,l_2 + l_1{}^2\right)s_1\right)s_2} - l_1 s_2 + l_2 s_1}{s_2 - s_1}.$$

Then we use the formula for the definite integral:

$$I(l, s, a, b) := \int_a^b p_{l,s}(x)\mathrm{d}x = \frac{1}{\pi}\left( \arctan\left( \frac{l-a}{s} \right) - \arctan\left( \frac{l-b}{s} \right) \right),$$

where $\arctan(-x) = -\arctan(x)$.

It follows that we have

$$D_{\text{TV}}(p_{l_1,s_1} : p_{l_2,s_2}) =$$
$$\frac{1}{\pi}\left( \arctan\left( \frac{l_2-r_1}{s_2} \right) - \arctan\left( \frac{l_2-r_2}{s_2} \right) + \arctan\left( \frac{l_1-r_1}{s_1} \right) - \arctan\left( \frac{l_1-r_2}{s_2} \right) \right).$$

Rearranging and simplifying the terms, we get:

$$D_{\text{TV}}(p_{l_1,s_1} : p_{l_2,s_1}) = \frac{2}{\pi}\arctan\left( \sqrt{\frac{\chi(p_{l_1,s_1} : p_{l_2,s_1})}{2}} \right),$$

$$= h_{f_{\text{TV}}}\left( \chi[p_{l_1,s_1}, p_{l_2,s_1}] \right),$$

with

$$h_{f_{\text{TV}}}(u) = \frac{2}{\pi}\arctan\left( \sqrt{\frac{u}{2}} \right).$$

### 2.4.4 $f$-divergences for polynomial generators

First, let us consider the $f$-divergence between two Cauchy densities for $f$ a (convex) monomial.

**Proposition 3** *Let $a \geq 2$ be an integer. Let $J_a$ be a function such that*

$$J_a(\chi(z, w)) = \int_{\mathbb{R}} p_z(x)^a p_w(x)^{1-a} \mathrm{d}x, \ z, w \in \mathbb{H}.$$

*Then, $J_a$ is a polynomial with degree $a - 1$.*

Proof. Let $\lambda \in (0, 1)$. Then,

$$J_a\left(\frac{(1 - \lambda)^2}{2\lambda}\right) = \frac{1}{\pi} \frac{1}{\lambda^{a-1}} \int_{\mathbb{R}} \frac{(x^2 + \lambda^2)^{a-1}}{(x^2 + 1)^a} \mathrm{d}x.$$

Hence it suffices to show that the right hand side is a polynomial of $\lambda + \lambda^{-1}$.

Let

$$R(a, i) := \int_{\mathbb{R}} \frac{x^{2i}}{(x^2 + 1)^a} dx, \ 0 \leq i \leq a - 1.$$

Then, by the change-of-variable that $x = 1/y$,

$$R(a, i) = R(a, a - 1 - i), \ 0 \leq i \leq a - 1.$$

By this and the binomial expansion,

$$\frac{1}{\lambda^{a-1}} \int_{\mathbb{R}} \frac{(x^2 + \lambda^2)^{a-1}}{(x^2 + 1)^a} dx = \sum_{i=0}^{a-1} \binom{a-1}{i} R(a, i) \lambda^{a-1-2i}$$

$$= \sum_{i=0}^{a-1} \binom{a-1}{i} R(a, i) \frac{\lambda^{a-1-2i} + \lambda^{2i-a+1}}{2}.$$

By induction in $n$, it is easy to see that $\lambda^n + \lambda^{-n}$ is a polynomial of $\lambda + \lambda^{-1}$ with degree $n$. QED.

It holds that

$$
\begin{aligned}
J_2(t) &= t + 1, \\
J_3(t) &= (3(t + 1)^2 - 1)/2 = \frac{3}{2}t^2 + 3t + 1, \\
J_4(t) &= (5(t + 1)^3 - 3(t + 1))/2 = \frac{5}{2}t^3 + \frac{15}{2}t^2 + 6t + 1, \\
J_5(t) &= (35(t + 1)^4 - 30(t + 1)^2 + 3)/8 = \frac{35}{8}t^4 + \frac{35}{2}t^3 + \frac{45}{2}t^2 + 10t + 1.
\end{aligned}
$$

Notice that the smallest degree coefficient $a_0$ of polynomial $J_d(t) = \sum_{i=0}^{d-1} a_i t^i$ is always one since when $\chi(z, w) = 0$, we have $z = w$ and therefore $J_d(0) = \int_{\mathbb{R}} p_z(x)^a p_w(x)^{1-a} \mathrm{d}x = \int_{\mathbb{R}} p_z(x)^a p_z(x)^{1-a} \mathrm{d}x = \int_{\mathbb{R}} p_z(x) \mathrm{d}x = 1 = a_0$.

The result extends for $f$-divergences between two Cauchy densities for $f(u) = P_d(u) = \sum_{i=1}^{d} a_i u^i - \sum_{i=1}^{d} a_i$ a convex polynomial in degree $d$ with $P_d(1) = 0$. Notice that the set of convex polynomials of degree $d$ can be characterized by the set of positive polynomials [44] of degree $d - 2$ since $P_d(u)$ is convex iff $P_d''(u) \geq 0$. A positive polynomial can always be decomposed as a sum of two squared polynomials [69, 7].

**Proposition 4** *The $f$-divergence between two Cauchy densities for a convex polynomial generator $P_d(u)$ of degree $d$ can be expressed as a $d-1$ dimensional polynomial $Q_{d-1}$ of the chi-squared divergence: $I_{P_d}(p_{\lambda_1}, p_{\lambda_2}) = Q_{d-1}(\chi(p_{\lambda_1}, p_{\lambda_2}))$.*

The proof follows from the fact that $I_{P_d}(p_{\lambda_1}, p_{\lambda_2}) = \sum_{i=0}^{d} I_{f_{a_i}}(p_{\lambda_1}, p_{\lambda_2})$ where $f_{a_i}(u) = a_i^u - a_i$ and Proposition 3. Notice that $\int_{\mathbb{R}} p_z(x)^a p_w(x)^{1-a} \mathrm{d}x = \int_{\mathbb{R}} p_z(x)^{1-a} p_w(x)^a \mathrm{d}x$ since $J_a(\chi(z,w)) = J_a(\chi(w,z))$.

**Remark 5** *In practice, we can estimate the coefficients of $J_d(t) = \sum_{i=0}^{d-1} a_i t^i$ using polynomial regression [21] as follows: Let $a = [a_0, \ldots, a_{d-1}]^\top$ denote the vector of polynomial coefficients of $J_d$. Let us draw $n \geq d$ random variates $\lambda_1, \ldots, \lambda_n$ and $\lambda'_1, \ldots, \lambda'_n$. Define the $n \times d$ matrix $M = [m_{ij}]$ with $m_{ij} = \chi(\lambda_i, \lambda'_i)^{j-1}$. Let $b = [b_1, \ldots, b_n]$ denote the vector with $b_i \simeq \int_{\mathbb{R}} p_{\lambda_i}(x)^d p_{\lambda'_i}(x)^{1-d} \mathrm{d}x$ is numerically approximated (e.g., using a quadrature integration rule or by stochastic Monte Carlo integration). Then we estimate $a$ by $\hat{a} = M^+ b$ where $M^+ := (M^\top M)^{-1} M^\top$ is the pseudo-inverse matrix (with $M^+ = M^{-1}$ when $n = d$). Notice that knowing that $\hat{a}_0$ should be close to one, allows to check the quality of the polynomial regression. In fact, we know that all coefficients $a_i$'s should be rational.*

*For example, we find that*

$$
\begin{aligned}
\hat{J}_6(t) \;=\; & 7.958522957345747t^5 + 39.020985312326296t^4 + 70.4495468953682t^3 + 52.37619399770375t^2 \\
& + 14.951303338589055t + 1.002873073997123
\end{aligned}
$$

*Running a second time, we find another estimate*

$$
\begin{aligned}
\hat{J}_6(t) \;=\; & 7.8720651949082665t^5 + 39.38158109425294t^4 + 70.00024065301261t^3 + 52.49185790967846t^2 \\
& + 15.004692984586242t + 0.9992248562053161
\end{aligned}
$$

*We can also estimate similarly the order-$k$ chi divergence [58]*

$$
D_{\chi,k}(p:q) = \int \frac{(p(x) - q(x))^k}{q(x)^{k-1}} \mathrm{d}x,
$$

*for even integers $k \geq 2$. The order-$k$ chi divergence $D_{\chi,k}(p:q)$ is an $f$-divergence obtained for the convex generator $f_{\chi,k}(u) = (u-1)^k$. Using the binomial expansion, we have [58]:*

$$
D_{\chi,k}(p:q) = \sum_{i=0}^{k} \binom{k}{i}(-1)^i \int q(x)^{i-k+1} p(x)^{k-i} \mathrm{d}x.
$$

*For example, we find using the polynomial regression for $k = 6$:*

$$
\begin{aligned}
\hat{h}_{f_{\chi,6}}(u) \;=\; & 7.875095431165917u^5 + 13.124758716080692u^4 + 2.4996228229861686u^3 \\
& + 0.0013068731474561446u^2 - 7.94214016995198310^{-4}u + 4.227071186713171610^{-5}.
\end{aligned}
$$

*Another run yields a close estimate:*

$$
\begin{aligned}
\hat{h}_{f_{\chi,6}}(u) \;=\; & 7.884522702454348u^5 + 13.081028848015308u^4 + 2.5649432632612843u^3 \\
& + -0.03537083264961893u^2 + 0.005359945026839341u - 1.90124993816098710^{-4}.
\end{aligned}
$$

*Since the first polynomial coefficient $a_0$ of $h_{f_{\chi,k}}(u)$ should be zero, we can assess the quality of the polynomial regression.*

*A different set of techniques consist in estimating symbolically the univariate functions $h_f$ and $k_f$ using symbolic regression [6, 15].*

### 2.4.5 The Jensen-Shannon divergence

Consider the Jensen-Shannon divergence [41, 25] (JSD) (a special case of Sibson's information radius [75] of order 1 between a 2-point set):

$$
\begin{aligned}
D_{\mathrm{JS}}(p:q) &= \frac{1}{2}\left(D_{\mathrm{KL}}\left(p:\frac{p+q}{2}\right) + D_{\mathrm{KL}}\left(q:\frac{p+q}{2}\right)\right), \\
&= h\left(\frac{p+q}{2}\right) - \frac{h(p)+h(q)}{2},
\end{aligned}
$$

where $h(p) = -\int p(x)\log p(x)\mathrm{d}x$ denotes Shannon entropy. The JSD can be rewritten as

$$
D_{\mathrm{JS}}(p:q) = \frac{1}{2}\left(D_K(p:q) + D_K(q:p)\right),
$$

where the divergence $D_K$ [41] is defined by

$$
D_K(p:q) := \int p(x)\log\frac{2p(x)}{p(x)+q(x)}\mathrm{d}x.
$$

The divergence $D_K$ is an $f$-divergence for the generator $f_K(u) = u\log\frac{2u}{1+u}$ such that the reverse $K$-divergence is $D_K{}^*(p:q) := D_K(q:p) = I_{f_K^*}(p:q)$ with conjugate generator $f_K^*(u) = -\log\frac{1+u}{2}$. Thus the JSD is an $f$-divergence for $f_{\mathrm{JS}}(u) = \frac{u}{2}u\log\frac{2u}{1+u} - \frac{1}{2}\log\frac{1+u}{2}$. Since $f_{\mathrm{JS}}(0) < \infty$, the JSD is upper bounded. It is bounded by $\log 2$ since $D_K(p:q) \le \log 2$.

Using the fact that $f$-divergences between Cauchy distributions are symmetric, we have

$$
D_{\mathrm{JS}}(p_{l_1,s_1}:p_{l_2,s_2}) = D_K(p_{l_1,s_1}:p_{l_2,s_2}) = D_K\left(p:p_{\frac{l_2-l_1}{s_1},\frac{s_2}{s_1}}\right).
$$

To get the JSD between two Cauchy distributions, we need to find a closed-form formula for $D_{\mathrm{JS}}(p:p_{l,s}) = D_K(p:p_{l,s})$. Let us skew the divergence $D_K$ [50] with a parameter $\alpha \in (0,1)$:

$$
D_{K_\alpha}(p:q) := D_{\mathrm{KL}}(p:(1-\alpha)p + \alpha q) = \int p(x)\log\frac{p(x)}{(1-\alpha)p(x)+\alpha q(x)}\mathrm{d}x. \tag{5}
$$

The divergence $D_{K_\alpha}$ is an $f$-divergence for the generator $f_{K_\alpha}(u) := -u\log\left((1-\alpha) + \frac{\alpha}{u}\right)$ [50].

Let $p_1(x) := p_{0,1}(x) = \frac{1}{\pi(x^2+1)}$, $p_2(x) := p_{l,s}(x)$ and $m_w(x) = (1-w)p_1(x) + wp_2(x) := \left(\frac{1-w}{\pi(x^2+1)} + \frac{ws}{\pi((x-l)^2+s^2)}\right)$.

In Proposition 1 of [11] (proven in Appendix A), a closed-form is reported for the following definite integral:

$$
\begin{aligned}
A(a,b,c;d,e,f) &= \int_{-\infty}^{\infty}\frac{\log\left(dx^2+ex+f\right)}{ax^2+bx+c}\mathrm{d}x, \\
&= \frac{2\pi\left(\log\left(2af - be + 2cd + \sqrt{4ac-b^2}\sqrt{4df-e^2}\right) - \log(2a)\right)}{\sqrt{4ac-b^2}}.
\end{aligned}
$$

16

Relying on this closed-form formula, we find after calculations that we have:

$$D_{\mathrm{KL}}(p_1 : m_w) = \log\left(\frac{l^2 + (s+1)^2}{(1-w)(l^2 + s^2 + 1) + 2ws + 2\sqrt{s^2 + s((1-s)^2 + l^2)w(1-w)}}\right).$$

We remark that $(1-w)(l^2 + s^2 + 1) + 2ws \geq 2s > 0$ and $s^2 + s((1-s)^2 + l^2)w(1-w) \geq s^2 > 0$. This is analytic with respect to $w$ on $(0,1)$, because there exists a holomorphic extension of this to an open neighborhood of the closed interval $[0,1]$ in $\mathbb{C}$.

We consider now the general case: Let $p_{l_1,s_1}(x) := \frac{s_1}{\pi((x-l_1)^2 + s_1^2)}$, $p_{l_2,s_2}(x) := \frac{s_2}{\pi((x-l_2)^2 + s_2^2)}$ and consider the mixture:

$$\begin{aligned} m(x) \quad &:= \quad (1-w)p_{l_1,s_1}(x) + wp_{l_2,s_2}(x), \\ &= \quad \left(\frac{(1-w)s_1}{\pi((x-l_1)^2 + s_1^2)} + \frac{ws_2}{\pi((x-l_2)^2 + s_2^2)}\right). \end{aligned}$$

Then we have:

$$D_{\mathrm{KL}}(p_{l_1,s_1} : m) =$$
$$\log\left(\frac{(l_1 - l_2)^2 + (s_1 + s_2)^2}{(1-w)(s_1^2 + s_2^2 + (l_1 - l_2)^2) + 2ws_1s_2 + 2\sqrt{s_1^2 s_2^2 + s_1 s_2((s_1 - s_2)^2 + (l_1 - l_2)^2)w(1-w)}}\right) \tag{6}$$

Let us report one example:

$$D_{K_w}(p : p_{1,1}) = D_{\mathrm{KL}}(p_1 : (1-w)p_1(x) + wp_2(x)) = \log 5 - \log\left(3 - w + 2\sqrt{1 + w - w^2}\right).$$

When $w = \frac{1}{2}$, we get $D_K(p_{l_1,s_1} : p_{l_2,s_2}) = D_{\mathrm{KL}}(p_{l_1,s_1} : m)$, and we get the JSD between Cauchy densities $p_{l_1,s_1}$ and $p_{l_2,s_2}$:

$$\begin{aligned} D_{\mathrm{JS}}(p_{l_1,s_1} : p_{l_2,s_2}) \quad &= \quad \log\left(\frac{2\sqrt{(l_1 - l_2)^2 + (s_1 + s_2)^2}}{\sqrt{(l_1 - l_2)^2 + (s_1 + s_2)^2} + 2\sqrt{s_1 s_2}}\right), \tag{7} \\ &=: \quad h_{\mathrm{JS}}(\chi(p_{l_1,s_1}, p_{l_2,s_2})), \end{aligned}$$

with

$$h_{\mathrm{JS}}(u) = \log\left(\frac{2\sqrt{2+u}}{\sqrt{2+u} + \sqrt{2}}\right),$$

since $\frac{(l_1-l_2)^2 + (s_1+s_2)^2}{2s_1 s_2} - 2 = \frac{(l_1-l_2)^2 + (s_1-s_2)^2}{2s_1 s_2}$.

Since $D_{\mathrm{JS}}(p_{l_1,s_1} : p_{l_2,s_2}) = h\left(\frac{p_{l_1,s_1} + p_{l_2,s_2}}{2}\right) - \frac{h(p_{l_1,s_1}) + h(p_{l_2,s_2})}{2}$ and $h(p_{l_s}) = \log(4\pi s)$ [11], we get a formula for the Shannon entropy of the mixture of two Cauchy densities:

$$h\left(\frac{p_{l_1,s_1} + p_{l_2,s_2}}{2}\right) = D_{\mathrm{JS}}(p_{l_1,s_1} : p_{l_2,s_2}) + \log(4\pi\sqrt{s_1 s_2}). \tag{8}$$

Notice that the JSD between two Gaussian distributions is not analytic [61].

**Remark 6** *Consider a mixture family [3, 60]*

$$\mathcal{M} := \left\{ m_\theta(x) = \sum_{i=1}^{D} \theta_i p_i(x) + \left(1 - \sum_{i=1}^{D} \theta_i\right) p_0(x) \; : \theta_i > 0, \sum_{i=1}^{D} \theta_i < 1 \right\}$$

*where the $p_i(x)$'s are linearly independent component distributions. The KLD between two densities $m_{\theta_1}$ and $m_{\theta_2}$ of $\mathcal{M}$ amount to a Bregman divergence [3, 60] for the Shannon negentropy $F(\theta) := -h(m_\theta)$:*

$$D_{\mathrm{KL}}(m_{\theta_1} : m_{\theta_2}) = B_F(\theta_1 : \theta_2),$$

*where*

$$B_F(\theta_1 : \theta_2) := F(\theta_1) - F(\theta_2) - (\theta_1 - \theta_2)^\top \nabla F(\theta_2).$$

*Since $\frac{1}{2}(m_{\theta_1} + m_{\theta_2}) = m_{\frac{\theta_1 + \theta_2}{2}}$, we have*

$$
\begin{aligned}
D_{\mathrm{JS}}(m_{\theta_1} : m_{\theta_2}) &= \frac{1}{2}\left( D_{\mathrm{KL}}\left(m_{\theta_1} : \frac{1}{2}(m_{\theta_1} + m_{\theta_2})\right) + D_{\mathrm{KL}}\left(m_{\theta_2} : \frac{1}{2}(m_{\theta_1} + m_{\theta_2})\right)\right), \\
&= \frac{1}{2}\left( B_F\left(\theta_1 : \frac{\theta_1 + \theta_2}{2}\right) + B_F\left(\theta_2 : \frac{\theta_1 + \theta_2}{2}\right)\right), \\
&= \frac{F(\theta_1) + F(\theta_2)}{2} - F\left(\frac{\theta_1 + \theta_2}{2}\right) := J_F(\theta_1 : \theta_2).
\end{aligned}
$$

*This last expression is called a Jensen divergence [56] $J_F(\theta_1 : \theta_2)$. In general, the Shannon entropy of a mixture is not available in closed-form. However, we have shown that the Shannon entropy of a mixture of two Cauchy distributions is available in closed form in Eq. 8.*

*For example, consider the family of mixtures of two Cauchy distributions with prescribed parameters $(l_0, s_0) = (0, 1)$ and $(l_1, s_1) = (1, 1)$. Then we have the following generator:*

$$F_{0,1,1,1}(\theta) = -h[(1-\theta)p_{0,1} + \theta p_{1,1}] = \theta \log \frac{2\sqrt{1 + \theta - \theta^2} + \theta + 2}{2\sqrt{1 + \theta - \theta^2} - \theta + 3} + \log \frac{2\sqrt{1 + \theta - \theta^2} - \theta + 3}{20\pi},$$

*and the derivative of $F_{0,1,1,1}(\theta)$ is*

$$\eta(\theta) = F'_{0,1,1,1}(\theta) = \log \frac{2\sqrt{1 + \theta - \theta^2} + \theta + 2}{2\sqrt{1 + \theta - \theta^2} - \theta + 3}.$$

*It follows that the Bregman divergence $B_{F_{0,1,1,1}}(\theta_1 : \theta_2)$ is*

$$B_{F_{0,1,1,1}}(\theta_1 : \theta_2) = D_{\mathrm{KL}}[m_{\theta_1} : m_{\theta_2}] = \theta_1 \log \frac{(2\sqrt{1+\theta_1-\theta_1^2}+\theta_1+2)(2\sqrt{1+\theta_2-\theta_2^2}-\theta_2+3)}{(2\sqrt{1+\theta_1-\theta_1^2}-\theta_1+3)(2\sqrt{1+\theta_2-\theta_2^2}+\theta_2+2)} + \log \frac{2\sqrt{1+\theta_1-\theta_1^2}-\theta_1+3}{2\sqrt{1+\theta_2-\theta_2^2}-\theta_2+2}.$$

*Let us define the skewed $\alpha$-Jensen-Shannon divergence:*

$$D_{\mathrm{JS},\alpha}(p : q) = (1-\alpha)D_{\mathrm{KL}}(p : (1-\alpha)p + \alpha q) + \alpha D_{\mathrm{KL}}(q : (1-\alpha)p + \alpha q). \tag{9}$$

*It is an f-divergence (i.e., $D_{\mathrm{JS},\alpha}(p : q) = I_{f_{\mathrm{JS},\alpha}}(p : q)$) for the convex generator:*

$$f_{\mathrm{JS},\alpha} = -(1-\alpha)\log(\alpha u + (1-\alpha)) - \alpha u \log\left(\frac{1-\alpha}{u} + \alpha\right). \tag{10}$$

18

When $\alpha = \frac{1}{2}$, we have $f_{\mathrm{JS}}(u) = f_{\mathrm{JS},\frac{1}{2}}(u) = -\frac{1}{2}\log\frac{1+u}{2} - \frac{1}{2}u\log\left(\frac{1}{2u} + \frac{1}{2}\right) = \frac{1}{2}u\log\frac{2u}{1+u} - \frac{1}{2}\log\frac{1+u}{2}$. The skewed $\alpha$-Jensen-Shannon divergence can be rewritten as

$$D_{\mathrm{JS},\alpha}(p:q) = h((1-\alpha)p + \alpha q) - ((1-\alpha)h(p) + \alpha h(q)). \tag{11}$$

Thus we have

$$h((1-\alpha)p + \alpha q) = D_{\mathrm{JS},\alpha}(p:q) + ((1-\alpha)h(p) + \alpha h(q)). \tag{12}$$

When $p = p_{l_1,s_1}$ and $q = p_{l_2,s_2}$, using Eq. 6, we get a closed-form for $D_{\mathrm{JS},\alpha}(p_{l_1,s_1} : p_{l_2,s_2})$, and hence we have a closed-form for the differential entropy of a mixture of two components $h((1-\alpha)p_{l_1,s_1} + \alpha p_{l_2,s_2})$. Let $m_\theta := (1-\theta)p_{l_1,s_1} + \theta p_{l_2,s_2}$.

The skewed $\alpha$-Jensen-Shannon divergence between two mixtures $m_{\theta_1}$ and $m_{\theta_2}$ amounts to

$$D_{\mathrm{JS},\alpha}(m_{\theta_1} : m_{\theta_2}) = h((1-\alpha)m_{\theta_1} + \alpha m_{\theta_2}) - ((1-\alpha)h(m_{\theta_1}) + \alpha h(m_{\theta_2})). \tag{13}$$

Since $(1-\alpha)m_{\theta_1} + \alpha m_{\theta_2} = m_{(1-\alpha)\theta_1 + \alpha\theta_2}$, we get a closed-form formula for the skewed $\alpha$-Jensen-Shannon divergence between two Cauchy mixtures with two prescribed component distributions.

Similarly, the KLD between two Cauchy mixtures $m_{\theta_1}$ and $m_{\theta_2}$ is available in closed-form using Eq. 6.

### 2.4.6   The Taneja divergence

The Taneja $T$-divergence [78] (Eq. 14) is a symmetric divergence defined by:

$$D_T(p,q) := \int \frac{p(x) + q(x)}{2} \log \frac{p(x) + q(x)}{2\sqrt{p(x)q(x)}} \mathrm{d}x.$$

The $T$-divergence can be rewritten as $D_T(p:q) = \int A(p(x), q(x)) \log \frac{A(p(x),q(x))}{G(p(x),q(x))} \mathrm{d}x$ where $A(a,b) := \frac{a+b}{2}$ and $G(a,b) := \sqrt{ab}$ are the arithmetic mean and the geometric mean of $a > 0$ and $b > 0$, respectively. (Thus the $T$-divergence is also called the arithmetic-geometric mean divergence in [78, 72].) In [1], Banerjee et al. proved that $\sqrt{\Delta(a,b)}$ with $\Delta(a,b) = \log \frac{A(a,b)}{G(a,b)}$ is a metric distance.

The $T$-divergence is an $f$-divergence for the generator:

$$f_T(u) = \frac{u+1}{2} \log \frac{u+1}{2\sqrt{u}}.$$

We have $D_T(p:q) = I_{f_T}(p:q)$ since $f_T(u)$ is convex $(f_T''(u) = \frac{u^2+1}{4u^2(u+1)})$.

The $T$-divergence satisfies $D_{\mathrm{JS}}(p:q) + D_T(p:q) = \frac{1}{4}D_J(p:q)$, where $D_J(p:q)$ is the Jeffreys divergence:

$$D_J(p:q) = D_{\mathrm{KL}}(p:q) + D_{\mathrm{KL}}(q:p).$$

Thus we have

$$D_T(p:q) = \frac{1}{4}D_J(p:q) - D_{\mathrm{JS}}(p:q).$$

Since the Jeffreys divergence is an $f$-divergence for the generator $f_J(u) = (u-1)\log u$, we get $f_T(u) = \frac{1}{4}f_J(u) - f_{\mathrm{JS}}(u)$ since $I_{f_T}(p,q) = I_{\frac{1}{4}f_J}(p,q) - I_{f_{\mathrm{JS}}}(p,q) = I_{\frac{1}{4}f_J - f_{\mathrm{JS}}}(p,q)$. (More generally, $I_{f_1 - f_2} = I_{f_1}(p:q) - I_{f_2}(p:q)$ is an $f$-divergence when $f_1 - f_2$ is convex and strictly convex at 1.)

It follows the following closed-form formula for the Taneja divergence between Cauchy densities:

$$D_T[p_{l_1,s_1}, p_{l_2,s_2}] = \log\left(\frac{1}{2}\left(1 + \sqrt{\frac{(s_1+s_2)^2 + (l_1-l_2)^2}{4s_1s_2}}\right)\right).$$

We can express the $T$-divergence between Cauchy densities as a function of the chi-squared divergence as follows:

$$h_T(u) = \frac{1}{2}h_{\mathrm{KL}}(u) - h_{\mathrm{JS}}(u) = \log\left(\frac{1 + \sqrt{1 + \frac{u}{2}}}{2}\right).$$

A related divergence to the $T$-divergence is the Kumar-Chhina divergence [37]:

$$D_{\mathrm{KC}}(p,q) = \int \frac{(p(x) + q(x))(p(x) - q(x))^2}{p(x)q(x)} \log \frac{p(x) + q(x)}{2\sqrt{p(x)q(x)}} \mathrm{d}x.$$

It is an $f$-divergence for the generator:

$$f_{\mathrm{KC}}(u) = \frac{(u+1)(u-1)^2}{u} \log \frac{u+1}{2\sqrt{u}},$$

since we $D_{\mathrm{KC}}(p,q) = I_{f_{\mathrm{KC}}}(p,q)$ for the convex generator $f_{\mathrm{KC}}$.

## 2.5 Maximal invariants (proof of Proposition 2)

This subsection gives details of arguments in the final part of [46, Section 1].

Proof.    First, let us show that

**Lemma 4** *For every $(z,w) \in \mathbb{H}^2$, there exist $\lambda \geq 1$ and $A \in \mathrm{SL}(2,\mathbb{R})$ such that $(A.z, A.w) = (\lambda i, i)$.*

Proof.    Since the special orthogonal group $\mathrm{SO}(2,\mathbb{R})$ is the isotropy subgroup of $\mathrm{SL}(2,\mathbb{R})$ for $i$ and the action is transitive, it suffices to show that for every $z \in \mathbb{H}$ there exist $\lambda \geq 1$ and $A \in \mathrm{SO}(2,\mathbb{R})$ such that $\lambda i = A.z$.

Since we have that for every $\lambda > 0$,

$$\begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}.\lambda i = \frac{i}{\lambda},$$

it suffices to show that for every $z \in \mathbb{H}$ there exist $\lambda > 0$ and $A \in \mathrm{SO}(2,\mathbb{R})$ such that $\lambda i = A.z$.

We have that

$$\begin{bmatrix} \cos\theta & -\sin\theta \\ \sin\theta & \cos\theta \end{bmatrix}.z = \frac{\frac{|z|^2-1}{2}\sin 2\theta + \mathrm{Re}(z)\cos 2\theta + i\mathrm{Im}(z)}{|z\sin\theta + \cos\theta|^2},$$

Therefore for some $\theta$, we have

$$\frac{|z|^2 - 1}{2}\sin 2\theta + \mathrm{Re}(z)\cos 2\theta = 0.$$

By this lemma, we have that for some $\lambda, \lambda' \geq 1$ and $A, A' \in \mathrm{SL}(2, \mathbb{R})$,

$$(\lambda i, i) = (A.z, A.w), \ (\lambda' i, i) = (A'.z', A'.w'),$$

We see that

$$\chi(z, w) = \chi(\lambda i, i) = \frac{(\lambda - 1)^2}{4\lambda} = \frac{1}{4}\left(\lambda + \frac{1}{\lambda} - 2\right),$$

and

$$\chi(z', w') = \chi(\lambda' i, i) = \frac{(\lambda' - 1)^2}{4\lambda'} = \frac{1}{4}(\lambda' + \frac{1}{\lambda'} - 2).$$

If $\chi(z', w') = \chi(z, w)$, then, $\lambda = \lambda'$ and hence $(A.z, A.w) = (A'.z', A'.w')$.    QED.

# 3    Invariance of $f$-divergences and $f$-divergences between distributions related to the Cauchy distributions

There are several distributions which are strongly related with the Cauchy distributions. In this section, we shall make use of the invariance properties of $f$-divergences to derive results for the circular Cauchy [31, 68], wrapped Cauchy [32] and log-Cauchy [43] families which are all related to the Cauchy distributions via various transformations either on the parameter space or on the observation space.

First, consider the family of circular Cauchy distributions parameterized by complex parameters $w$ belonging to the unit disk $\mathbb{D} = \{w \in \mathbb{C} \ : \ |w| < 1\}$. A Circular Cauchy distribution (CC) is an angular distribution [68] playing an important role in circular and directional statistics [42] with the following probability density function:

$$p_w^{cc}(\phi) := \frac{1}{2\pi} \frac{1 - |w|^2}{|e^{i\phi} - w|^2} \, \mathrm{d}z, \quad \phi \in [-\pi, \pi),$$

where $z := e^{i\phi} \in \mathbb{C}$. Let $w = \rho e^{i\phi_0}$ be the polar form of $w$. The circular Cauchy density can be rewritten [31] as:

$$p_{\rho, \phi_0}^{cc}(\phi) = \frac{1}{2\pi} \frac{1 - \rho^2}{1 + \rho^2 - 2\rho \cos(\phi - \phi_0)} \, \mathrm{d}\phi, \quad \phi \in [-\pi, \pi).$$

Consider the subgroup of Möbius transformations $\mathrm{SL}_2(\mathbb{C})$ that maps $\mathbb{D}$ onto itself via transformations of the holomorphic automorphism group of the complex unit disk [79, 49] (informally speaking, hyperbolic motions):

$$w \mapsto t_{\phi, a}(w) := e^{i\phi} \frac{w + a}{\bar{a}w + 1}, \quad \phi \in [-\pi, \pi), a \in \mathbb{C}.$$

The following invariance of $f$-divergences with respect to non-degenerate holomorphic mappings $t_{\phi, a}$ of parameters holds:

**Proposition 5** *We have $I_f(p_{w_1}^{cc} : p_{w_2}^{cc}) = I_f(p_{t_{\phi,a}(w_1)}^{cc} : p_{t_{\phi,a}(w_2)}^{cc})$ for all $\phi \in [-\pi, \pi)$ and $a \in \mathbb{C}$.*

This proposition relies on the fact that $I_f(p_{\theta_1} : p_{\theta_2}) = I_f(p_{\eta_1} : p_{\eta_2})$ for any smooth invertible transformations $\eta(\theta)$ (with smooth inverse $\theta(\eta)$). Here, however the distribution parameters are complex numbers.

Next, McCullagh [45] noticed that if $X \sim \text{Cauchy}(\theta)$ then $Y = \frac{1+iX}{1-iX}$ follows CCauchy $\left(\frac{1+i\theta}{1-i\theta}\right)$ with parameter complex $w = \frac{1+i\theta}{1-i\theta}$. Denote the complex parameter reciprocal conversion functions $\theta \leftrightarrow w$ by $w(\theta) = \frac{1+i\theta}{1-i\theta}$ and $\theta(w) = i\frac{1-w}{(1+w)}$. Let us write $w = a + ib$ for $a, b \in \mathbb{R}$.

**Theorem 4 ($f$-divergences between circular Cauchy distributions)** *The $f$-divergence between two circular Cauchy distributions amounts to the $f$-divergence between two corresponding Cauchy distributions: $I_f(p_{w_1}^{\text{cc}} : p_{w_2}^{\text{cc}}) = I_f(p_{\theta(w_1)} : p_{\theta(w_2)})$. It follows that all $f$-divergences between circular Cauchy distributions are symmetric and can be expressed as scalar functions of the chi square divergence.*

This theorem follows from the invariance of $f$-divergences [3, 53] and Theorem 1. That is, let $Y = m(X)$ for $m$ a diffeomorphism between continuous random variables $X$ and $Y$. Denote by $p_X$ and $q_Y$ the probability densities functions with support $\mathcal{X}$. It is a key property of $f$-divergences that $f$-divergences are invariant under diffeomorphic transformations [71, 55]:

$$I_f(p_{X_1} : p_{X_2}) = I_f(q_{Y_1} : q_{Y_2}).$$

This invariance of $f$-divergences further holds for non-deterministic mappings called sufficiency of stochastic kernels [40]. This result is related to the the result obtained for the Kullback-Leibler divergence in [2] (Lemma 5.1). It is worth noting that the circular Cauchy distribution can be interpreted as the exit distribution of a Brownian motion starting at $w \in \mathbb{D}$ when reaching the unit boundary circle, see [45].

Next, consider the wrapped Cauchy distributions (WC) [32] with probability density functions:

$$p_{\mu,\gamma}^{\text{wc}}(\phi) = \sum_{n=-\infty}^{\infty} \frac{\gamma}{\pi\left(\gamma^2 + (\phi - \mu + 2\pi n)^2\right)}, \quad -\pi \leq \phi < \pi,$$

where $\mu \in \mathbb{R}$ denotes the peak position of the unwrapped distribution and $\gamma > 0$ the scale parameter. Let $\eta = \mu + i\gamma$.

The density can be rewritten equivalently as

$$p_{\mu,\gamma}^{\text{wc}}(\phi) = \frac{1}{2\pi}\frac{\sinh(\gamma)}{\cosh(\gamma) - \cos(\phi - \mu)}.$$

Since we have the following identity:

$$p_w^{\text{cc}}(\phi) = p^{\text{wc}}(\phi, \eta(w)), \quad \eta(w) = \frac{w - i}{w + i}$$

it follows the following theorem:

**Theorem 5 ($f$-divergences between wrapped Cauchy distributions)** *The $f$-divergence between two wrapped Cauchy distributions amounts to the $f$-divergence between two corresponding Cauchy distributions: $I_f(p_{\eta_1}^{\text{wc}} : p_{\eta_2}^{\text{wc}}) = I_f(p_{\theta(\eta_1)} : p_{\theta(\eta_2)})$. It follows that the $f$-divergence between wrapped Cauchy distributions is symmetric and can be expressed as a scalar function of the chi square divergence.*

Finally, consider the family $\mathcal{LC}$ of Log-Cauchy (LC) distributions (see [43], p. 443) and [65], p. 329):

$$\mathcal{LC} := \left\{ p^{\text{lc}}_{\mu,\sigma}(y) = \frac{1}{y\pi} \left[ \frac{\sigma}{(\log y - \mu)^2 + \sigma^2} \right], \quad \mu > 0, \sigma > 0 \right\},$$

defined on the positive real support $\mathcal{Y} = \mathbb{R}_{++}$.

If $X \sim \text{Cauchy}(l, s)$ is a random variable following a Cauchy distribution then $Y = \exp(X)$ is a random variable following a log-Cauchy distribution with $\mu = l$ and $\sigma = s$. Reciprocally, if $Y$ follows a log-Cauchy distribution $\text{LogCauchy}(\mu, \sigma)$, then $X = \log(Y)$ follows a Cauchy distribution with $l = \mu$ and $s = \sigma$. In particular, if $Y \sim \text{LogCauchy}(0, 1)$ then $X = \log(Y) \sim \text{Cauchy}(0, 1)$.

We state the symmetric property of $f$-divergences between log-Cauchy distributions:

**Theorem 6** *The $f$-divergences between two Log-Cauchy distributions* $\text{LogCauchy}(\mu_1, \sigma_1)$ *and* $\text{LogCauchy}(\mu_2, \sigma_2)$ *amount to the $f$-divergences between the two corresponding Cauchy distributions:* $I_f(p^{\text{lc}}_{\mu_1,\sigma_1} : p^{\text{lc}}_{\mu_2,\sigma_2}) = I_f(p_{\mu_1,\sigma_1} : p_{\mu_2,\sigma_2})$. *It follows that the $f$-divergences between two Log-Cauchy distributions are symmetric and can be expressed as a scalar function of the chi square divergence.*

Proof.     First, let us recall that the generic relationships between the probability density functions $p_X$ and $q_Y$ with corresponding real-valued random variables satisfying $Y = m(X)$ for a differentiable and invertible function $m$ with $m'(x) \neq 0$ is

$$
\begin{aligned}
p_X(x) &= m'(x) \times q_Y(m(x)) = m'(x) \times q_Y(y), \\
q_Y(y) &= (m^{-1})'(y) \times p_X(m^{-1}(y)) = (m^{-1})'(y) \times p_X(x).
\end{aligned}
$$

Now consider the case $y = m(x) = \exp(x)$ with $m^{-1}(y) = \log(y)$, and $m'(x) = \exp(x)$ and $(m^{-1})'(y) = 1/y$. Let us make a change of variable in the $f$-divergence integral with $y = \exp(x)$ and $dy = \exp(x)dx$. We have $p_{l,s}(x)dx = p^{\text{lc}}_{\mu,\sigma}(y)dy$, with $\frac{dx}{dy} = \frac{1}{y}$ and $\frac{dy}{dx} = e^y$. Let $q_{Y_i} \sim \text{LogCauchy}(\mu_i, \sigma_i)$ and $p_{X_i} \sim \text{Cauchy}(\mu_i, \sigma_i)$ for $i \in \{1, 2\}$. By a change of variable, we have:

$$
\begin{aligned}
I_f(q_{Y_1} : q_{Y_2}) &:= \int_{\mathbb{R}_{++}} q_{Y_1}(y) f\left(\frac{q_{Y_2}(y)}{q_{Y_1}(y)}\right) dy \\
&= \int_{\mathbb{R}_{++}} (m^{-1})'(y) \times p_{X_1}(m^{-1}(y)) f\left(\frac{(m^{-1})'(y) \times p_{X_2}(m^{-1}(y))}{(m^{-1})'(y) \times p_{X_1}(m^{-1}(y))}\right) dy, \\
&= \int_{\mathbb{R}} p_{X_1}(x) f\left(\frac{p_{X_2}(x)}{p_{X_1}(x)}\right) dx, \\
&=: I_f(p_{X_1} : p_{X_2}).
\end{aligned}
$$

Then we use the symmetric property of the $f$-divergences of the Cauchy distributions to deduce the symmetry of the $f$-divergences between log-Cauchy distributions: $I_f(p^{\text{lc}}_{\mu_1,\sigma_1} : p^{\text{lc}}_{\mu_2,\sigma_2}) = I_f(p^{\text{lc}}_{\mu_2,\sigma_2} : p^{\text{lc}}_{\mu_1,\sigma_1})$. It follows that we have $I_f(p^{\text{lc}}_{\mu_1,\sigma_1} : p^{\text{lc}}_{\mu_2,\sigma_2}) = h_f(\chi((\mu_1, \sigma_1), (\mu_2, \sigma_2)))$.     QED.

# 4 Asymmetric Kullback-Leibler divergence between multivariate Cauchy distributions

For a symmetric positive-definite $d \times d$ matrix $P \succ 0$ and a $d$-dimensional location vector $\mu$, the density of a random variable [55] $X_{\mu,P} := PX + \mu$ with $X \sim p(x)$ (standard density) is

$$p_{\mu,P}(x) := |P|^{-1} p(P^{-1}(x - \mu)). \tag{14}$$

A $d$-dimensional location scale family is formed by the set of densities $\{p_{\mu,P}(x) \; : \; P \succ 0, \mu \in \mathbb{R}^d\}$. For example, the set of multivariate normal distributions (MVNs) form a multidimensional location-scale family [55].

The probability density function of a $d$-dimensional Cauchy distribution [70] (MVCs) with parameters $\mu \in \mathbb{R}^d$ and $\Sigma \succ 0$ be a $d \times d$ positive-definite symmetric matrix is defined by:

$$p_{\mu,\Sigma}(x) := \frac{C_d}{(\det \Sigma)^{1/2}} \left(1 + (x - \mu)^\top \Sigma^{-1} (x - \mu)\right)^{-(d+1)/2}, \; x \in \mathbb{R}^d,$$

where $C_d = \frac{\Gamma\left(\frac{d+1}{2}\right)}{\pi^{\frac{d+1}{2}}}$ is a normalizing constant, and $\Gamma(\cdot)$ denotes the gamma function. The MVCs form a multivariate location-scale family with standard density:

$$p(x) := \frac{\Gamma\left(\frac{d+1}{2}\right)}{\pi^{\frac{d+1}{2}}} \left(1 + x^\top x\right)^{-(d+1)/2},$$

where matrix parameter $P = \Sigma^{\frac{1}{2}}$ denotes the symmetric positive-definite square root matrix of $\Sigma \succ 0$.

In this section, we shall prove that the $f$-divergences between any two densities of a multidimensional location-scale family with prescribed scale root matrix $P$ and even standard density (i.e., $p(x) = p(-x)$) is symmetric, and then show that the KLD between bivariate Cauchy distributions is asymmetric in general.

First, let us consider the case $\Sigma = I$: The corresponding set of multivariate Cauchy distributions yields a multivariate location subfamily $\{p_\mu(x) = p_{\mu,I}(x) \; : \; \mu \in \mathbb{R}^d\}$ with standard distribution $p(x) = p_{0,I}(x) = \frac{C_d}{(\det \Sigma)^{1/2}} \left(1 + x^\top x\right)^{-(d+1)/2}$. Since the standard density is even (i.e., $p(x) = p(-x)$), we can extend straightforwardly the result of Proposition 1 using a multidimensional change of variable in the integrals of $f$-divergences:

**Proposition 6** *The $f$-divergences between any two densities of the multivariate location Cauchy family is symmetric: $I_f(p_{\mu_1}, p_{\mu_2}) = I_f(p_{\mu_2}, p_{\mu_1})$.*

Next, we consider the case of MVC location subfamilies with prescribed matrix $\Sigma$ (or equivalently $P$).

**Proposition 7** *The $f$-divergences between any two densities of the multivariate location Cauchy family $\{p_{\mu,\Sigma} \; : \; \mu \in \mathbb{R}^d\}$ with prescribed matrix $\Sigma$ is symmetric: $I_f(p_{\mu_1,\Sigma}, p_{\mu_2,\Sigma}) = I_f(p_{\mu_2,\Sigma}, p_{\mu_1,\Sigma})$.*

Proof.    We shall use the following identities of $f$-divergences arising from the location-scale family group structure [55]:

$$I_f\left(p_{l_1,P_1} : p_{l_2,P_2}\right) = I_f\left(p : p_{P_1^{-1}(l_2-l_1),P_1^{-1}P_2}\right) = I_f\left(p_{P_2^{-1}(l_1-l_2),P_2^{-1}P_1} : p\right).$$

Thus for the MVCs, we have:

$$I_f\left(p_{\mu_1,\Sigma_1} : p_{\mu_2,\Sigma_2}\right) = I_f\left(p : p_{\Sigma_1^{-\frac{1}{2}}(\mu_2-\mu_1),\Sigma_1^{-\frac{1}{2}}\Sigma_2^{\frac{1}{2}}}\right) = I_f\left(p_{\Sigma_2^{-\frac{1}{2}}(\mu_1-\mu_2),\Sigma_2^{-\frac{1}{2}}\Sigma_1^{\frac{1}{2}}} : p\right).$$

It follows that when $\Sigma_1 = \Sigma_2 = \Sigma$, we get:

$$I_f\left(p_{\mu_1,\Sigma} : p_{\mu_2,\Sigma}\right) = I_f\left(p : p_{\Sigma^{-\frac{1}{2}}(\mu_2-\mu_1),I}\right) = I_f\left(p_{\Sigma^{-\frac{1}{2}}(\mu_1-\mu_2),I} : p\right).$$

Recasting the equalities using the multivariate location Cauchy family, we obtain:

$$I_f\left(p_{\mu_1,\Sigma} : p_{\mu_2,\Sigma}\right) = I_f\left(p : p_{\Sigma^{-\frac{1}{2}}(\mu_2-\mu_1)}\right) = I_f\left(p_{\Sigma^{-\frac{1}{2}}(\mu_1-\mu_2)} : p\right).$$

Since we proved in Proposition 6 for the multivariate Cauchy location family that $I_f(p_{\mu_1}, p_{\mu_2}) = I_f(p_{\mu_2}, p_{\mu_1})$ (with $p_\mu(x) := p_{\mu,I}(x)$), it follows that we have:

$$I_f\left(p_{\mu_1,\Sigma} : p_{\mu_2,\Sigma}\right) = I_f\left(p : p_{\Sigma^{-\frac{1}{2}}(\mu_2-\mu_1)}\right) = I_f\left(p_{\Sigma^{-\frac{1}{2}}(\mu_2-\mu_1)} : p\right) = I_f\left(p_{\mu_2,\Sigma} : p_{\mu_1,\Sigma}\right).$$

$$\text{QED.}$$

However, contrary to the family of univariate Cauchy distributions, we have the following result:

**Proposition 8** *There exist two bivariate Cauchy densities $p_{\mu_1,\Sigma_1}$ and $p_{\mu_2,\Sigma_2}$ such that $D_{\mathrm{KL}}\left(p_{\mu_1,\Sigma_1} : p_{\mu_2,\Sigma_2}\right) \neq D_{\mathrm{KL}}\left(p_{\mu_2,\Sigma_2} : p_{\mu_1,\Sigma_1}\right).$*

Proof.    We let $d = 2$. By the change of variable in the integral [55], we have

$$D_{\mathrm{KL}}\left(p_{\mu_1,\Sigma_1} : p_{\mu_2,\Sigma_2}\right) = D_{\mathrm{KL}}\left(p_{0,I_2} : p_{\Sigma_1^{-1/2}(\mu_2-\mu_1),\Sigma_1^{-1/2}\Sigma_2\Sigma_1^{-1/2}}\right),$$

where $I_2$ denotes the unit $2 \times 2$ matrix.

Let

$$\mu_1 = 0, \Sigma_1 = I_2, \ \mu_2 = (0,1)^\top, \Sigma_2 = \begin{bmatrix} n & 0 \\ 0 & \frac{1}{n} \end{bmatrix},$$

where $n$ is a natural number. We will show that $D_{\mathrm{KL}}\left(p_{\mu_1,\Sigma_1} : p_{\mu_2,\Sigma_2}\right) \neq D_{\mathrm{KL}}\left(p_{\mu_2,\Sigma_2} : p_{\mu_1,\Sigma_1}\right)$ for sufficiently large $n$. Then,

$$D_{\mathrm{KL}}\left(p_{\mu_1,\Sigma_1} : p_{\mu_2,\Sigma_2}\right) = \frac{3C_2}{2} \int_{\mathbb{R}^2} \frac{\log(1+x_1^2/n+nx_2^2) - \log(1+x_1^2+x_2^2)}{(1+x_1^2+x_2^2)^{3/2}} \mathrm{d}x_1 \mathrm{d}x_2$$

and

$$\begin{aligned} D_{\mathrm{KL}}\left(p_{\mu_2,\Sigma_2} : p_{\mu_1,\Sigma_1}\right) &= D_{\mathrm{KL}}\left(p_{0,I_2} : p_{-\Sigma_1^{-1/2}\mu_1,\Sigma_1^{-1}}\right), \\ &= \frac{3C_2}{2} \int_{\mathbb{R}^2} \frac{\log(1+x_1^2/n+n(x_2+\sqrt{n})^2) - \log(1+x_1^2+x_2^2)}{(1+x_1^2+x_2^2)^{3/2}} \mathrm{d}x_1 \mathrm{d}x_2. \end{aligned}$$

Hence it suffices to show that

$$\int_{\mathbb{R}^2} \frac{\log(1+x_1^2/n+n(x_2+\sqrt{n})^2) - \log(1+x_1^2/n+nx_2^2)}{(1+x_1^2+x_2^2)^{3/2}} \mathrm{d}x_1 \mathrm{d}x_2 \neq 0$$

25

for some $n$.

We see that $\log(1 + x_1^2/n + n(x_2 + \sqrt{n})^2) > \log(1 + x_1^2/n + nx_2^2)$ if and only if $x_2 > -\sqrt{n}/2$. Since $\{(x_1, x_2) : x_2 > -\sqrt{n}/2\} \to \mathbb{R}^2, n \to \infty$, we see that by Fatou's lemma [33] (p. 93),

$$\lim_{n\to\infty} \int_{x_2 > -\sqrt{n}/2} \frac{\log(1 + x_1^2/n + n(x_2 + \sqrt{n})^2) - \log(1 + x_1^2/n + nx_2^2)}{(1 + x_1^2 + x_2^2)^{3/2}} dx_1 dx_2 = +\infty.$$

Hence it suffices to show that

$$\liminf_{n\to\infty} \int_{x_2 \leq -\sqrt{n}/2} \frac{\log(1 + x_1^2/n + n(x_2 + \sqrt{n})^2) - \log(1 + x_1^2/n + nx_2^2)}{(1 + x_1^2 + x_2^2)^{3/2}} dx_1 dx_2 > -\infty. \quad (15)$$

If $x_2 \leq -\sqrt{n}/2$, then,

$$\log(1 + x_1^2/n + n(x_2 + \sqrt{n})^2) - \log(1 + x_1^2/n + nx_2^2) = \log\left(1 + \frac{n^{3/2}(n^{1/2} + 2x_2)}{1 + x_1^2/n + nx_2^2}\right)$$

$$\geq \log\left(1 + \frac{n^{3/2}(n^{1/2} + 2x_2)}{1 + nx_2^2}\right).$$

Let $f(x) := \frac{n^{1/2} + 2x}{1 + nx^2}$, $x < -\sqrt{n}/2$. Then, $f$ is decreasing on $\left(-\infty, -\frac{\sqrt{n}}{2} - \sqrt{\frac{n^2+4}{4n}}\right]$ and increasing on $\left[-\frac{\sqrt{n}}{2} - \sqrt{\frac{n^2+4}{4n}}, -\frac{\sqrt{n}}{2}\right]$. Since $-\frac{\sqrt{n}}{2} - \sqrt{\frac{n^2+4}{4n}} > -\frac{3}{2}\sqrt{n}$ for $n \geq 2$, it holds that for $n \geq 2$,

$$\int_{x_2 \leq -3\sqrt{n}/2} \frac{\log(1 + x_1^2/n + n(x_2 + \sqrt{n})^2) - \log(1 + x_1^2/n + nx_2^2)}{(1 + x_1^2 + x_2^2)^{3/2}} dx_1 dx_2$$

$$\geq \log\left(\frac{4 + n^2}{4 + 9n^2}\right) \int_{\mathbb{R}^2} \frac{dx_1 dx_2}{(1 + x_1^2 + x_2^2)^{3/2}} \geq -2\pi \log 5. \quad (16)$$

If $x_2 = -\frac{\sqrt{n}}{2} - \sqrt{\frac{n^2+4}{4n}}$, then,

$$\log\left(1 + \frac{n^{3/2}(n^{1/2} + 2x_2)}{1 + nx_2^2}\right) = 2\log 2 - 2\log\left(n + \sqrt{n^2 + 4}\right) \geq -\log(n^2 + 4).$$

Hence,

$$\int_{-3\sqrt{n}/2 \leq x_2 \leq -\sqrt{n}/2} \frac{\log(1 + x_1^2/n + n(x_2 + \sqrt{n})^2) - \log(1 + x_1^2/n + nx_2^2)}{(1 + x_1^2 + x_2^2)^{3/2}} dx_1 dx_2$$

$$\geq -\log(n^2 + 4) \int_{-3\sqrt{n}/2 \leq x_2 \leq -\sqrt{n}/2} \frac{dx_1 dx_2}{(1 + x_1^2 + x_2^2)^{3/2}}$$

$$\geq -\sqrt{n}\log(n^2 + 4) \int_{\mathbb{R}} \frac{dx_1}{(1 + x_1^2 + n^2/2)^{3/2}} = -\frac{4\sqrt{n}\log(n^2 + 4)}{n^2 + 2} \to 0, \ n \to \infty. \quad (17)$$

By Eq. (16) and (17), we have Eq. (15).

QED.

**Remark 7** *By numerical computations, we have that*

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \frac{\log(1 + x^2/100 + 100(y + 10)^2)}{(1 + x^2 + y^2)^{3/2}} dx dy = 57.953$$

*and*

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \frac{\log(1 + x^2/100 + 100y^2)}{(1 + x^2 + y^2)^{3/2}} dx dy = 30.1523.$$

# 5   Taylor series of $f$-divergences

In this section, we aim at rewriting the $f$-divergences as converging infinite series of power chi divergences [58, 57]. The Pearson power chi divergence $D_{\chi,k}^P$ of order $k$ (for any integer $k \in \{2, \ldots, \}$) is a dissimilarity obtained for the generator $f_{\chi,k}^P(u) = (u - 1)^k$ which generalizes the Pearson $\chi_2$-divergence ($k = 2$):

$$
\begin{aligned}
D_{\chi,k}^P(p : q) &= \int p(x) f_{\chi,k}^P\left(\frac{q(x)}{p(x)}\right) d\mu(x), \\
&= \int p(x) \left(\frac{q(x)}{p(x)} - 1\right)^k d\mu(x), \\
&= \int \frac{(q(x) - p(x))^k}{p(x)^{k-1}} d\mu(x).
\end{aligned}
$$

We have $D_{\chi,2}^P(p : q) = D_\chi^P(p : q) := \int \frac{(p(x) - q(x))^2}{p(x)} d\mu(x)$. For even integers $k \geq 4$, the Pearson power chi divergence are non-negative dissimilarities since $f_{\chi,k}^P(u)$ is strictly convex (we have $f_{\chi,k}^{P \prime\prime}(u) = k(k-1)(u-1)^{k-2} \geq 0$). For odd integers $k \geq 3$, the Pearson power chi divergence may be negative. Similarly, we can define the Neyman power chi divergence $D_{\chi,k}^N$ of order $k$:

$$D_{\chi,k}^N(p : q) = D_{\chi,k}(q : p) = \int \frac{(p(x) - q(x))^k}{q(x)^{k-1}} d\mu(x).$$

We have $D_{\chi,2}^N(p : q) = D_\chi^N(p : q) := \int \frac{(p(x) - q(x))^2}{q(x)} d\mu(x)$. When $k$ is even it is a $f$-divergence, otherwise $D_{\chi,k}^N$ may fail the positive-definiteness property of $f$-divergences. We note $D_{\chi,k}(p : q) = D_{\chi,k}^P(p : q)$ below.

We first state a general framework to obtain power chi divergence expansions of $f$-divergences.

**Theorem 7** *Let $X$ be a topological space and $\mu$ be a Borel measure on $X$ with full support. Let $\{p_\theta(x)\}_\theta$ be a family of probability density functions on $(X, \mu)$. Assume that for each $\theta$, $p_\theta(x)$ is positive and continuous with respect to $x$. We also assume that for each $\theta_1$ and $\theta_2$ there exists $C = C(\theta_1, \theta_2)$ such that $p_{\theta_1}(x) \leq C p_{\theta_2}(x)$ for every $x \in X$. Let $f(z) = \sum_{n=1}^{\infty} a_n(z - 1)^n$ be an analytic function ($f \in C^\omega$), and denote by $r_f$ be the convergence radius of $f$. Assume that $r_f \geq 1$. Let $I_f$ be the induced $f$-divergence. Then,*
*(i) If $\frac{p_{\theta_2}(x)}{p_{\theta_1}(x)} < 1 + r_f$ for every $x$, then,*

$$I_f(p_{\theta_1} : p_{\theta_2}) = \sum_{n=2}^{\infty} a_n \int_X \left(\frac{p_{\theta_2}(x)}{p_{\theta_1}(x)} - 1\right)^n p_{\theta_1}(x) d\mu(x) = \sum_{n=2}^{\infty} a_n D_{\chi,n}(p_{\theta_1} : p_{\theta_2}).$$

*(ii) If $\frac{p_{\theta_2}(x)}{p_{\theta_1}(x)} > 1 + r_f$ for some $x$, then, the infinite sum*
$\sum_{n=2}^{\infty} a_n \int_X \left( \frac{p_{\theta_2}(x)}{p_{\theta_1}(x)} - 1 \right)^n p_{\theta_1}(x)\mu(dx)$ *diverges.*

Proof.    (i) By the assumption and $r_f \geq 1$, $\inf_{x \in X} \frac{p_{\theta_2}(x)}{p_{\theta_1}(x)} > 1 - r_f$. Hence, $\sup_{x \in X} \left| \frac{p_{\theta_2}(x)}{p_{\theta_1}(x)} - 1 \right| < r_f$. Thus we have the Taylor series:

$$ f\left( \frac{p_{\theta_2}(x)}{p_{\theta_1}(x)} \right) = \sum_{n=2}^{\infty} a_n \left( \frac{p_{\theta_2}(x)}{p_{\theta_1}(x)} - 1 \right)^n, $$

and the convergence is uniform with respect to $x$. By noting that $p_\theta(x)$ is a probability density function, we have the assertion.

(ii) Since $\frac{p_{\theta_2}(x)}{p_{\theta_1}(x)}$ is continuous with respect to $x$, there exist $\delta_0 > 0$ and an open set $U_0$ such that

$$ \inf_{x \in U_0} \frac{p_{\theta_2}(x)}{p_{\theta_1}(x)} \geq \delta_0 + 1 + r_f \geq \delta_0 + 2. $$

Then,

$$ a_n \int_{\frac{p_{\theta_2}(x)}{p_{\theta_1}(x)} \geq 1} \left( \frac{p_{\theta_2}(x)}{p_{\theta_1}(x)} - 1 \right)^n p_{\theta_1}(x)\mu(dx) \geq a_n(\delta_0 + r_f)^n \int_{U_0} p_{\theta_1}(x)\mu(dx). $$

Since $r_f \geq 1$,

$$ a_n \int_{\frac{p_{\theta_2}(x)}{p_{\theta_1}(x)} < 1} \left| \frac{p_{\theta_2}(x)}{p_{\theta_1}(x)} - 1 \right|^n p_{\theta_1}(x)\mu(dx) \leq a_n \left( 1 - \inf_{x \in \mathbb{R}} \frac{p_{\theta_2}(x)}{p_{\theta_1}(x)} \right)^n \to 0, \ n \to \infty. $$

By the assumptions, $\int_{U_0} p_{\theta_1}(x)\mu(dx) > 0$. Thus we see that

$$ \lim_{n \to \infty} a_n \int_X \left( \frac{p_{\theta_2}(x)}{p_{\theta_1}(x)} - 1 \right)^n p_{\theta_1}(x)\mu(dx) = +\infty. $$

QED.

Now we deal with the particular case of Cauchy distributions. We first remark that for every $(l_1, s_1)$ and $(l_2, s_2)$,

$$ \max_{x \in \mathbb{R} \cup \{\pm\infty\}} \frac{p_{l_2,s_2}(x)}{p_{l_1,s_1}(x)} = \max_{x \in \mathbb{R} \cup \{\pm\infty\}} \frac{p_{l_1,s_1}(x)}{p_{l_2,s_2}(x)}, $$

because there exists $A \in \mathrm{SL}(2,\mathbb{R})$ such that $\theta_1 = A.\theta_2$ and $\theta_2 = A.\theta_1$ where $\theta_j = \ell_j + is_j, \ j = 1, 2$.

We first deal with the case that the convergence radius is 1. We denote the Kullback-Leibler, $\alpha$-divergence, Jensen-Shannon and squared Hellinger divergences by $D_{\mathrm{KL}}$, $I_\alpha$, $D_{\mathrm{JS}}$ and $D_H^2$, respectively.

**Lemma 5** *(i) If $l^2 + (s - 4/5)^2 < 9/16$, then, $\sup_{x \in \mathbb{R}} \frac{p_{0,1}(x)}{p_{l,s}(x)} < 2$, and hence,*

$$D_{\mathrm{KL}}(p_{l,s} : p_{0,1}) = \sum_{n=2}^{\infty} \frac{(-1)^n}{n} D_{\chi,n}(p_{l,s} : p_{0,1}),$$

$$I_\alpha(p_{l,s} : p_{0,1}) = \sum_{n=2}^{\infty} \frac{-4}{1 - \alpha^2} \binom{(1 + \alpha)/2}{n} D_{\chi,n}(p_{l,s} : p_{0,1})$$

$$D_{\mathrm{JS}}(p_{l,s} : p_{0,1}) = \sum_{n=2}^{\infty} \frac{(-1)^n (2^{n-1} - 1)}{n(n - 1)2^{n-1}} D_{\chi,n}(p_{l,s} : p_{0,1}),$$

$$D_H^2(p_{l,s} : p_{0,1}) = \sum_{n=2}^{\infty} \frac{(-1)^n (2n - 3)!!}{2^{n-1} n!} D_{\chi,n}(p_{l,s} : p_{0,1}),$$

*where we used the generalized binomial coefficient for the $\alpha$-divergences.*
*(ii) If $l^2 + (s - 4/5)^2 > 9/16$, then, $\sup_{x \in \mathbb{R}} \frac{p_{0,1}(x)}{p_{l,s}(x)} > 2$, and hence, all of the infinite sums in (i) diverge.*

We now deal with the case that the convergence radius is 2. Let $D_{\mathrm{HM}}(p : q) = \int \frac{2p(x)q(x)}{p(x)+q(x)} \mathrm{d}x$ be the harmonic (mean) divergence [29, 16].

**Lemma 6** *(i) If $l^2 + (s - 5/3)^2 < 16/9$, then, $\sup_{x \in \mathbb{R}} \frac{p_{0,1}(x)}{p_{l,s}(x)} < 3$ and hence,*

$$D_{\mathrm{HM}}(p_{l,s} : p_{0,1}) = \sum_{n=2}^{\infty} \frac{(-1)^{n+1}}{2^n} \int_{\mathbb{R}} \left( \frac{p_{0,1}(x)}{p_{l,s}(x)} - 1 \right)^n p_{l,s}(x) \mathrm{d}x = \sum_{n=2}^{\infty} \frac{(-1)^{n+1}}{2^n} D_{\chi,n}(p_{l,s} : p_{0,1}).$$

*(ii) If $l^2 + (s - 5/3)^2 > 16/9$, then, $\sup_{x \in \mathbb{R}} \frac{p_{0,1}(x)}{p_{l,s}(x)} > 3$ and hence, the infinite sum in (i) diverges.*

Other expansions are available in Table 3 of [57] (e.g., Jeffreys' divergence). We refer to the Appendix H for an implementation of the calculation of $f$-divergences using these series.

We finally consider the total variation distance between the Cauchy distributions. Then, we *cannot* expect power chi expansions.

**Proposition 9** *Let $f(u) := \frac{|u-1|}{2}$. Then, for every $a_1, \cdots, a_n$,*

$$\lim_{(l,s) \to (l_0,s_0)} \frac{I_f(p_{l,s}, p_{l_0,s_0}) - \sum_{j=2}^{n} a_j \int_{\mathbb{R}} \left( \frac{p_{l,s}(x)}{p_{l_0,s_0}(x)} - 1 \right)^j p_{l_0,s_0}(x) dx}{\left| \int_{\mathbb{R}} \left( \frac{p_{l,s}(x)}{p_{l_0,s_0}(x)} - 1 \right)^n p_{l_0,s_0}(x) \mathrm{d}x \right|} = +\infty.$$

Proof.

**Lemma 7**

$$\sup_{x \in \mathbb{R}} \left| \frac{p_{l,s}(x)}{p_{l_0,s_0}(x)} - 1 \right| = O\left( \sqrt{(l - l_0)^2 + (s - s_0)^2} \right), \ (l, s) \to (l_0, s_0).$$

Proof. We see that

$$\frac{p_{l,s}(x)}{p_{l_0,s_0}(x)} - 1 = \frac{s}{s_0} - 1 + \left(\frac{s}{s_0} - 1\right)\left(\frac{(x-l_0)^2 + s_0^2}{(x-l)^2 + s^2} - 1\right) + \frac{(x-l_0)^2 + s_0^2}{(x-l) + s^2} - 1.$$

Since

$$\frac{(x-l_0)^2 + s_0^2}{(x-l)^2 + s^2} - 1 = \frac{2(l-l_0)(x-l) + (l-l_0)^2 + s_0^2 - s^2}{(x-l)^2 + s^2} = O\left(\sqrt{(l-l_0)^2 + (s-s_0)^2}\right),$$

we have the assertion. QED.

By this lemma, we see that

$$\sum_{j=2}^{n} a_j \int_{\mathbb{R}} \left(\frac{p_{l,s}(x)}{p_{l_0,s_0}(x)} - 1\right)^j p_{l_0,s_0}(x)dx = O\left((l-l_0)^2 + (s-s_0)^2\right).$$

On the other hand,

$$I_f(p_{l,s}, p_{l_0,s_0}) = \frac{2}{\pi} \arctan\left(\frac{1}{2}\sqrt{\frac{(l-l_0)^2 + (s-s_0)^2}{ss_0}}\right).$$

Hence,

$$\lim_{(l,s)\to(l_0,s_0)} \frac{I_f(p_{l,s}, p_{l_0,s_0})}{(l-l_0)^2 + (s-s_0)^2} = +\infty.$$

Thus we see that

$$\lim_{(l,s)\to(l_0,s_0)} \frac{I_f(p_{l,s}, p_{l_0,s_0}) - \sum_{j=2}^{n} a_j \int_{\mathbb{R}} \left(\frac{p_{l,s}(x)}{p_{l_0,s_0}(x)} - 1\right)^j p_{l_0,s_0}(x)\mathrm{d}x}{(l-l_0)^2 + (s-s_0)^2} = +\infty.$$

By Lemma 7, we see that for $n \geq 2$,

$$\int_{\mathbb{R}} \left(\frac{p_{l,s}(x)}{p_{l_0,s_0}(x)} - 1\right)^n p_{l_0,s_0}(x)\mathrm{d}x = O\left(\left((l-l_0)^2 + (s-s_0)^2\right)^{n/2}\right), \ (l,s)\to(l_0,s_0).$$

Thus we have the assertion. QED.

**Remark 8** *Consider the exponential family of exponential distributions $\{p_\lambda(x) = \lambda\exp(-\lambda x), \ \lambda \in \mathbb{R}_{++}\}$ defined on the positive half-line support $\mathcal{X} = \mathbb{R}_+$. The criterion $\frac{p_{\theta_2}}{p_{\theta_1}} < 1 + r_f$ is satisfied for $\lambda_1 < \lambda_2 < (1+r_f)\lambda_1$. Moreover the Pearson order-k power chi divergences are available in closed form for integers $k > 1$ since $\lambda_1 < \lambda_2$ by adapting Lemma 3 of [58] (i.e., when $\lambda_1 < \lambda_2$, it is enough to have conic natural parameter spaces instead of affine spaces). Thus we can calculate the KLD between $p_{\lambda_1}$ and $p_{\lambda_2}$ as converging Taylor chi series. In this case, the KLD is also known to be in closed-form as a Bregman divergence for exponential distributions:*

$$D_{\mathrm{KL}}(p_{\lambda_1} : p_{\lambda_2}) = \frac{\lambda_2}{\lambda_1} - \log\frac{\lambda_2}{\lambda_1} - 1.$$

*However, if we choose the exponential family of normal distributions, we cannot bound their density ratio, and therefore the Taylor chi series diverge.*

Notice that even if the series diverge, the $f$-divergences may be finite (e.g., when the ratio of densities fails to be bounded by $1 + r_f$). In that case, we cannot represent $I_f$ by a Taylor series. By truncating the distributions, we may potentially find a validity range where to apply the Taylor expansion.

# 6 Metrization of $f$-divergences between Cauchy densities

Recall that $f$-divergences can always be symmetrized by taking the generator $s(u) = f(u) + uf(1/u)$. Metrizing $f$ divergences consists in finding the largest exponent $\alpha$ such that $I_s^\alpha$ is a metric distance satisfying the triangle inequality [30, 66, 80]. For example, the square root of the Jensen-Shannon divergence [25] yields a metric distance which is moreover Hilbertian [1], i.e., meaning that there is an embedding $\phi(\cdot)$ into a Hilbert space $\mathcal{H}$ such that $D_{JS}(p:q) = \|\phi(p) - \phi(q)\|_{\mathcal{H}}$. That is, $\sqrt{JSD}$ admits of Hilbert embedding.

We will show that the square roots of the Kullback-Leibler divergence and the Bhattacharyya divergence are distances on the upper-half plane in Theorems 8 and 9 below respectively. We also show that the square root of the KLD is isometrically embeddable into a Hilbert space in Theorem 11.

## 6.1 Metrization of the Kullback-Leibler diveregnce

The following is a generalization of Theorem 3 in [54].

**Theorem 8** *Let $0 < \alpha \leq 1$. Then $D_{KL}(p_{\theta_1} : p_{\theta_2})^\alpha$ is a metric on $\mathbb{H}$ if and only if $0 < \alpha \leq 1/2$.*

In the following we also give full details of the proof of Theorem 3 in [54].
Proof. We proceed as in [54] by letting

$$t(u) := \log\left(\frac{1 + \cosh(\sqrt{2}u)}{2}\right), u \geq 0.$$

Let us consider the properties of $F_2(u) := t(u)^\alpha/u$.

$$F_2'(u) = -2\frac{t(u)^{\alpha-1}}{u^2}G(u/\sqrt{2}),$$

where

$$G_2(w) := (2 + e^{2w} + e^{-2w})\log\left(\frac{e^w + e^{-w}}{2}\right) - \alpha w(e^{2w} - e^{-2w}).$$

If we let $x := e^w$, then,

$$G_2(w) = (x + x^{-1})\left((x + x^{-1})\log(\frac{x^2 + 1}{2x}) - \alpha(x - x^{-1})\log x\right).$$

Let

$$H_2(x) := x\left((x + x^{-1})\log(\frac{x^2 + 1}{2x}) - \alpha(x - x^{-1})\log x\right).$$

Then, $H_2(1) = 0$ and

$$H_2'(x) = 4\left(x\log(\frac{x^2 + 1}{2}) - (1 + \alpha)x\log x + \frac{x^3}{x^2 + 1} - \alpha x\right).$$

Let

$$I_2(x) := x\log(\frac{x^2 + 1}{2}) - (1 + \alpha)x\log x + \frac{x^3}{x^2 + 1} - \alpha x.$$

Then, $I_2(1) = 1/2 - \alpha$ and

$$I_2'(x) = \log(\frac{x^2+1}{2}) - (1+\alpha)\log x + \frac{x^2(3x^2+5)}{(x^2+1)^2} - (1+2\alpha).$$

Consider the case that $\alpha > 1/2$. Then, $I_2(x) < 0$ for every $x > 1$ which is sufficiently close to 1. Hence, $G_2(w) < 0$ for every $w > 0$ which is sufficiently close to 0. Hence, $F_2'(u) > 0$ for every $u > 0$ which is sufficiently close to 0. This means that $F_2$ is strictly increasing near the origin.

Hence there exists $u_0 > 0$ such that

$$2t(u_0)^\alpha < t(2u_0)^\alpha.$$

Take $x_0, z_0 \in \mathbb{H}$ such that $\rho_{\mathrm{FR}}(x_0, z_0) = 2u_0$, where $\rho_{\mathrm{FR}}$ is the Fisher metric distance on $\mathbb{H}$. By considering the geodesic between $x_0$ and $z_0$, we can take $y_0 \in \mathbb{H}$ such that $\rho_{\mathrm{FR}}(x_0, y_0) = \rho_{\mathrm{FR}}(y_0, z_0) = u_0$.

Finally we consider the case that $\alpha = 1/2$. Let

$$J_2(x) := (x^2+1)^2 \log(\frac{x^2+1}{2}) - \frac{3}{2}(x^2+1)^2 \log x + x^2(3x^2+5) - 2(x^2+1)^2.$$

Then, $J_2(1) = 0$. If we let $y := x^2$, then,

$$J_2(x) = (y+1)^2 \log\left(\frac{y+1}{2}\right) - \frac{3}{4}(y+1)^2 \log y + (y^2 + y - 2).$$

Let $K_2(y) := J(\sqrt{y})$. Then,

$$
\begin{aligned}
K_2'(y) &= 2(y+1)(\log(\frac{y+1}{2})+1) - \frac{3}{2}(y+1)\log y - \frac{3(y+1)^2}{4y} + (2y+1), \\
&= y + (y+1)\left(2\log(y+1) - \frac{3}{2}\log y + \frac{9}{4} - \frac{3}{4y} - 2\log 2\right).
\end{aligned}
$$

If $y > 1$, then,

$$2\log(y+1) > \frac{3}{2}\log y$$

and

$$\frac{9}{4} - \frac{3}{4y} - 2\log 2 > \frac{3}{2} - 2\log 2 > 0.$$

Then, $J_2(x) > J(1) = 0$ for every $x > 1$. Hence, $I_2(x) > I(1) = 0$ for every $x > 1$. Hence, $G_2(w) > 0$ for every $w > 0$. Hence, $F_2'(u) < 0$ for every $u > 0$. This means that $F_2$ is strictly decreasing on $[0, \infty)$. Thus we proved that $D_{\mathrm{KL}}(p_{\theta_1} : p_{\theta_2})^{1/2}$ gives a distance, hence $D_{\mathrm{KL}}(p_{\theta_1} : p_{\theta_2})^\alpha$ is also a distance for every $\alpha \in (0, 1/2)$. QED.

## 6.2 Metrization of the Bhattacharyya divergence

The Bhattacharyya divergence [5] is defined by

$$D_{\mathrm{Bhat}}(p : q) := -\log\left(\int \sqrt{p(x)q(x)}\mathrm{d}x\right).$$

The term $\int \sqrt{p(x)q(x)}\mathrm{d}x$ is called the Bhattacharyya coefficient. It is easy to see that $D_{\mathrm{Bhat}}(p : q) = 0$ iff $p = q$, and $D_{\mathrm{Bhat}}(p : q) = D_{\mathrm{Bhat}}(q : p)$.

**Theorem 9** $\sqrt{D_{\mathrm{Bhat}}(p_{\theta_1} : p_{\theta_2})}$ *is a distance on* $\mathbb{H}$.

For exponential families, see [54, Proposition 2] and [56]. We cannot apply the method of [54, Proposition 2] in a direct manner. We state the reason in the end of this section. We can also show that $D_{\mathrm{Bhat}}(p_{\theta_1} : p_{\theta_2})^{\alpha}$ is not a metric if $\alpha > 1/2$ in the same manner as in the proof of Theorem 8.

Proof.    We show the triangle inequality.  We follow the idea in the proof of Theorem 3 in [54]. We construct the metric transform $t_{\mathrm{FR}\to\mathrm{Bhat}}$ and show that $t_{\mathrm{FR}\to\mathrm{Bhat}}(s)$ is increasing and $\sqrt{t_{\mathrm{FR}\to\mathrm{Bhat}}(s)}/s$ is decreasing.

Let $\rho_{\mathrm{FR}}$ be the Fisher-Rao distance.  Then, by following the argument in the proof of [54, Theorem 3],

$$\chi(z, w) = F_3(\rho_{\mathrm{FR}}(z, w)),$$

where we let

$$F_3(s) := \cosh(\sqrt{2}s) - 1.$$

Let

$$I_3(z, w) := \int \sqrt{p_z(x) p_w(x)} dx.$$

Then, by the invariance of the $f$-divergences,

$$I_3(A.z, A.w) = I_3(z, w).$$

Hence we have that for some function $J_3$, $J_3(\chi(z, w)) = I_3(z, w)$. Hence,

$$\sqrt{D_{\mathrm{Bhat}}(p_{\theta_1} : p_{\theta_2})} = \sqrt{-\log J_3\left(F_3(\rho_{\mathrm{FR}}(\theta_1, \theta_2))\right)}.$$

We have that

$$t_{\mathrm{FR}\to\mathrm{Bhat}}(s) = -\log J_3(F_3(s)).$$

It holds that for every $a \in (0, 1)$,

$$J\left(\chi(ai, i)\right) = I(ai, i).$$

By the change-of-variable $x = \tan\theta$ in the integral of $I(ai, i)$, it is easy to see that

$$I_3(ai, i) = \frac{2\sqrt{a}\mathbf{K}(1 - a^2)}{\pi},$$

where $\mathbf{K}$ is the elliptic integral of the first kind. It is defined by[1]

$$\mathbf{K}(t) := \int_0^{\pi/2} \frac{1}{\sqrt{1 - t\sin^2\theta}} d\theta, \ 0 \le t < 1.$$

Hence,

$$J_3\left(\frac{(1 - a)^2}{2a}\right) = \frac{2\sqrt{a}\mathbf{K}(1 - a^2)}{\pi}.$$

---

[1]This is a little different from the usual definition. The usual one is $\mathbf{K}(t) = \int_0^{\pi/2} \frac{1}{\sqrt{1 - t^2\sin^2\theta}} d\theta$.

Since

$$F_3(s) = \cosh(\sqrt{2}s) - 1 = \frac{(1 - e^{-\sqrt{2}s})^2}{2e^{-\sqrt{2}s}},$$

we have that

$$J_3(F_3(s)) = \frac{2e^{-s/\sqrt{2}}\mathbf{K}(1 - e^{-2\sqrt{2}s})}{\pi}.$$

Since the above function is decreasing with respect to $s$, $t_{\mathrm{FR}\to\mathrm{Bhat}}(s)$ is increasing.

Furthermore, we have that

$$\frac{\sqrt{t_{\mathrm{FR}\to\mathrm{Bhat}}(s)}}{s} = \sqrt{-\frac{1}{s^2}\log\left(\frac{2e^{-s/\sqrt{2}}\mathbf{K}(1 - e^{-2\sqrt{2}s})}{\pi}\right)}. \tag{18}$$

This function is decreasing with respect to $s$. See Figure 1. We can show this fact by using the results for the complete elliptic integrals. The full proof is somewhat complicated. See Section F. QED.

**Remark 9** *It holds that*

$$\lim_{s\to+0} \frac{\sqrt{t_{\mathrm{FR}\to\mathrm{Bhat}}(s)}}{s} = \frac{1}{8}, \text{ and } \lim_{s\to+\infty} \frac{\sqrt{t_{\mathrm{FR}\to\mathrm{Bhat}}(s)}}{s} = 0.$$

**Remark 10** *The squared Hellinger distance $H^2(p : q) := \frac{1}{2}\int\left(\sqrt{p(x)} - \sqrt{q(x)}\right)^2 \mathrm{d}x$ (an $f$-divergence for $f_{\mathrm{Hellinger}}(u) = \frac{1}{2}(\sqrt{u} - 1)^2$) satisfies that*

$$H^2(p_{\theta_1} : p_{\theta_2}) = 1 - \exp\left(-D_{\mathrm{Bhat}}(p_{\theta_1} : p_{\theta_2})\right) = 1 - J_3(F_3(\rho_{\mathrm{FR}}(\theta_1, \theta_2)))$$

$$= 1 - \frac{2e^{-\rho_{\mathrm{FR}}(\theta_1,\theta_2)/\sqrt{2}}\mathbf{K}(1 - e^{-2\sqrt{2}\rho_{\mathrm{FR}}(\theta_1,\theta_2)})}{\pi}$$

$$= 1 - \frac{2K\left(1 - \left(1 + \chi(\theta_1,\theta_2) + \sqrt{\chi(\theta_1,\theta_2)(2 + \chi(\theta_1,\theta_2))}\right)^{-2}\right)}{\pi\sqrt{1 + \chi(\theta_1,\theta_2) + \sqrt{\chi(\theta_1,\theta_2)(2 + \chi(\theta_1,\theta_2))}}}.$$

*The Hellinger distance $H(p_{\theta_1} : p_{\theta_2})$ is known to be a metric distance. Notice that*

$$h_{f_{\mathrm{Hellinger}}}(u) = 1 - \frac{2\mathbf{K}\left(1 - \left(1 + u + \sqrt{u(2 + u)}\right)^{-2}\right)}{\pi\sqrt{1 + u + \sqrt{u(2 + u)}}}$$

*and we check that $h_{f_{\mathrm{Hellinger}}}(0) = 0$ since $\mathbf{K}(0) = \frac{\pi}{2}$.*

**Remark 11** *More generally, let $\mathrm{BC}_\alpha[p : q] := \int_{\mathbb{R}} p(x)^\alpha q(x)^{1-\alpha}\mathrm{d}x$ denote the $\alpha$-skewed Bhattacharyya coefficient for $\alpha \in \mathbb{R}\backslash\{0, 1\}$ (also called the $\alpha$-Chernoff coefficient [51, 52]). The $\alpha$-skewed Bhattacharyya divergence is defined by*

$$D_{\mathrm{Bhat},\alpha}(p : q) := -\log\mathrm{BC}_\alpha[p : q] = -\log\int_{\mathbb{R}} p(x)^\alpha q(x)^{1-\alpha}\mathrm{d}x.$$

Figure 1: Graph of $\dfrac{\sqrt{t_{\text{FR}\to\text{Bhat}}(s)}}{s}$

Using a computer algebra system[2], we can compute the $\alpha$-skewed Bhattacharyya coefficients for integers $\alpha$ in closed form. For example, we find the following closed-form for the definite integrals:

$$\text{BC}_2[p:p_{l,s}] = \frac{s^2 + l^2 + 1}{2s},$$

$$\text{BC}_3[p:p_{l,s}] = \frac{3s^4 + \left(6l^2 + 2\right)s^2 + 3l^4 + 6l^2 + 3}{8s^2},$$

$$\text{BC}_4[p:p_{l,s}] = \frac{5s^6 + \left(15l^2 + 3\right)s^4 + \left(15l^4 + 18l^2 + 3\right)s^2 + 5l^6 + 15l^4 + 15l^2 + 5}{16s^3}, \quad and,$$

$$\text{BC}_5[p:p_{l,s}] = \frac{35s^8 + \left(140l^2 + 20\right)s^6 + \left(210l^4 + 180l^2 + 18\right)s^4 + \left(140l^6 + 300l^4 + 180l^2 + 20\right)s^2 + 35l^8 + 140l^6 + 210l^4 + 140l^2 + 35}{128s^4}.$$

Furthermore, we give some remarks about the complete elliptic integrals of the first and second kinds.

**Remark 12** *(i) In practice, we can calculate efficiently $\mathbf{K}(t)$ using the arithmetic-geometric mean (AGM):*

$$\mathbf{K}(t) = \frac{\pi}{2\text{AGM}(1, \sqrt{1 - t^2})}$$

*where $\text{AGM}(a,b) = \lim_{n\to\infty} a_n = \lim_{n\to\infty} g_n$ with $a_0 = a$, $g_0 = b$, $a_{n+1} = \frac{a_n + g_n}{2}$ and $g_{n+1} = \sqrt{a_n g_n}$. The mean is called the arithmetic-geometric mean because it falls in-between the geometric mean and the arithmetic mean: $g_n \leq \text{AGM}(a,b) \leq a_n$, where $g_n$ is an increasing sequence and $a_n$ is a decreasing sequence. We see that*

$$\text{AGM}(a,b) = \frac{\pi}{4} \frac{a+b}{\mathbf{K}\left(\frac{a-b}{a+b}\right)}.$$

*One way to show this relation is using the invariance of the Cauchy distribution with respect to the Boole transform which is mentioned in Section A.*

---

[2]`https://maxima.sourceforge.io/`

*(ii) Let $\mathbf{K}$ and $\mathbf{E}$ be the complete elliptic integrals of the first and second kinds respectively. We let[3]*

$$\mathbf{E}(t) := \int_0^{\pi/2} \sqrt{1 - t\sin^2\theta}\, d\theta.$$

*The following expansion by C. F. Gauss in 1818 is well-known:*

$$1 - \frac{\mathbf{E}(x)}{\mathbf{K}(x)} = \frac{x}{2} + \sum_{n\geq 1} 2^{n-1}(a_n - b_n)^2, \quad x \in (0,1),$$

*where $(a_0, b_0) = (1, \sqrt{1-x})$ and $(a_{n+1}, b_{n+1}) = \left(\frac{a_n + b_n}{2}, \sqrt{a_n b_n}\right)$, $n \geq 0$. See [73] for more details.*

*By investigating of the behaviors of $\frac{\sqrt{t_{\mathrm{FR}\to\mathrm{Bhat}}(s)}}{s}$ in Eq. 18, we get some approximation formulae of $1 - \frac{\mathbf{E}(x)}{\mathbf{K}(x)}$. See Lemma 16 below for example. By numerical computations, it holds that*

$$1 - \frac{\mathbf{E}(x)}{\mathbf{K}(x)} = \frac{x}{2} + \frac{x^2}{16} + \frac{x^3}{32} + \frac{41}{2048}x^4 + \frac{59}{4096}x^5 + \frac{727}{65536}x^6 + O(x^7),$$

$$x\left(\frac{3}{2} + 4\frac{\log(2K(x)/\pi)}{\log(1-x)}\right) = \frac{x}{2} + \frac{x^2}{16} + \frac{x^3}{32} + \frac{251}{12288}x^4 + \frac{123}{8192}x^5 + \frac{34781}{2949120}x^6 + O(x^7)$$

*and*

$$\frac{x\left(4 - x - \sqrt{(4-3x)^2 + 4(2-x)(1-x)\log(1-x)}\right)}{4x + 2(x-1)\log(1-x)}$$

$$= \frac{x}{2} + \frac{x^2}{16} + \frac{x^3}{32} + \frac{49}{3072}x^4 + \frac{41}{6144}x^5 + \frac{259}{491520}x^6 + O(x^7).$$

*See also [36, Lemma 6.2]. They are very close to each other if $x > 0$ is close to 0. For just a few of recent results about complete elliptic integrals and its applications, see [36], [82] and the references therein.*

*Table 1 summarizes the symmetric closed-form $f$-divergences $I_f(p_\lambda : p_{\lambda'}) = h_f(\chi[p_\lambda : p_{\lambda'}])$ between two univariate Cauchy densities $p_\lambda$ and $p_{\lambda'}$ that we obtained as a function $h_f$ of the chi-squared divergence $\chi[p_\lambda : p_{\lambda'}] = \frac{\|\lambda - \lambda'\|^2}{2\lambda_2 \lambda_2'}$ (with $h_f(0) = 0$).*

**Remark 13** *The proof of [54, Proposition 2] is not applicable to the proof of Theorem 9 above, because it cannot be a Bregman divergence. See [1].*

## 6.3 The Chernoff information

The Chernoff information [52] between two densities $p_1$ and $p_2$ is defined by:

$$C(p_1 : p_2) := -\log \min_{a\in(0,1)} \int p_1(x)^a p_2(x)^{1-a} dx.$$

The Chernoff information provides an upper bound for the error probabilities of Bayes hypothesis testing [13] (Chapter 11).

---

[3]This is also a little different from the usual definition. The usual one is $E(t) := \int_0^{\pi/2} \sqrt{1 - t^2 \sin^2\theta}\, d\theta$.

| $f$-divergence name | $f(u)$ | $h_f(u)$ for $I_f[p_{\lambda_1} : p_{\lambda_2}] = h_f(\chi[p_{\lambda_1} : p_{\lambda_2}])$ |
|---|---|---|
| Chi squared divergence | $(u-1)^2$ | $u$ |
| Total variation distance | $\frac{1}{2}\|u-1\|$ | $\frac{2}{\pi}\arctan\left(\sqrt{\frac{u}{2}}\right)$ |
| Kullback-Leibler divergence | $-\log u$ | $\log(1+\frac{1}{2}u)$ |
| Jensen-Shannon divergence | $\frac{u}{2}\log\frac{2u}{1+u} - \frac{1}{2}\log\frac{1+u}{2}$ | $\log\left(\frac{2\sqrt{2+u}}{\sqrt{2+u}+\sqrt{2}}\right)$ |
| Taneja $T$-divergence | $\frac{u+1}{2}\log\frac{u+1}{2\sqrt{u}}$ | $\log\left(\frac{1+\sqrt{1+\frac{u}{2}}}{2}\right),$ |
| LeCam-Vincze divergence | $\frac{(u-1)^2}{1+u}$ | $2 - 4\sqrt{\frac{1}{2(u+2)}}$ |
| squared Hellinger divergence | $\frac{1}{2}(\sqrt{u}-1)^2$ | $1 - \dfrac{2K\left(1-\left(1+u+\sqrt{u(2+u)}\right)^{-2}\right)}{\pi\sqrt{1+u+\sqrt{u(2+u)}}}$ |

Table 1: Closed-form $f$-divergences between two univariate Cauchy densities expressed as a function $h_f$ of the chi-squared divergence $\chi[p_\lambda : p_{\lambda'}] = \frac{\|\lambda-\lambda'\|^2}{2\lambda_2\lambda_2'}$. The square root of the KLD, LeCam and squared Hellinger divergences between Cauchy densities yields metric distances.

**Theorem 10** *For the univariate Cauchy location-scale families, the Chernoff information is equal to the Bhattacharyya divergence.*

Proof.    Let
$$\Lambda(a) := \log\int_{\mathbb{R}} p_{\theta_1}(x)^a p_{\theta_2}(x)^{1-a}\mathrm{d}x.$$
This is finite for every $\mathbb{R}$, and is in $C^\infty$ class on $\mathbb{R}$.

We see that for every $a \in \mathbb{R}$,

$$\Lambda'(a) = \frac{\int_{\mathbb{R}} p_{\theta_1}(x)^a p_{\theta_2}(x)^{1-a}\log\frac{p_{\theta_1}(x)}{p_{\theta_2}(x)}\mathrm{d}x}{\int_{\mathbb{R}} p_{\theta_1}(x)^a p_{\theta_2}(x)^{1-a}\mathrm{d}x}.$$

By the symmetry of $f$-divergences,

$$\int_{\mathbb{R}} p_{\theta_1}(x)^a p_{\theta_2}(x)^{1-a}\log\frac{p_{\theta_1}(x)}{p_{\theta_2}(x)}\mathrm{d}x = \int_{\mathbb{R}} p_{\theta_2}(x)^a p_{\theta_1}(x)^{1-a}\log\frac{p_{\theta_2}(x)}{p_{\theta_1}(x)}\mathrm{d}x.$$

Hence, for $a = 1/2$,
$$\int_{\mathbb{R}} p_{\theta_1}(x)^{1/2} p_{\theta_2}(x)^{1/2}\log\frac{p_{\theta_1}(x)}{p_{\theta_2}(x)}dx = 0. \tag{19}$$

Hence, $\Lambda'(1/2) = 0$. By the Cauchy-Schwarz inequality, $\Lambda''(a) \geq 0$. Hence $\Lambda(a)$ takes its minimum at $a = 1/2$.                                                QED.

Thus the Chernoff information between two Cauchy distributions $p_{\lambda_1}$ and $p_{\lambda_2}$ can be computed from the Bhattacharyya coefficient $\mathrm{BC}_\alpha[p_{\lambda_1} : p_{\lambda_2}] := \int \sqrt{p_{\lambda_1}(x)p_{\lambda_2}(x)}\mathrm{d}x$:

$$C(p_{\lambda_1} : p_{\lambda_2}) = -\log\mathrm{BC}_\alpha[p_{\lambda_1} : p_{\lambda_2}].$$

Since the Bhattacharyya coefficient can be recovered from the squared Hellinger divergence:

$$\text{BC}_\alpha[p_{\lambda_1} : p_{\lambda_2}] = 1 - H^2(p_{\lambda_1} : p_{\lambda_2}),$$

we use the closed-form of the squared Hellinger divergence (Eq. 10) to recover the closed-form formula of the Bhattacharyya coefficient. The Bhattacharyya and Chernoff divergences are not $f$-divergences because they are not separable divergences. Nevertheless, by abuse of notation, let us write $h_{\text{Chernoff}}(u) = -\log(1 - h_{\text{Hellinger}}(u))$.

**Remark 14** *We can compute Eq. 19 by using the two formulas 4.386.3 and 4.386.4 in p. 588 of [27]. However such approach is much more tedious than the above proof.*

Additional material is available at `https://franknielsen.github.io/CauchyFdivergences/`

# 7 Geometric properties of the metrizations of $f$-divergences

If a divergence $D$ is given, then we can define an associated Riemannian metric $g_D$ on the parameter space by following Eguchi [19, 20]. (See also Remark 1.) Specifically, by regarding $D$ is a smooth function on $M \times M$ where $M$ is the space of parameters, we let

$$(g_D)_r(X_r, Y_r) := -X_p Y_q D(p, q)|_{p=q=r}, \ r \in M,$$

where $X, Y$ are vector fields on $M$.

It is known that if $D$ is the Kullback-Leibler divergence, then, $g_D$ is the Fisher metric. If $D$ is not the Kullback-Leibler divergence, then, we are not sure whether $g_D$ is the Fisher metric. However, $g_D$ is the Fisher metric for every smooth $f$-divergence between the Cauchy distribution.

**Proposition 10** *Let $D_f$ be the $f$-divergence between the univariate Cauchy densities. Let $F$ be a function such that*
$$D_f(p_{\theta_1} : p_{\theta_2}) = F(\chi(\theta_1, \theta_2)), \ \theta_1, \theta_2 \in \mathbb{H}.$$
*Assume that $F$ is in $C^2([0, \infty))$. Then, the Riemannian metric $g_D$ is $F'(0)\rho$, where $\rho$ is the Poincaré metric on $\mathbb{H}$.*

For the (dual) connections induced by the $f$-divergence, see Remark 1. We remark that $\sqrt{2}\rho_{\text{FR}}$ is identical with the Poincaré distance on $\mathbb{H}$.

**Proposition 11** *Let $\sqrt{D_{\text{Bhat}}}$ and $\sqrt{D_{\text{KL}}}$ be the distances between Cauchy densities. Then, neither $(\mathbb{H}, \sqrt{D_{\text{Bhat}}})$ nor $(\mathbb{H}, \sqrt{D_{\text{KL}}})$ is a geodesic metric space.*

Proof. Recall that $p_z(x) = \frac{\text{Im}(z)}{\pi|x-z|^2}, z \in \mathbb{H}$.

Assume that $(\mathbb{H}, \sqrt{D_{\text{KL}}})$ is a geodesic metric space. Then, for every $A > 0$, there exists a continuous map $\gamma : [0, 1] \to \mathbb{H}$ such that $\gamma(0) = i, \gamma(1) = Ai$, and

$$\sqrt{D_{\text{KL}}(p_i : p_{Ai})} = \sqrt{D_{\text{KL}}(p_i : p_{\gamma(t)})} + \sqrt{D_{\text{KL}}(p_{\gamma(t)} : p_{Ai})}$$

for every $t \in (0, 1)$.

Let $\widetilde{\gamma}(t) := \mathrm{Im}(\gamma(t))i$. Then,

$$\chi(\widetilde{\gamma}(t_1), \widetilde{\gamma}(t_2)) \leq \chi(\gamma(t_1), \gamma(t_2)), \quad t_1, t_2 \in [0, 1].$$

Since $\sqrt{D_{\mathrm{KL}}(p_z : p_w)}$ is increasing as a function of $\chi(z, w)$,

$$D_{\mathrm{KL}}\left(p_{\widetilde{\gamma}(t_1)} : p_{\widetilde{\gamma}(t_2)}\right) \leq D_{\mathrm{KL}}(p_{\gamma(t_1)} : p_{\gamma(t_2)}), \quad t_1, t_2 \in [0, 1].$$

Since $\sqrt{D_{\mathrm{Bhat}}}$ is a distance,

$$\sqrt{D_{\mathrm{KL}}(p_i : p_{Ai})} = \sqrt{D_{\mathrm{KL}}(p_i : p_{\widetilde{\gamma}(t)})} + \sqrt{D_{\mathrm{KL}}(p_{\widetilde{\gamma}(t)} : p_{Ai})}$$

for every $t \in (0, 1)$. Since $\widetilde{\gamma}$ is continuous, we see that

$$\sqrt{D_{\mathrm{KL}}(p_i : p_{Ai})} = \sqrt{D_{\mathrm{KL}}(p_i : p_{Bi})} + \sqrt{D_{\mathrm{KL}}(p_{Bi} : p_{Ai})}, \quad B \in (1, A),$$

by the intermediate value theorem.

Let $a > 0$. Then,

$$\rho_{\mathrm{FR}}(i, a^2 i) = \rho_{\mathrm{FR}}(i, ai) + \rho_{\mathrm{FR}}(ai, a^2 i) = 2\rho_{\mathrm{FR}}(i, ai).$$

Hence,

$$\frac{7}{5}\sqrt{\rho_{\mathrm{FR}}(i, a^2 i)} < \sqrt{\rho_{\mathrm{FR}}(i, ai)} + \sqrt{\rho_{\mathrm{FR}}(ai, a^2 i)}.$$

Since

$$\lim_{\chi(z,w)\to\infty} \frac{D_{\mathrm{KL}}(p_z, p_w)}{\rho_{\mathrm{FR}}(z, w)} = \frac{1}{\sqrt{2}}, \tag{20}$$

we see that

$$\sqrt{D_{\mathrm{KL}}(p_i : p_{a^2 i})} < \sqrt{D_{\mathrm{KL}}(p_i : p_{ai})} + \sqrt{D_{\mathrm{KL}}(p_{ai} : p_{a^2 i})}$$

for sufficiently large $a > 0$. Thus we see that $(\mathbb{H}, \sqrt{D_{\mathrm{KL}}})$ is not a geodesic metric space.

The proof for $\sqrt{D_{\mathrm{Bhat}}}$ goes in the same manner, because

$$\lim_{\chi(z,w)\to\infty} \frac{D_{\mathrm{Bhat}}(p_z : p_w)}{\rho_{\mathrm{FR}}(z, w)} = \frac{1}{\sqrt{2}}. \tag{21}$$

QED.

**Proposition 12** *The metric spaces* $(\mathbb{H}, \sqrt{D_{\mathrm{KL}}})$ *and* $(\mathbb{H}, \sqrt{D_{\mathrm{Bhat}}})$ *are both complete.*

Proof. Assume that $(z_n)_n$ is a Cauchy sequence with respect to $\sqrt{D_{\mathrm{KL}}}$. Since $\sqrt{D_{\mathrm{KL}}(p_z : p_w)}$ is increasing as a function of $\chi(z, w)$, we see that $\chi(z_n, z_m) \to 0, n, m \to \infty$. We see that $\chi(z, w) \leq \delta$ if and only if

$$|w - (\mathrm{Re}(z) + i(1 + \delta)\mathrm{Im}(z))| \leq \sqrt{\delta(\delta + 2)}\mathrm{Im}(z).$$

Hence $(z_n)_n$ is bounded. Let $z$ be an accumulation point of $(z_n)_n$. Then, $z_{k_n} \to z$, $n \to \infty$ with respect to the Euclid distance. Hence, $\chi(z_{k_n}, z) \to 0, n \to \infty$. Hence, $\sqrt{D_{\mathrm{KL}}(p_{z_{k_n}} : p_z)} \to 0, n \to \infty$. Since $(z_n)_n$ is a Cauchy sequence with respect to $\sqrt{D_{\mathrm{KL}}}$, we see that $\sqrt{D_{\mathrm{KL}}(p_{z_n} : p_z)} \to 0, n \to \infty$.

QED.

Now by Hopf-Rinow's theorem (see [67, Theorem 16]) and Propositions 11 and 12,

**Proposition 13** *Let $\sqrt{D_{\mathrm{Bhat}}}$ and $\sqrt{D_{\mathrm{KL}}}$ be the distances between Cauchy densities. Then, neither $\sqrt{D_{\mathrm{Bhat}}}$ or $\sqrt{D_{\mathrm{KL}}}$ between Cauchy densities is a Riemannian distance.*

**Remark 15 (alternative proof of Proposition 13)** *For $A \in SL(2, \mathbb{R})$, let $\varphi_A(\theta) := A.\theta$, $\theta \in \mathbb{H}$. We first remark that every Riemannian distance $d$ on $\mathbb{H}$ which is preserved by every $\varphi_A$ has a form of $c\rho$ for some non-negative constant $c$. This is shown by two classical results in Riemannian geometry. We remark that every $SL(2, \mathbb{R})$ action to $\mathbb{H}$ is smooth and bijective. By the Myers-Steenrod theorem (see [67, Theorem 18]), every $\varphi_A$ is a Riemannian isometry with respect to the Riemannian metric associated with $d$. It is well-known that if a Riemannian metric on $\mathbb{H}$ is a Riemannian isometry for every $\varphi_A$, then, it has a form of $c\rho_{\mathrm{FR}}$ for some constant $c$. This is usually stated in the much more general framework for homogeneous spaces. See [35, Proposition X.3.1 and Theorem XI.8.6] for example. By (20) and (21), neither $\sqrt{D_{\mathrm{KL}}}$ or $\sqrt{D_{\mathrm{Bhat}}}$ has a form of $c\rho_{\mathrm{FR}}$ for some constant $c$.*

We finally consider isometric embedding into a Hilbert space.

**Theorem 11** *The square root of the Kullback-Leibler divergence between Cauchy densities is isometrically embeddable into a Hilbert space.*

Proof. By [74], it suffices to show that

$$\sum_{i,j=1}^{n} c_i c_j D_{\mathrm{KL}}(p_{z_i} : p_{z_j}) \leq 0$$

for every $(c_1, \cdots, c_n)$ such that $\sum_{i=1}^{n} c_i = 0$ and every $z_1, \cdots, z_n \in \Theta$.

Let the hyperboloid model be

$$\mathbb{L} := \{(x, y, z) \in \mathbb{R}^3 : z > 0, x^2 + y^2 - z^2 = -1\}.$$

Let

$$d_{\mathbb{L}}\left((x_1, y_1, z_1), (x_2, y_2, z_2)\right) := \cosh^{-1}\left(z_1 z_2 - x_1 x_2 - y_1 y_2\right), \quad (x_1, y_1, z_1), (x_2, y_2, z_2) \in \mathbb{L}.$$

Let $\phi_1 : \mathbb{L} \to \mathbb{D}$ be the map defined by

$$\phi_1(x, y, z) = \left(\frac{x}{1+z}, \frac{y}{1+z}\right).$$

Let $\phi_2 : \mathbb{D} \to \mathbb{H}$ be the map defined by

$$\phi_2(x, y) = \left(-\frac{2y}{(1-x)^2 + y^2}, \frac{1 - x^2 - y^2}{(1-x)^2 + y^2}\right).$$

Then, $\phi_1$ and $\phi_2$ are both bijective. Hence $\phi_2 \circ \phi_1$ is a bijection between $\mathbb{H}$ and $\mathbb{L}$.

Hence it suffices to show that for $(x_1, y_1, z_1), \cdots, (x_n, y_n, z_n) \in \mathbb{L}$,

$$\sum_{i,j=1}^{n} c_i c_j \log\left(1 + \frac{\chi\left(\phi_2(\phi_1(x_i, y_i, z_i)), \phi_2(\phi_1(x_j, y_j, z_j))\right)}{2}\right) \leq 0.$$

Since
$$\chi(\phi_2(w_1), \phi_2(w_2)) = \frac{2|w_1 - w_2|^2}{(1 - |w_1|^2)(1 - |w_2|^2)}, \quad w_1, w_2 \in \mathbb{D},$$
we see that
$$\chi\left(\phi_2(\phi_1(x_1, y_1, z_1)), \phi_2(\phi_1(x_2, y_2, z_2))\right) = z_1 z_2 - x_1 x_2 - y_1 y_2 - 1$$
$$= \cosh\left(d_{\mathbb{L}}\left((x_1, y_1, z_1), (x_2, y_2, z_2)\right)\right) - 1, \quad (x_1, y_1, z_1), (x_2, y_2, z_2) \in \mathbb{L}.$$

Hence, it suffices to show that for $(x_1, y_1, z_1), \cdots, (x_n, y_n, z_n) \in \mathbb{L}$,
$$\sum_{i,j=1}^{n} c_i c_j \log\left(\frac{1 + \cosh\left(d_{\mathbb{L}}\left((x_i, y_i, z_i), (x_j, y_j, z_j)\right)\right)}{2}\right) \leq 0.$$

Since $2(\cosh(x/2))^2 = 1 + \cosh(x), x \in \mathbb{R}$, it suffices to show that for $(x_1, y_1, z_1), \cdots, (x_n, y_n, z_n) \in \mathbb{L}$,
$$\sum_{i,j=1}^{n} c_i c_j 2 \log\left(\cosh\left(\frac{d_{\mathbb{L}}\left((x_i, y_i, z_i), (x_j, y_j, z_j)\right)}{2}\right)\right) \leq 0.$$

Now we can apply Theorem 7.5 in Faraut-Harzallah [23] in order to show the last inequality. QED.

**Remark 16** *(i) The proof of Theorem 7.5 in Faraut-Harzallah [23] heavily depends on Takahashi's long paper [77] in representation theory. Faraut-Harzallah [22] gave another derivation of Theorem 7.5 in Faraut-Harzallah [23]. However it heavily depends on Helgason's long paper [28] in representation theory. By following the outline of [22], we give an elementary proof of Theorem 11 without using the terminologies of representation theory. See Appendix G.*
*(ii) It is natural to consider whether the square root of the Bhattacharyya divergence $\sqrt{D_{\mathrm{Bhat}}}$ is isometrically embeddable into a Hilbert space. The squared Hellinger distance $H^2$ satisfies that*
$$D_{\mathrm{Bhat}}(p_z : p_w) = -\log\left(1 - H^2(p_z : p_w)\right) = -\log\left(\int_{\mathbb{R}} \sqrt{p_z(x)}\sqrt{p_w(x)}dx\right).$$

*$\sqrt{D_{\mathrm{Bhat}}}$ is isometrically embeddable into a Hilbert space if and only if for every $s > 0$, As a function of $(z, w)$, $\left(\int_{\mathbb{R}} p_z(x)p_w(x)dx\right)^s$ is a positive definite kernel on $\mathbb{H}$. By the definition of the squared Hellinger distance, $\int_{\mathbb{R}} p_z(x)p_w(x)dx$ is positive definite. However, to our knowledge, it is not known whether $\left(\int_{\mathbb{R}} p_z(x)p_w(x)dx\right)^s$ is positive definite or not for $s \neq 1$. For Cauchy densities, we can show that*
$$\int_{\mathbb{R}} p_z(x)p_w(x)dx = \frac{1}{\pi}\int_0^{\pi}\left(\cosh(d(z, w)) + \cos\theta\sinh(d(z, w))\right)^{-1/2}d\theta, \quad z, w \in \mathbb{H},$$
*where $d$ is the Poincaré distance. See Appendix G for more details.*

It is also natural to consider whether $(\mathbb{H}, \sqrt{D_{\mathrm{KL}}})$ or $(\mathbb{H}, \sqrt{D_{\mathrm{Bhat}}})$ is Gromov-hyperbolic.

**Definition 1** *Let $(M, d)$ be a metric space.*
*(i) Let the Gromov product be*
$$(x|y)_z := \frac{d(x, z) + d(y, z) - d(x, y)}{2}, \quad x, y, z \in M.$$

(ii) Let $\delta > 0$. We say that $(M, d)$ is $\delta$-hyperbolic if

$$(x|z)_w \geq \min\{(x|y)_w, (y|z)_w\} - \delta, \quad x, y, z, w \in M.$$

We say that $(M, d)$ is Gromov-hyperbolic if it is $\delta$-hyperbolic for some $\delta > 0$.

It is known that $\mathbb{H}$ equipped with the Poincaré metric is Gromov-hyperbolic. (see Proposition 1.4.3 in [12])

**Theorem 12** *Neither $(\mathbb{H}, \sqrt{D_{\mathrm{KL}}})$ or $(\mathbb{H}, \sqrt{D_{\mathrm{Bhat}}})$ is Gromov-hyperbolic.*

Proof.    By Proposition 1.6 in [12], $(M, d)$ is not Gromov-hyperbolic if and only if

$$\sup_{x,y,z,w \in M} (d(x,y) + d(z,w) - \max\{d(x,z) + d(y,w), d(x,w) + d(y,z)\}) = +\infty.$$

We first consider $(\mathbb{H}, \sqrt{D_{\mathrm{KL}}})$. For $0 < a < b$,

$$\sqrt{D_{\mathrm{KL}}(p_{ai} : p_{bi})} = \sqrt{\log\left(\frac{b}{4a} + \frac{a}{4b} + \frac{1}{2}\right)}.$$

Hence, for $k \geq 1$,

$$\lim_{n\to\infty} \sup_{a>0} \left| \sqrt{D_{\mathrm{KL}}(p_{ai} : p_{an^k i})} - \sqrt{k \log n} \right| = \lim_{n\to\infty} \left| \sqrt{D_{\mathrm{KL}}(p_i : p_{n^k i})} - \sqrt{k \log n} \right| = 0.$$

Hence,

$$\lim_{n\to\infty} \left( \sqrt{D_{\mathrm{KL}}(p_i : p_{n^2 i})} + \sqrt{D_{\mathrm{KL}}(p_{ni} : p_{n^3 i})} \right.$$

$$\left. - \max\{ \sqrt{D_{\mathrm{KL}}(p_i : p_{ni})} + \sqrt{D_{\mathrm{KL}}(p_{n^2 i} : p_{n^3 i})}, \sqrt{D_{\mathrm{KL}}(p_i : p_{n^3 i})} + \sqrt{D_{\mathrm{KL}}(p_{ni} : p_{n^2 i})} \} \right)$$

$$= \lim_{n\to\infty} \sqrt{D_{\mathrm{KL}}(p_i : p_{n^2 i})} + \sqrt{D_{\mathrm{KL}}(p_{ni} : p_{n^3 i})} - \sqrt{D_{\mathrm{KL}}(p_i : p_{n^3 i})} - \sqrt{D_{\mathrm{KL}}(p_{ni} : p_{n^2 i})} = +\infty.$$

We second consider $(\mathbb{H}, \sqrt{D_{\mathrm{Bhat}}})$. For $0 < a < b$,

$$\sqrt{D_{\mathrm{Bhat}}(p_{ai} : p_{bi})} = \sqrt{\frac{1}{2} \log \frac{b}{a} - \log\left( \frac{2}{\pi} \mathbf{K}\left( 1 - \frac{a^2}{b^2} \right) \right)}.$$

By Lemma 18 in Appendix,

$$\log(4m) \leq \mathbf{K}\left( 1 - \frac{1}{m^2} \right) \leq 2 \log(4m), \ m \geq 2.$$

Hence, for $k \geq 1$,

$$\lim_{n\to\infty} \sup_{a>0} \left| \sqrt{D_{\mathrm{Bhat}}(p_{ai} : p_{an^k i})} - \sqrt{\frac{k}{2} \log n} \right| = \lim_{n\to\infty} \left| \sqrt{D_{\mathrm{Bhat}}(p_i : p_{n^k i})} - \sqrt{\frac{k}{2} \log n} \right| = 0.$$

Hence,

$$\lim_{n \to \infty} \left( \sqrt{D_{\text{Bhat}}(p_i : p_{n^2 i})} + \sqrt{D_{\text{Bhat}}(p_{ni} : p_{n^3 i})} - \right.$$

$$\left. \max\{ \sqrt{D_{\text{Bhat}}(p_i : p_{ni})} + \sqrt{D_{\text{Bhat}}(p_{n^2 i} : p_{n^3 i})}, \sqrt{D_{\text{Bhat}}(p_i : p_{n^3 i})} + \sqrt{D_{\text{Bhat}}(p_{ni} : p_{n^2 i})} \} \right) = +\infty.$$

QED.

Now we see that both of the metrics $\sqrt{D_{\text{KL}}}$ and $\sqrt{D_{\text{Bhat}}}$ are locally related with the Poincaré metric, however, in global, they are completely different from the Poincaré metric.

# References

[1] Sreangsu Acharyya, Arindam Banerjee, and Daniel Boley. Bregman divergences and triangle inequality. In *Proceedings of the 2013 SIAM International Conference on Data Mining*, pages 476–484. SIAM, 2013.

[2] Yuichi Akaoka, Kazuki Okamura, and Yoshiki Otobe. Bahadur efficiency of the maximum likelihood estimator and one-step estimator for quasi-arithmetic means of the cauchy distribution. *Annals of the Institute of Statistical Mathematics (to appear)*, 2021.

[3] Shun-ichi Amari. *Information Geometry and Its Applications*. Applied Mathematical Sciences. Springer Japan, 2016.

[4] Glen D. Anderson, Mavina Krishna Vamanamurthy, and Matti Vuorinen. Functional inequalities for hypergeometric functions and complete elliptic integrals. *SIAM journal on mathematical analysis*, 23(2):512–524, 1992.

[5] Anil Bhattacharyya. On a measure of divergence between two statistical populations defined by their probability distributions. *Bull. Calcutta Math. Soc.*, 35:99–109, 1943.

[6] Lynne Billard and Edwin Diday. Symbolic regression analysis. In *Classification, Clustering, and Data Analysis*, pages 281–288. Springer, 2002.

[7] Martin Campos-Pinto, Frédérique Charles, and Bruno Després. Algorithms for positive polynomial approximation. *SIAM Journal on Numerical Analysis*, 57(1):148–172, 2019.

[8] Çagatay Candan. Chebyshev center computation on probability simplex with $\alpha$-divergence measure. *IEEE Signal Processing Letters*, 27:1515–1519, 2020.

[9] Limei Cao, Didong Li, Erchuan Zhang, Zhenning Zhang, and Huafei Sun. A statistical cohomogeneity one metric on the upper plane with constant negative curvature. *Advances in Mathematical Physics*, 2014.

[10] B. C. Carlson and John L. Gustafson. Asymptotic expansion of the first elliptic integral. *SIAM Journal on Mathematical Analysis*, 16:1072–1092, 1985.

[11] Frédéric Chyzak and Frank Nielsen. A closed-form formula for the Kullback-Leibler divergence between Cauchy distributions. *arXiv preprint arXiv:1905.10965*, 2019.

[12] Michel Coornaert, Thomas Delzant, and Athanase Papadopoulos. *Géométrie et théorie des groupes. Les groupes hyperboliques de Gromov. (Geometry and group theory. The hyperbolic groups of Gromov)*, volume 1441. 1990.

[13] Thomas M Cover and Joy A Thomas. *Elements of information theory.* John Wiley & Sons, 2012.

[14] Imre Csiszár. Information-type measures of difference of probability distributions and indirect observation. *Studia Scientiarum Mathematicarum Hungarica*, 2:229–318, 1967.

[15] Fabricio Olivetti de Franca and Maira Zabuscha de Lima. Interaction-transformation symbolic regression with extreme learning machine. *Neurocomputing*, 423:609–619, 2021.

[16] Sever S Dragomir et al. A refinement of Jensen's inequality with applications for $f$-divergence measures. *Taiwanese Journal of Mathematics*, 14(1):153–164, 2010.

[17] Jean-Louis Dunau and Henri Senateur. An elementary proof of the Knight-Meyer characterization of the Cauchy distribution. *Journal of multivariate analysis*, 22(1):74–78, 1987.

[18] Morris L Eaton. *Group invariance applications in statistics.* Institute of Mathematical Statistics Hayward, California, 1989.

[19] Shinto Eguchi. Second order efficiency of minimum contrast estimators in a curved exponential family. *The Annals of Statistics*, pages 793–803, 1983.

[20] Shinto Eguchi. Geometry of minimum contrast. *Hiroshima Mathematical Journal*, 22(3):631–647, 1992.

[21] Jianqing Fan and Irene Gijbels. *Local polynomial modelling and its applications*, volume 66. CRC Press, 1996.

[22] Jacques Faraut and Khélifa Harzallah. Fonctions sphériques de type positif sur les espaces hyperboliques. *C. R. Acad. Sci. Paris Sér. A-B*, 274:A1396–A1398, 1972.

[23] Jacques Faraut and Khélifa Harzallah. Distances hilbertiennes invariantes sur un espace homogène. *Ann. Inst. Fourier (Grenoble)*, 24(3):xiv, 171–217, 1974.

[24] Lado Filipovic and Siegfried Selberherr. A Two-Dimensional Lorentzian Distribution for an Atomic Force Microscopy Simulator. In *Monte Carlo Methods and Applications*, pages 97–104. De Gruyter, 2012.

[25] Bent Fuglede and Flemming Topsoe. Jensen-Shannon divergence and Hilbert space embedding. In *International Symposium onInformation Theory, 2004. ISIT 2004. Proceedings.*, page 31. IEEE, 2004.

[26] Shin-itiro Goto and Ken Umeno. Maps on statistical manifolds exactly reduced from the perron-frobenius equations for solvable chaotic maps. *Journal of Mathematical Physics*, 59(3):032701, 2018.

[27] Izrail Solomonovich Gradshteyn and Iosif Moiseevich Ryzhik. *Table of integrals, series, and products.* Academic press, 2014.

[28] Sigurdur Helgason. A duality for symmetric spaces with applications to group representations. *Advances in Math.*, 5:1–154, 1970.

[29] KC Jain and Amit Srivastava. On symmetric information divergence measures of Csiszar's $f$-divergence class. *Journal of Applied Mathematics, Statistics and Informatics (JAMSI)*, 3(1):85–102, 2007.

[30] Peter Kafka, Ferdinand Österreicher, and István Vincze. On powers of $f$-divergences defining a distance. *Studia Sci. Math. Hungar*, 26(4):415–422, 1991.

[31] Shogo Kato, MC Jones, et al. An extended family of circular distributions related to wrapped Cauchy distributions via Brownian motion. *Bernoulli*, 19(1):154–171, 2013.

[32] John T. Kent and David E. Tyler. Maximum likelihood estimation for the wrapped Cauchy distribution. *Journal of Applied Statistics*, 15(2):247–254, 1988.

[33] Srinivasan Kesavan. *Measure and Integration.* Springer, 2019.

[34] Frank B Knight. A characterization of the Cauchy type. *Proceedings of the American Mathematical Society*, 55(1):130–135, 1976.

[35] Shoshichi Kobayashi and Katsumi Nomizu. *Foundations of differential geometry. Vol. II.* Interscience Tracts in Pure and Applied Mathematics, No. 15 Vol. II. Interscience Publishers John Wiley & Sons, Inc., New York-London-Sydney, 1969.

[36] Satoshi Kosugi, Yoshihisa Morita, and Shoji Yotsutani. Stationary solutions to the one-dimensional Cahn-Hilliad equations: Proof by the complete elliptic integrals. *Discrete and Continuous Dynamical Systems*, 19(4):609–629, 2007.

[37] Pranesh Kumar and S Chhina. A symmetric information divergence measure of the Csiszár's $f$-divergence class and its bounds. *Computers & Mathematics with applications*, 49(4):575–588, 2005.

[38] Lucien Le Cam. *Asymptotic methods in statistical decision theory.* Springer Science & Business Media, 2012.

[39] Gérard Letac. Which functions preserve Cauchy laws? *Proceedings of the American Mathematical Society*, 67(2):277–286, 1977.

[40] Friedrich Liese and Igor Vajda. On divergences and informations in statistics and information theory. *IEEE Transactions on Information Theory*, 52(10):4394–4412, 2006.

[41] Jianhua Lin. Divergence measures based on the Shannon entropy. *IEEE Transactions on Information theory*, 37(1):145–151, 1991.

[42] Kanti V Mardia and Peter E Jupp. *Directional statistics*, volume 494. John Wiley & Sons, 2009.

[43] Albert W. Marshall and Ingram Olkin. *Life Distributions: Structure of Nonparametric, Semi-parametric, and Parametric Families.* Springer, 2007.

[44] Murray Marshall. *Positive polynomials and sums of squares.* American Mathematical Soc., 2008. Volume 146.

[45] Peter McCullagh. Conditional inference and Cauchy models. *Biometrika*, 79(2):247–259, 1992.

[46] Peter McCullagh. On the distribution of the Cauchy maximum-likelihood estimator. *Proceedings of the Royal Society of London. Series A: Mathematical and Physical Sciences*, 440(1909):475–479, 1993.

[47] Peter McCullagh. Möbius transformation and Cauchy parameter estimation. *Annals of statistics*, 24(2):787–808, 1996.

[48] Ann F. S. Mitchell. Statistical manifolds of univariate elliptic distributions. *International Statistical Review/Revue Internationale de Statistique*, pages 1–16, 1988.

[49] Tristan Needham. *Visual complex analysis.* Oxford University Press, 1998.

[50] Frank Nielsen. A family of statistical symmetric divergences based on Jensen's inequality. *arXiv preprint arXiv:1009.4004*, 2010.

[51] Frank Nielsen. Chernoff information of exponential families. *arXiv preprint arXiv:1102.2684*, 2011.

[52] Frank Nielsen. An information-geometric characterization of Chernoff information. *IEEE Signal Processing Letters*, 20(3):269–272, 2013.

[53] Frank Nielsen. An elementary introduction to information geometry. *Entropy*, 22(10):1100, 2020.

[54] Frank Nielsen. On Voronoi diagrams on the information-geometric Cauchy manifolds. *Entropy*, 22(7):713, 2020.

[55] Frank Nielsen. On information projections between multivariate elliptical and location-scale families. Technical report, arXiv, 2021. 2101.03839.

[56] Frank Nielsen and Sylvain Boltz. The Burbea-Rao and Bhattacharyya centroids. *IEEE Transactions on Information Theory*, 57(8):5455–5466, 2011.

[57] Frank Nielsen and Gaëtan Hadjeres. On power chi expansions of $f$-divergences. *arXiv preprint arXiv:1903.05818*, 2019.

[58] Frank Nielsen and Richard Nock. On the chi square and higher-order chi distances for approximating $f$-divergences. *IEEE Signal Processing Letters*, 21(1):10–13, 2013.

[59] Frank Nielsen and Richard Nock. Total Jensen divergences: definition, properties and clustering. In *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 2016–2020. IEEE, 2015.

[60] Frank Nielsen and Richard Nock. On the geometry of mixtures of prescribed distributions. In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 2861–2865. IEEE, 2018.

[61] Frank Nielsen and Ke Sun. Guaranteed bounds on information-theoretic measures of univariate mixtures using piecewise log-sum-exp inequalities. *Entropy*, 18(12):442, 2016.

[62] Richard Nock, Frank Nielsen, and Shun-ichi Amari. On conformal divergences and their population minimizers. *IEEE Transactions on Information Theory*, 62(1):527–538, 2015.

[63] Tomonori Noda. Symplectic structures on statistical manifolds. *Journal of the Australian Mathematical Society*, 90(3):371–384, 2011.

[64] Kazuki Okamura. An equivalence criterion for infinite products of Cauchy measures. *Statistics & Probability Letters*, 163(108797):1–5, 2020.

[65] David J Olive. *Statistical theory and inference*. Springer, 2014.

[66] Ferdinand Österreicher and Igor Vajda. A new class of metric divergences on probability spaces and its applicability in statistics. *Annals of the Institute of Statistical Mathematics*, 55(3):639–653, 2003.

[67] Peter Petersen. *Riemannian geometry*, volume 171 of *Graduate Texts in Mathematics*. Springer, New York, second edition, 2006.

[68] Arthur Pewsey, Markus Neuhäuser, and Graeme D Ruxton. *Circular statistics in R*. Oxford University Press, 2013.

[69] Victoria Powers and Bruce Reznick. Polynomials that are positive on an interval. *Transactions of the American Mathematical Society*, 352(10):4677–4692, 2000.

[70] S James Press. Multivariate stable distributions. *Journal of Multivariate Analysis*, 2(4):444–462, 1972.

[71] Yu Qiao and Nobuaki Minematsu. A study on invariance of $f$-divergence and its application to speech recognition. *IEEE Transactions on Signal Processing*, 58(7):3884–3890, 2010.

[72] Mark D. Reid and Robert C. Williamson. Generalised Pinsker Inequalities. In *Conference on Learning Theory*, 2009.

[73] Eugene Salamin. Computation of $\pi$ using arithmetic-geometric mean. *Mathematics of Computation*, 30:565–570, 1976.

[74] I. J. Schoenberg. Metric spaces and positive definite functions. *Transactions of the American Mathematical Society*, 44(3):522–536, 1938.

[75] Robin Sibson. Information radius. *Zeitschrift für Wahrscheinlichkeitstheorie und verwandte Gebiete*, 14(2):149–160, 1969.

[76] Thomas Simon. On the relative entropy between two cauchy distributions. *personal communication*, 2020.

[77] Reiji Takahashi. Sur les représentations unitaires des groupes de Lorentz généralisés. *Bull. Soc. Math. France*, 91:289–433, 1963.

[78] Inder Jeet Taneja. New developments in generalized information measures. In *Advances in Imaging and Electron Physics*, volume 91, pages 37–135. Elsevier, 1995.

[79] Abraham A Ungar. The holomorphic automorphism group of the complex disk. *Aequationes mathematicae*, 47(2-3):240–254, 1994.

[80] Igor Vajda. On metric divergences of probability measures. *Kybernetika*, 45(6):885–900, 2009.

[81] István Vincze. On the concept and measure of information contained in an observation. In *Contributions to Probability*, pages 207–214. Elsevier, 1981.

[82] Zhen-Hang Yang, Wei-Mao Qian, Yu-Ming Chu, and Wen Zhang. On approximating the arithmetic-geometric mean and complete elliptic integral of the first kind. *Journal of Mathematical Analysis and Applications*, 462(2):1714–1726, 2018.

## A    Information geometry of location-scale families

The Fisher information matrix [48, 54] (FIM) of a location-scale family with continuously differentiable standard density $p(x)$ with full support $\mathbb{R}$ is

$$I(\lambda) = \frac{1}{s^2} \begin{bmatrix} a^2 & c \\ c & b^2 \end{bmatrix},$$

where

$$
\begin{aligned}
a^2 &= E_p\left[\left(\frac{p'(x)}{p(x)}\right)^2\right], \\
b^2 &= E_p\left[\left(1 + x\frac{p'(x)}{p(x)}\right)^2\right], \\
c &= E_p\left[\frac{p'(x)}{p(x)}\left(1 + x\frac{p'(x)}{p(x)}\right)\right].
\end{aligned}
$$

When the standard density is even (i.e., $p(x) = p(-x)$), we get a diagonal Fisher matrix that can reparameterize with

$$\theta(\lambda) = \left(\frac{a}{b}\lambda_1, \lambda_2\right)$$

so that the Fisher matrix with respect to $\theta$ becomes

$$I_\theta(\theta) = \frac{b^2}{\theta_2^2} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}.$$

It follows that the Fisher-Rao geometry is hyperbolic with curvature $\kappa = -\frac{1}{b^2} < 0$, and that the Fisher-Rao distance is

$$\rho_p(\lambda_1, \lambda_2) = b\,\rho_U\left(\left(\frac{a}{b}l_1, s_1\right), \left(\frac{a}{b}l_2, s_2\right)\right)$$

where

$$\rho_U(\theta_1, \theta_2) = \operatorname{arccosh}\left(1 + \chi(\theta_1, \theta_2)\right),$$

where $\operatorname{arccosh}(u) = \log(u + \sqrt{u^2 - 1})$ for $u > 1$.

For the Cauchy family, we have $a^2 = b^2 = \frac{1}{2}$ (curvature $\kappa = -\frac{1}{b^2} = -2$) and the Fisher-Rao distance is

$$\rho_{\mathrm{FR}}(p_{\lambda_1} : p_{\lambda_2}) = \frac{1}{\sqrt{2}} \operatorname{arccosh}(1 + \chi(\lambda_1, \lambda_2)).$$

Notice that if we let $\theta = l + is$ then the metric in the complex upper plane $\mathbb{H}$ is $\frac{|\mathrm{d}\theta|^2}{\operatorname{Im}(\theta)^2}$ where $|x + iy| = \sqrt{x^2 + y^2}$ denotes the complex modulus, and $\theta \in \mathbb{H} := \{x + iy \ : \ x \in \mathbb{R}, y \in \mathbb{R}_{++}\}$.

It has been shown that Amari's dual $\pm\alpha$-connections [3] $^\alpha\Gamma$ all coincide with the Levi-Civita metric connection [48] $\Gamma = {}^g\Gamma$ for the Cauchy family since the Amari-Chentsov's totally symmetric cubic tensor $T$ vanishes (i.e., $T_{ijk} = 0$). That is, the $\alpha$-geometry coincides with the Fisher-Rao geometry for the Cauchy family [54], for all $\alpha \in \mathbb{R}$. The $2^3 = 8$ Christoffel functions defining the Levi-Civita metric connection [48] for the Cauchy family are:

$$\begin{aligned}
\Gamma_{11}^1 &= \Gamma_{22}^1 = \Gamma_{12}^2 = \Gamma_{21}^2 = 0, \\
\Gamma_{12}^1 &= \Gamma_{21}^1 = \Gamma_{22}^2 = -\frac{1}{s}, \\
\Gamma_{11}^2 &= \frac{1}{s}.
\end{aligned}$$

Next, we recall the symplectic manifold construction of Goto and Umeno [26] for the family of Cauchy distributions (see also [63] for additional details): The Fisher information metric tensor (FIm) is

$$g_{l,s} = \frac{\mathrm{d}l^2 + \mathrm{d}s^2}{2s^2}.$$

A vector field $K$ is a Killing vector field when the Lie derivative $\mathcal{L}$ of the metric $g$ with respect to $K$ is zero: $\mathcal{L}_K g = 0$, i.e. the vector field $K$ preserves the metric (the flow induced by Killing vector field $K$ is a continuous isometry). The three Killing vector fields on $TM$ are

$$\begin{aligned}
K_1 &= (l^2 - s^2)\partial_l + 2ls\partial_v, \\
K_2 &= l\partial_l + s\partial_s, \\
K_3 &= \partial_l.
\end{aligned}$$

Consider the almost complex structure $J = \mathrm{d}s \otimes \partial_l - \mathrm{d}l \otimes \partial_s$ and the Levi-Civita connection $\nabla^{\mathrm{LC}}$ induced by the Fisher information metric. Then $(M, g, J, \nabla^{\mathrm{LC}})$ is a symplectic statistical manifold (Definition 4.14 of [26], see also [63]) equipped with the symplectic form $\omega = -\frac{1}{2s^2}\mathrm{d}l \wedge \mathrm{d}s$ with the set of canonical coordinates $(l, \frac{1}{2s})$. We have $\mathcal{L}_{K_1}\omega = \mathcal{L}_{K_2}\omega = \mathcal{L}_{K_3}\omega = 0$.

The information geometry of the wrapped Cauchy family is investigated in [9]. Goto and Umeno [26] regards the Cauchy distribution as an invariant measure of the generalized Boole transforms and they model the Cauchy manifold is modeled as a symplectic statistical manifold. The Boole transform $\frac{1}{2}\left(X - \frac{1}{X}\right)$ of a standard Cauchy random variable $X$ yields a standard Cauchy random variable. See Subsection B.3 below. See [39] for a description of the functions preserving Cauchy distributions.

# B Relationship between the parametric family

We can interpret that the invariance of Cauchy $f$-divergence in Lemma 1 arises from a relationship between the parametric family as in Assumption 1 below rather than the definition of the Cauchy

density itself, although it is shown that they are equivalent to each other by [47, 26]. This measure-theoretic viewpoint is clear and useful. As an application, we can give a simple, alternative proof of [26, Proposition 3.1 and Theorem 3.1].

## B.1 measure-theoretic framework

Let $(X, \mu)$ be a measure space. Let $\varphi : \Theta \cup X \to \Theta \cup X$ be a map such that $\varphi(\Theta) \subset \Theta$ and $\varphi(X) \subset X$. Assume that $\varphi|_X$ is measurable. For $\theta \in \mathbb{H}$, let $P_\theta(dx) := p_\theta(x)\mu(dx)$, where $p_\theta$ is non-negative measurable function on $X$ and $P_\theta(dx)$ is a probability measure on $X$.

**Assumption 1** $P_{\varphi(\theta)} = P_\theta \circ \varphi^{-1}$ for every $\theta$ and $\varphi$.

We consider one-dimensional location-scale families. We assume that $X = \mathbb{R}$, $\Theta = \mathbb{H}$ and $\mu$ is the Lebesgue measure. Let $(U_i)_i$ be at most countable disjoint open sets of $\mathbb{R}$ such that $\mu(\mathbb{R} \setminus (\cup_i U_i)) = 0$ and $\varphi|_{U_i}$ is smooth and injective for each $i$.

**Lemma 8**

$$p_{\varphi(\theta)}(x) = \sum_i \frac{p_\theta(\varphi_i^{-1}(x))}{|\varphi'(\varphi_i^{-1}(x))|} 1_{\varphi(U_i)}(x), \quad \text{a.e. } x.$$

Proof.    By Assumption 1 and the change of variable formula, it holds that for every nonnegative measurable function $f$,

$$\int_{\mathbb{R}} f(x) p_{\varphi(\theta)}(x)dx = \int_{\mathbb{R}} f(\varphi(x)) p_\theta(x)dx = \sum_i \int_{U_i} f(\varphi(x)) p_\theta(x)dx$$

$$= \sum_i \int_{\varphi(U_i)} f(y) \frac{p_\theta(\varphi_i^{-1}(y))}{|\varphi'(\varphi_i^{-1}(y))|} dy. \tag{22}$$

Thus we have the assertion.                                                                                 QED.

**Proposition 14** *Assume that $f : (0, \infty) \to \mathbb{R}$ is smooth and $f(1) = 0$ and convex. Let $D_f(p_{\theta_1} : p_{\theta_2})$ be the $f$-divergence between $p_{\theta_1}$ and $p_{\theta_2}$, that is,*

$$D_f(p_{\theta_1} : p_{\theta_2}) := \int_{\mathbb{R}} f\left(\frac{p_{\theta_1}(x)}{p_{\theta_2}(x)}\right) p_{\theta_1}(x)dx.$$

*Assume that $\{\varphi_i(U_i)\}_i$ are disjoint. Then,*

$$D_f\left(p_{\varphi(\theta_1)} : p_{\varphi(\theta_2)}\right) = D_f(p_{\theta_1} : p_{\theta_2}).$$

The assumption that $\{\varphi_i(U_i)\}_i$ are disjoint is crucial. See Remark 18.
Proof.    By using (22) and the fact that $\varphi$ is bijective except a measure zero set, we see that

$$p_\theta(x) = p_{\varphi(\theta)}(\varphi(x))|\varphi'(x)|, \quad \text{a.e. } x.$$

Hence,

$$D_f\left(p_{\varphi(\theta_1)} : p_{\varphi(\theta_2)}\right) = \int_{\mathbb{R}} f\left(\frac{p_{\varphi(\theta_1)}(\varphi(x))}{p_{\varphi(\theta_2)}(\varphi(x))}\right) p_{\theta_1}(x)dx = \int_{\mathbb{R}} f\left(\frac{p_{\theta_1}(x)}{p_{\theta_2}(x)}\right) p_{\theta_1}(x)dx.$$

QED.

## B.2 Möbius transformations

For $A = \begin{pmatrix} a & b \\ c & d \end{pmatrix} \in SL(2, \mathbb{R})$, let $\varphi_A(z) = A \cdot z := \dfrac{az + b}{cz + d}$. This is well-defined on $\overline{\mathbb{H}}$ if $c = 0$, and on $\overline{\mathbb{H}} \setminus \{-d/c\}$ if $c \neq 0$. If $c \neq 0$, then we let $\varphi_A(-d/c) := a/c$. Then, $\varphi_A$ is a bijection on $\mathbb{R}$. This also holds if $c = 0$.

For $\theta \in \mathbb{H}$, let $P_\theta(dx) := p_\theta(x)dx$. Assume that $p_\theta(x) > 0$ for every $x \in \mathbb{R}$. This is a probability measure on $\mathbb{R}$.

**Lemma 9** $P_{\varphi_A(\theta)} = P_\theta \circ \varphi_A^{-1}$ *for every* $A \in SL(2, \mathbb{R})$ *and* $\theta \in \mathbb{H}$.

By [34], such parametric location-scale family is restricted to the univariate Cauchy distribution. Let $\chi$ be the maximal invariant. By Proposition 14, we have Theorem 1.

**Proposition 15** *Let*

$$C(x; \ell, s) := \frac{s}{\pi} \frac{1}{(x - \ell)^2 + s^2}, \quad x, \ell \in \mathbb{R}, s > 0.$$

*Assume that* $x \neq -d/c$ *if* $c \neq 0$. *Let* $x' := \varphi_A(x)$,

$$\ell' := \mathrm{Re}(\varphi_A(\ell + is)) = \frac{(a\ell + b)(c\ell + d) + acs^2}{(c\ell + d)^2 + c^2 s^2}$$

*and*

$$s' := \mathrm{Im}(\varphi_A(\ell + is)) = \frac{s}{(c\ell + d)^2 + c^2 s^2}.$$

*Then,* $\varphi_A^{-1}(x') = \{x\}$, *and*

$$C(x'; \ell', s') = \frac{C(x; \ell, s)}{|\varphi_A'(x)|}.$$

This assertion essentially corresponds to [26, Proposition 3.1 and Theorem 3.1].

Proof.    We remark that $\varphi_A$ is bijective. Hence $\varphi_A^{-1}(x') = \{x\}$. By Lemma 8,

$$C(x'; \ell', s') = \frac{C(x; \ell, s)}{|\varphi_A'(x)|}, \quad \text{a.e. } x.$$

Since the functions in the left and right hand sides in the above display are both continuous on $\mathbb{R} \setminus \{-d/c\}$, we have the assertion.                                                          QED.

**Remark 17** *(i) Since* $\varphi_A(\mathbb{H}) \subset \mathbb{H}$, $\varphi_A$ *defines a flow on* $\mathbb{H}$.
*(ii) If* $c \neq 0$, *then,* $\{-d/c + yi : y \in \mathbb{R}\}$ *and* $\mathbb{R}$ *are invariant manifolds. The map restricted on* $\{-d/c + yi : y \in \mathbb{R}\}$ *is* $s \mapsto 1/(c^2 s)$, *and the map restricted on* $\mathbb{R}$ *is* $\ell \mapsto \varphi_A(\ell)$.

## B.3 Boole transformations

We give an alternative simultaneous proof of [26, Proposition 3.1 and Theorem 3.1] themselves.

For $a > 0$, let

$$\varphi_a(z) := \begin{cases} a(z - z^{-1}) & z \in \overline{\mathbb{H}} \setminus \{0\} \\ 0 & z = 0 \end{cases}.$$

**Proposition 16** *Assume that $x \neq 0$. Let $x' := \varphi_a(x)$,*

$$\ell' := \mathrm{Re}(\varphi_a(\ell + is)) = a\ell\frac{\ell^2 + s^2 - 1}{\ell^2 + s^2}$$

*and*

$$s' := \mathrm{Im}(\varphi_a(\ell + is)) = as\frac{\ell^2 + s^2 + 1}{\ell^2 + s^2}.$$

*Then, $\varphi_a^{-1}(x') = \{x, -1/x\}$, and*

$$C(x'; \ell', s') = \frac{C(x; \ell, s)}{|\varphi_a'(x)|} + \frac{C(-1/x; \ell, s)}{|\varphi_a'(-1/x)|}.$$

Proof.    For ease of notation we let $\theta := \ell + si$. We remark that

$$\varphi_a(y) = \varphi_a(-1/y), \quad y \neq 0, \tag{23}$$

and

$$|\varphi_a(x) - \varphi_a(\theta)| = a|x - \theta|\frac{|x\theta + 1|}{|x\theta|}, \quad x \neq 0, \theta \in \mathbb{H}. \tag{24}$$

Let $F$ be a non-negative Borel measurable function. By (24) and the change of variable formula with $y = \varphi_a(x)$,

$$\int_{\mathbb{R}} F(y)C(y; \varphi_a(\theta))dy = \int_0^\infty F(\varphi_a(x))C(x; \theta)\frac{(x^2 + 1)(|\theta|^2 + 1)}{|x\theta + 1|^2}dx$$

$$= \int_{-\infty}^0 F(\varphi_a(x))C(x; \theta)\frac{(x^2 + 1)(|\theta|^2 + 1)}{|x\theta + 1|^2}dx.$$

Hence,

$$\int_{\mathbb{R}} F(y)C(y; \varphi_a(\theta))dy = \int_{\mathbb{R}} F(\varphi_a(x))C(x; \theta)\frac{(x^2 + 1)(|\theta|^2 + 1)}{2|x\theta + 1|^2}dx$$

$$= \int_{\mathbb{R}} F(\varphi_a(x))C(x; \theta)dx + \int_{\mathbb{R}} F(\varphi_a(x))C(x; \theta)\left(\frac{(x^2 + 1)(|\theta|^2 + 1)}{2|x\theta + 1|^2} - 1\right)dx.$$

Since

$$\frac{(x^2 + 1)(|\theta|^2 + 1)}{|x\theta + 1|^2} = 1 + \frac{|x - \theta|^2}{|x\theta + 1|^2},$$

it holds that

$$\int_{\mathbb{R}} F(\varphi_a(x))C(x; \theta)\left(\frac{(x^2 + 1)(|\theta|^2 + 1)}{2|x\theta + 1|^2} - 1\right)dx$$

$$= \frac{1}{2}\int_{\mathbb{R}} F(\varphi_a(x))C(x; \theta)\left(\frac{|x - \theta|^2}{|x\theta + 1|^2} - 1\right)dx$$

$$= \frac{s}{2\pi}\int_{\mathbb{R}} \frac{F(\varphi_a(x))}{|x\theta + 1|^2}dx - \frac{1}{2}\int_{\mathbb{R}} F(\varphi_a(x))C(x; \theta)dx.$$

By the change of variable formula and (23),

$$\frac{s}{\pi}\int_{\mathbb{R}} \frac{F(\varphi_a(x))}{|x\theta + 1|^2}dx = \int_{\mathbb{R}} F(\varphi_a(x))C(x; \theta)dx,$$

and hence,

$$\int_{\mathbb{R}} F(\varphi_a(x)) C(x;\theta) \left( \frac{(x^2+1)(|\theta|^2+1)}{2|x\theta+1|^2} - 1 \right) dx = 0.$$

Thus we obtain that

$$\int_{\mathbb{R}} F(y) C(y; \varphi_a(\theta)) dy = \int_{\mathbb{R}} F(\varphi_a(x)) C(x;\theta) dx.$$

Let two functions $\varphi_{a,\pm}$ be the restrictions of $\varphi_a$ to $(0,\infty)$ and $(-\infty,0)$ respectively. Then, by Lemma 8,

$$C(y; \varphi_a(\theta)) = \frac{C(\varphi_{a,+}^{-1}(y); \theta)}{\varphi_a'(\varphi_{a,+}^{-1}(y))} + \frac{C(\varphi_{a,-}^{-1}(y); \theta)}{\varphi_a'(\varphi_{a,-}^{-1}(y))}, \quad \text{a.e. } y.$$

Since the functions in the left and right hand sides in the above display are both continuous on $\mathbb{R} \setminus \{0\}$,

$$C(y; \varphi_a(\theta)) = \frac{C(\varphi_{a,+}^{-1}(y); \theta)}{\varphi_a'(\varphi_{a,+}^{-1}(y))} + \frac{C(\varphi_{a,-}^{-1}(y); \theta)}{\varphi_a'(\varphi_{a,-}^{-1}(y))}, \quad \text{for every } y \neq 0.$$

<div align="right">QED.</div>

**Remark 18**

$$D_f \left( p_{\varphi_a(\theta_1)} : p_{\varphi_a(\theta_2)} \right) \neq D_f(p_{\theta_1} : p_{\theta_2})$$

for $a = 2, \theta_1 = i$ and $\theta_2 = 2i$.

# C  Revisiting the KLD between Cauchy densities

We shall prove the following result [11] using complex analysis:

$$D_{\mathrm{KL}}(p_{l_1,s_1} : p_{l_2,s_2}) = \log \left( \frac{(s_1+s_2)^2 + (l_1-l_2)^2}{4s_1 s_2} \right).$$

Proof.

$$D_{\mathrm{KL}}(p_{l_1,s_1} : p_{l_2,s_2}) = \frac{s_1}{\pi} \int_{\mathbb{R}} \frac{\log((z-l_2)^2 + s_2^2)}{(z-l_1)^2 + s_1^2} \mathrm{d}z$$

$$- \frac{s_1}{\pi} \int_{\mathbb{R}} \frac{\log((z-l_1)^2 + s_1^2)}{(z-l_1)^2 + s_1^2} \mathrm{d}z + \log \frac{s_1}{s_2}. \tag{25}$$

As a function of $z$,

$$\frac{\log(z - l_2 + is_2)}{z - l_1 + is_1}$$

is holomorphic on the upper-half plane $\{x + yi : y > 0\}$. By the Cauchy integral formula [49], we have that for sufficiently large $R$,

$$\frac{1}{2\pi i} \int_{C_R^+} \frac{\log(z - l_2 + is_2)}{(z - l_1)^2 + s_1^2} \mathrm{d}z = \frac{\log(l_1 - l_2 + i(s_2 + s_1))}{2s_1 i},$$

where

$$C_R^+ := \{z : |z| = R, \operatorname{Im}(z) > 0\} \cup \{z : \operatorname{Im}(z) = 0, |\operatorname{Re}(z)| \leq R\}.$$

Hence, by $R \to +\infty$, we get

$$\frac{s_1}{\pi} \int_{\mathbb{R}} \frac{\log(z - l_2 + is_2)}{(z - l_1)^2 + s_1^2} dz = \log(l_1 - l_2 + i(s_2 + s_1)). \tag{26}$$

As a function of $z$,

$$\frac{\log(z - l_2 - is_2)}{z - l_1 - is_1}$$

is holomorphic on the lower-half plane $\{x + yi : y < 0\}$. By the Cauchy integral formula again, we have that for sufficiently large $R$,

$$\frac{1}{2\pi i} \int_{C_R^-} \frac{\log(z - l_2 - is_2)}{(z - l_1)^2 + s_1^2} dz = \frac{\log(l_1 - l_2 - i(s_2 + s_1))}{-2s_1 i},$$

where

$$C_R^- := \{z : |z| = R, \mathrm{Im}(z) < 0\} \cup \{z : \mathrm{Im}(z) = 0, |\mathrm{Re}(z)| \le R\}.$$

Hence, by $R \to +\infty$, we get

$$\frac{s_1}{\pi} \int_{\mathbb{R}} \frac{\log(z - l_2 - is_2)}{(z - l_1)^2 + s_1^2} dz = \log(l_1 - l_2 - i(s_2 + s_1)). \tag{27}$$

By Eq. 26 and Eq. 27, we have that

$$\frac{s_1}{\pi} \int_{\mathbb{R}} \frac{\log((z - l_2)^2 + s_2^2)}{(z - l_1)^2 + s_1^2} dz = \log\left((l_1 - l_2)^2 + (s_1 + s_2)^2\right). \tag{28}$$

In the same manner, we have that

$$\frac{s_1}{\pi} \int_{\mathbb{R}} \frac{\log((z - l_1)^2 + s_1^2)}{(z - l_1)^2 + s_1^2} dz = \log(4s_1^2). \tag{29}$$

By substituting Eq. 28 and Eq. 29 into Eq. 25, we obtain the formula Eq. 2.  QED.

**Remark 19** *Thomas Simon [76] also obtained an alternative proof of [11], which uses the Lévy-Khintchine formula and the potential formula for the infinitely divisible distributions, and the Frullani integral.*

# D   Revisiting the chi-squared divergence between Cauchy densities

**Proposition 17**

$$D_\chi^N(p_{l_1,s_1} : p_{l_2,s_2}) = \frac{(l_1 - l_2)^2 + (s_1 - s_2)^2}{2s_1 s_2}. \tag{30}$$

Proof.   We first remark that

$$D_\chi^N(p_{l_1,s_1} : p_{l_2,s_2}) = \int_{\mathbb{R}} \frac{p_{l_2,s_2}^2(x)}{p_{l_1,s_1}(x)} dx - 1.$$

Let $F(z) := \frac{(z-l_1)^2 + s_1^2}{(z - l_2 + is_2)^2}$. Then, this is holomorphic on the upper-half plane $\mathbb{H}$, and,

$$\frac{p_{l_2, s_2}(x)^2}{p_{l_1, s_1}(x)} = \frac{s_2^2}{\pi s_1} \frac{F(x)}{(x - l_2 - is_2)^2}.$$

By the Cauchy integral formula [49], we have that for sufficiently large $R$,

$$\frac{1}{2\pi i} \int_{C_R^+} \frac{F(z)}{(z - l_2 - is_2)^2} \mathrm{d}z = F'(l_2 + is_2),$$

where $C_R^+ := \{z : |z| = R, \mathrm{Im}(z) > 0\} \cup \{z : \mathrm{Im}(z) = 0, |\mathrm{Re}(z)| \leq R\}$.

Since

$$F'(z) = 2\frac{(z - s_1)(z - l_2 + is_2) - (z - l_1)^2 - s_1^2}{(z - l_2 + is_2)^3},$$

we have that

$$\int_{C_R^+} \frac{F(z)}{(z - l_2 - is_2)^2} \mathrm{d}z = \frac{\pi}{2} \frac{(l_1 - l_2)^2 + s_1^2 + s_2^2}{s_2^3}.$$

Now, by $R \to \infty$, we obtain the formula Eq. 30. $\hfill$ QED.

# E  Total variation between densities of a location family

Consider a location family with *even* standard density $p(-x) = p(x)$. Then $p(x - l_1) = p(x - l_2) = p(l_2 - x)$ when $x = \frac{l_1 + l_2}{2}$. Let $\Phi(a) = \int_{-\infty}^{a} p(x)\mathrm{d}x$ denote the standard cumulative density function, $\Phi_{l,s}(a) = \int_{-\infty}^{a} p(\frac{x-l}{s})\mathrm{d}x = \Phi(\frac{a-l}{s})$ with $\Phi_{l,s}(-\infty) = 0$ and $\Phi_{l,s}(+\infty) = 1$. We have $\int_{a}^{b} p(x)\mathrm{d}x = \Phi(b) - \Phi(a)$ and $\int_{a}^{+\infty} p_{l,s}(x)\mathrm{d}x = 1 - \Phi(\frac{a-l}{s})$.

Then the total variation distance between $p_{l_1}$ and $p_{l_2}$ is

$$
\begin{aligned}
D_{\mathrm{TV}}(p_{l_1} : p_{l_2}) &= \frac{1}{2}\left( \int_{-\infty}^{\frac{l_1 + l_2}{2}} |p_{l_1}(x) - p_{l_2}(x)|\mathrm{d}x + \int_{\frac{l_1 + l_2}{2}}^{+\infty} |p_{l_2}(x) - p_{l_1}(x)|\mathrm{d}x \right) \\
&= 2\Phi\left(\frac{|l_1 - l_2|}{2s}\right) - 1 \leq 1
\end{aligned}
$$

**Proposition 18** *The total variation between two densities $p_{l_1}$ and $p_{l_2}$ of a location family with even standard density is $2\Phi\left(\frac{|l_1 - l_2|}{2s}\right) - 1$.*

For the Cauchy distribution, since we have

$$\Phi_{l,s}(x) = \frac{1}{\pi}\arctan\left(\frac{x - l}{s}\right) + \frac{1}{2},$$

we recover $D_{\mathrm{TV}}(p_{l_1} : p_{l_2}) = \frac{2}{\pi}\arctan\left(\frac{|l_2 - l_1|}{2s}\right)$.

The total variation formula extends to any fixed scale location families.

# F Complete elliptic integrals

This section is devoted to the details of the proof of (18) in the proof of Theorem 9.

Proof. Let

$$F_4(u) := \frac{-\log\left(2e^{-u/4}\mathbf{K}(1-e^{-u})/\pi\right)}{u^2}.$$

We consider the derivative.

$$F_4'(u) = \frac{-1}{u^2}\left(\frac{1}{4} + e^{-u}\frac{\mathbf{K}'(1-e^{-u})}{\mathbf{K}(1-e^{-u})} - \frac{2}{u}\log\left(2\mathbf{K}(1-e^{-u})/\pi\right)\right).$$

Now it suffices to show that for every $u > 0$,

$$\frac{1}{4} + e^{-u}\frac{\mathbf{K}'(1-e^{-u})}{\mathbf{K}(1-e^{-u})} - \frac{2}{u}\log\left(2\mathbf{K}(1-e^{-u})/\pi\right) > 0.$$

Let $x := 1 - e^{-u}$. Then, it suffices to show that for every $x \in (0,1)$,

$$\frac{1}{4} + (1-x)\frac{\mathbf{K}'(x)}{\mathbf{K}(x)} + \frac{2}{\log(1-x)}\log\left(2\mathbf{K}(x)/\pi\right) > 0.$$

Let

$$G_4(x) := \log\left(2\mathbf{K}(x)/\pi\right) + (\log(1-x))\left(\frac{1}{8} + \frac{1-x}{2}\frac{\mathbf{K}'(x)}{\mathbf{K}(x)}\right).$$

It suffices to show that $G_4(x) < 0$ for every $x \in (0,1)$.

We see that $G_4(0) = 0$. Hence it suffices to show that $G_4'(x) < 0$ for every $x \in (0,1)$. By Lemma 14 below,

$$G_4(x) = \log\left(2\mathbf{K}(x)/\pi\right) + (\log(1-x))\left(\frac{3}{8} + \frac{1}{4x}\left(\cdot\frac{\mathbf{E}(x)}{\mathbf{K}(x)} - 1\right)\right).$$

By Lemmas 14 and 15 below,

$$G_4'(x) = -\frac{H_4(x)}{8x^2(1-x)},$$

where we let

$$H_4(x) := (x(2-x) + (x-1)\log(1-x))\mathbf{K}(x)^2 - 2x\mathbf{K}(x)\mathbf{E}(x) + \log(1-x)\mathbf{E}(x)^2.$$

Then it suffices to show that $H_4(x) > 0$ for every $x \in (0,1)$. Since $-2x < 0$ and $\log(1-x) < 0$, by noting Lemma 16 below, it holds that

$$\frac{H_4(x)}{\mathbf{K}(x)^2} \geq (x(2-x) + (x-1)\log(1-x)) - 2xI_4(x) + \log(1-x)I_4(x)^2,$$

where we let

$$I_4(x) := \frac{1}{2} - \frac{x}{4} + \frac{\sqrt{1-x}}{2}.$$

Our main idea is to use different estimates for $H(x)/\mathbf{K}(x)^2$ on a neighborhood of 1 and on the compliment of it.

**Lemma 10** *For $x \leq 0.998$,*

$$(x(2 - x) + (x - 1)\log(1 - x)) - 2xI(x) + \log(1 - x)I(x)^2 > 0.$$

Proof.    Let $y := \sqrt{1 - x}$. Then,

$$(x(2 - x) + (x - 1)\log(1 - x)) - 2xI(x) + \log(1 - x)I(x)^2 > 0$$

is equivalent with

$$\log y > 4\frac{y^2 - 1}{y^2 + 6y + 1}.$$

Let

$$P_4(y) := \log y - 4\frac{y^2 - 1}{y^2 + 6y + 1}.$$

Then, $P_4(1) = 0$. By considering the derivative of $P$, it is increasing $y < 5 - 2\sqrt{6}$ and decreasing $y > 5 - 2\sqrt{6}$.

We see that $P_4(y) > 0, \quad y > 0.041$. Now the assertion follows from the fact that

$$0.998 < 1 - (0.041)^2.$$

QED.

Now it suffices to show that $H(x) > 0$ for $x > 0.998$.

**Lemma 11**

$$x(2 - x) + (x - 1)\log(1 - x) \geq 1, \quad x \in (0.998, 1).$$

Proof.    Let $g_4(x) := x(2 - x) + (x - 1)\log(1 - x)$. Then, $g_4(1) = 1$ and

$$g_4'(x) = 3 - 2x + \log(1 - x).$$

This is negative if $x > 0.9$. QED.

**Lemma 12**

$$2x.\frac{\mathbf{E}(x)}{\mathbf{K}(x)} < \frac{1}{2}, \quad x \in (0.998, 1).$$

Proof.    We see that

$$\frac{d}{dx}\left(x.\frac{\mathbf{E}(x)}{\mathbf{K}(x)}\right) \leq 2.\frac{\mathbf{E}(x)}{\mathbf{K}(x)} - \frac{1}{2}.$$

By Lemma 15 below and the fact that

$$\frac{\mathbf{E}(0.995)}{\mathbf{K}(0.995)} < \frac{1}{4},$$

we see that

$$2.\frac{\mathbf{E}(x)}{\mathbf{K}(x)} \leq \frac{1}{2}, \quad x > 0.995.$$

Hence,

$$2x.\frac{\mathbf{E}(x)}{\mathbf{K}(x)} < 2\frac{\mathbf{E}(0.995)}{\mathbf{K}(0.995)} < \frac{1}{2}.$$

QED.

**Lemma 13**

$$-\log(1-x)\left(\cdot\frac{\mathbf{E}(x)}{\mathbf{K}(x)}\right)^2 < \frac{1}{2}, \quad x \in (0.998, 1).$$

Proof. We use Lemma 17 below. It suffices to show that

$$\frac{2x^{1/2}}{\log(1+x^{1/2}) - \log(1-x^{1/2})} \le \sqrt{\frac{1}{-2\log(1-x)}}, \quad x \in (0.998, 1).$$

This is equivalent with

$$h_4(x) := \left(\log(1+x^{1/2}) - \log(1-x^{1/2})\right)^2 + 8x\log(1-x) \ge 0, \quad x \in (0.998, 1).$$

We see that

$$h_4'(x) = -2\frac{\log(1-\sqrt{x}) - \log(1+\sqrt{x}) + 2\sqrt{x}(x + (x-1)\log(1-x))}{(1-x)\sqrt{x}}.$$

It is easy to see that

$$\log(1-\sqrt{x}) - \log(1+\sqrt{x}) + 2\sqrt{x}(x + (x-1)\log(1-x)) < 0, \quad x \in (0.998, 1).$$

Hence $h_4$ is increasing at least on $(0.998, 1)$. Now use the fact that $h_4(0.998) > 0$. QED.

By Lemmas 11, 12 and 13, we see that $H_4(x) > 0$ for $x > 0.998$. The proof of Eq. 18 is completed. QED.

## F.1 Some Lemmas concerning the complete elliptic integrals

In this subsection, we collect standard results about the complete elliptic integrals.

**Lemma 14**

$$\mathbf{K}'(x) = -\frac{\mathbf{K}(x)}{2x} + \frac{\mathbf{E}(x)}{2x(1-x)}.$$

**Lemma 15**

$$\frac{d}{dx}\left(\cdot\frac{\mathbf{E}(x)}{\mathbf{K}(x)}\right) = -\frac{1}{2x} + \frac{1}{x}\cdot\frac{\mathbf{E}(x)}{\mathbf{K}(x)} - \frac{1}{2x(1-x)}\left(\cdot\frac{\mathbf{E}(x)}{\mathbf{K}(x)}\right)^2 \le 0.$$

*In particular, $\mathbf{E}/\mathbf{K}$ is strictly decreasing.*

**Lemma 16**

$$\cdot\frac{\mathbf{E}(x)}{\mathbf{K}(x)} \le \frac{1}{2} - \frac{x}{4} + \frac{\sqrt{1-x}}{2}, \quad x \in [0, 1).$$

The following is due to Anderson, Vamanamurthy, and Vuorinen [4].

**Lemma 17 ([4, Theorem 3.6])**

$$\cdot\frac{\mathbf{E}(x)}{\mathbf{K}(x)} \le \frac{2x^{1/2}}{\log(1+x^{1/2}) - \log(1-x^{1/2})}, \quad x \in [0, 1).$$

The following is due to Eq. (1.1) in [10]. See also Eq. (6.2) in [4].

**Lemma 18**

$$\log\left(\frac{4}{\sqrt{1-x}}\right) \le \mathbf{K}(x) \le \frac{4}{3+x}\log\left(\frac{4}{\sqrt{1-x}}\right), \quad x \in [0, 1).$$

# G   Negative definiteness of the KLD between Cauchy densities

In this section, we give an elementary proof of Theorem 11. We first give an outline of the proof. Our proof follows the strategy of [22] and consists of three steps. We do not need to introduce the hyperboloid space $\mathbb{L}$.

Step 1.    Let $d$ be the Poincaré distance on $\mathbb{H}$.   We remark that $d = \sqrt{2}\rho_{\mathrm{FR}}$.    Then, $\cosh(d(z,w)) = 1 + \chi(z,w)$ and

$$2\log\cosh\left(\frac{d(z,w)}{2}\right) = \log\left(1 + \frac{\chi(z,w)}{2}\right).$$

We see that for every $r \geq 0$,

$$2\log\cosh\left(\frac{r}{2}\right) = \lim_{s\to+0}\frac{1}{s}\left(1 - \frac{1}{2\pi}\int_{-\pi}^{\pi}(\cosh(r) + \cos\theta\sinh(r))^{-s}\,d\theta\right).$$

Hence it suffices to show that

$$H_s(z,w) := \frac{1}{2\pi}\int_{-\pi}^{\pi}(\cosh(d(z,w)) + \cos\theta\sinh(d(z,w)))^{-s}\,d\theta, \ z,w \in \mathbb{H},$$

is positive definite for every $s \in (0,1)$.

Step 2.   Let

$$P(z,x) := \frac{\mathrm{Im}(z)}{|x-z|^2}(x^2 + 1), \ z \in \mathbb{H}, x \in \mathbb{R},$$

and $\mu(dx) := \dfrac{dx}{\pi(x^2 + 1)}.$

Then we see that

$$H_s(z,w) = \int_{\mathbb{R}} P(z,x)^s P(w,x)^{1-s}\mu(dx)$$

$$= C(s)\iint_{\mathbb{R}^2} P(w,x)^{1-s}P(z,y)^{1-s}\left(\frac{(x-y)^2}{(x^2+1)(y^2+1)}\right)^{-s}\mu(dx)\mu(dy),$$

where $C(s)$ is a positive constant depending only on $s$.

Step 3.   Let $z_1,\cdots,z_n \in \mathbb{H}$ and $c_1,\cdots,c_n \in \mathbb{R}$ with $\sum_{i=1}^{n}c_i = 0$. Let

$$\varphi_s(x) := \sum_{i=1}^{n}c_i P(z_i,x)^{1-s}, \ x \in \mathbb{R},$$

which is continuous on $\mathbb{R}$.

Let $k_s(x,y) := \left(\frac{(x-y)^2}{(x^2+1)(y^2+1)}\right)^{-s}$, which is a positive definite kernel on $\mathbb{R}$.

Thus we see that

$$\sum_{i.j=1}^{n}c_i c_j H_s(z_i,z_j) = \frac{C(s)}{\pi^2}\iint_{\mathbb{R}^2}\frac{\varphi_s(x)\varphi_s(y)k_s(x,y)}{(x^2+1)(y^2+1)}dxdy \geq 0.$$

Now we proceed to the full proof.

Step 1. It is known that (see formula no.4.224.9 in [27])

$$2 \log \cosh \left( \frac{r}{2} \right) = \frac{1}{2\pi} \int_{-\pi}^{\pi} \log \left( \cosh(r) + \cos \theta \sinh(r) \right) d\theta, \quad r \geq 0.$$

We see that for $r \geq 0$,

$$|\log(\cosh(r) + \cos \theta \sinh(r))| \leq r.$$

Since for $t > 0$, $\lim_{s \to +0} \frac{1-t^{-s}}{s} = \log t$ and $|\frac{1-t^{-s}}{s}| \leq |\log t|$,

$$\int_{-\pi}^{\pi} \log \left( \cosh(r) + \cos \theta \sinh(r) \right) d\theta = \lim_{s \to +0} \int_{-\pi}^{\pi} \frac{1 - (\cosh(r) + \cos \theta \sinh(r))^{-s}}{s} d\theta, \quad r > 0,$$

by the Lebesgue convergence theorem. This convergence also holds for $r = 0$.

Step 2.

**Lemma 19**

$$\frac{1}{2\pi} \int_{-\pi}^{\pi} (\cosh(r) + \cos \theta \sinh(r))^{-s} d\theta = \int_{\mathbb{R}} P(e^r i, x)^s \mu(dx).$$

Proof. Let $x = \tan \frac{\theta}{2}$. Then, $d\theta = \frac{2}{1 + x^2} dx$ and

$$\cosh(r) + \cos \theta \sinh(r) = \frac{e^{2r} + x^2}{e^r (1 + x^2)} = \frac{1}{P(e^r i, x)}.$$

QED.

**Lemma 20** *For $A \in SO(2)$ and $z \in \mathbb{H}$,*

$$\int_{\mathbb{R}} P(A.z, x)^s \mu(dx) = \int_{\mathbb{R}} P(z, x)^s \mu(dx).$$

Proof. Let $A = \begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix}$. Let $y \in \mathbb{R}$ such that $x = A.y$. Then,

$$P(A.z, A.y) = P(z, y)$$

and

$$\mu(dx) = \frac{1}{\pi} \frac{1}{(A.y)^2 + 1} \frac{dx}{dy} dy = \frac{1}{\pi} \frac{1}{y^2 + 1} dy = \mu(dy).$$

QED.

Now we introduce a group structure on $\mathbb{H}$. For $z = z_1 + iz_2$ and $w = w_1 + iw_2$, let

$$zw := (z_1 + z_2 w_1) + iz_2 w_2.$$

This gives a group structure on $\mathbb{H}$. It holds that

$$w^{-1} = \frac{-w_1 + i}{w_2}, \quad w = w_1 + iw_2$$

and the unit element is the imaginary unit $i$.

We see that $\chi(w^{-1}z, i) = \chi(z, w), z, w \in \mathbb{H}$ and hence

$$d(w^{-1}z, i) = d(z, w), z, w \in \mathbb{H}. \tag{31}$$

60

**Lemma 21** *For $z, w \in \mathbb{H}$,*

$$H_s(z, w) = \int_{\mathbb{R}} P(w^{-1}z, x)^s \mu(dx).$$

Proof.    By (31), we can assume that $w = i$. Then there exists $A \in SO(2)$ such that $e^{d(z,i)}i = A.z$. Now the assertion follows from Lemmas 19 and 20.                QED.

For $w = w_1 + iw_2 \in \mathbb{H}$ and $x \in \mathbb{R}$, we let $wx := w_2 x + w_1$.

**Lemma 22**

$$P(w^{-1}z, x)P(w, wx) = P(z, wx), \quad z, w \in \mathbb{H}, x \in \mathbb{R}.$$

Proof.    Since

$$w^{-1}z = \frac{z_1 - w_1 + iz_2}{w_2}, \quad z = z_1 + iz_2, w = w_1 + iw_2,$$

we see that

$$P(w^{-1}z, x) = \frac{z_2 w_2}{(z_1 - w_1 - w_2 x)^2 + z_2^2}(x^2 + 1).$$

We also see that

$$P(z, wx) = \frac{z_2((w_2 x + w_1)^2 + 1)}{(z_1 - w_1 - w_2 x)^2 + z_2^2}$$

and

$$P(w, wx) = \frac{(w_2 x + w_1)^2 + 1}{w_2(x^2 + 1)}.$$

The assertion follows from these identities.                QED.

**Proposition 19**

$$H_s(z, w) = \int_{\mathbb{R}} P(z, x)^s P(w, x)^{1-s} \mu(dx), \quad z, w \in \mathbb{H}.$$

Proof.    By Lemmas 21 and 22,

$$H_s(z, w) = \int_{\mathbb{R}} P(z, wx)^s P(w, wx)^{-s} \mu(dx).$$

Let $y = wx = w_2 x + w_1$. Then, $\mu(dx) = \dfrac{w_2}{\pi|y - w|^2} dy$. Hence,

$$\int_{\mathbb{R}} P(z, wx)^s P(w, wx)^{-s} \mu(dx) = \int_{\mathbb{R}} P(z, y)^s P(w, y)^{1-s} \mu(dy).$$

QED.

**Lemma 23** *For every $s \in (0, 1/2)$, there exists a positive constant $C(s)$ such that for every $a \in \mathbb{R}$*

$$(1 + a^2)^{-s} = \frac{C(s)}{\pi} \int_{\mathbb{R}} \frac{|x + a|^{-2s}}{(1 + x^2)^{1-s}} dx.$$

**Proof.** Let $x = \tan\theta$, $|\theta| < \pi/2$. Then, $d\theta = \cos^2\theta \, dx = \frac{1}{1+x^2} dx$ and

$$\frac{(x+a)^2}{1+x^2} = (\sin\theta + a\cos\theta)^2.$$

Hence,

$$\int_{\mathbb{R}} \frac{|x|^{-2s}}{(1+(x-a)^2)^{1-s}} dx = \int_{-\pi/2}^{\pi/2} |\sin\theta + a\cos\theta|^{-2s} \, d\theta.$$

By symmetry,

$$\int_{-\pi/2}^{\pi/2} |\sin\theta + a\cos\theta|^{-2s} \, d\theta = \frac{1}{2} \int_{-\pi}^{\pi} |\sin\theta + a\cos\theta|^{-2s} \, d\theta = \pi(1+a^2)^{-s} \int_{-\pi}^{\pi} |\cos\theta|^{-2s} \, d\theta.$$

The assertion holds if we let $C(s) := \left( \int_{-\pi}^{\pi} |\cos\theta|^{-2s} \, d\theta \right)^{-1}$. QED.

**Lemma 24 (intertwining formula)** *For every* $s \in (0, 1/2)$, $w \in \mathbb{H}$ *and* $y \in \mathbb{R}$,

$$P(w,y)^s = C(s) \int_{\mathbb{R}} P(w,x)^{1-s} \left( \frac{(x-y)^2}{(x^2+1)(y^2+1)} \right)^{-s} \mu(dx). \tag{32}$$

**Proof.** Let $\xi := w - y$ and $t := x - y$. Then, (32) holds if and only if

$$\left( \frac{\operatorname{Im}(\xi)}{|\xi|^2} \right)^s = \frac{C(s)}{\pi} \int_{\mathbb{R}} \left( \frac{\operatorname{Im}(\xi)}{|\xi - t|^2} \right)^{1-s} |t|^{-2s} dt. \tag{33}$$

Let $u := (t - \operatorname{Re}(\xi))/\operatorname{Im}(\xi)$. Then,

$$\int_{\mathbb{R}} \left( \frac{\operatorname{Im}(\xi)}{|\xi - t|^2} \right)^{1-s} |t|^{-2s} dt = (\operatorname{Im}(\xi))^{-s} \int_{\mathbb{R}} \left( \frac{1}{1+u^2} \right)^{1-s} \left| u + \frac{\operatorname{Re}(\xi)}{\operatorname{Im}(\xi)} \right|^{-2s} du.$$

Hence (33) holds if and only if

$$\left( \left( \frac{\operatorname{Re}(\xi)}{\operatorname{Im}(\xi)} \right)^2 + 1 \right)^{-s} = \frac{C(s)}{\pi} \int_{\mathbb{R}} \left( \frac{1}{1+u^2} \right)^{1-s} \left| u + \frac{\operatorname{Re}(\xi)}{\operatorname{Im}(\xi)} \right|^{-2s} du,$$

which follows from Lemma 23. QED.

By Proposition 19 and Lemma 24,

**Proposition 20** *For every* $s \in (0, 1/2)$,

$$H_s(z, w) = C(s) \iint_{\mathbb{R}^2} P(w,x)^{1-s} P(z,y)^{1-s} \left( \frac{(x-y)^2}{(x^2+1)(y^2+1)} \right)^{-s} \mu(dx)\mu(dy), \quad z, w \in \mathbb{H}.$$

Step 3.

**Lemma 25** $k_s(x,y)$ *is a positive definite kernel on* $\mathbb{R}$.

Proof. For $r \in (0,1)$, let

$$k_s^{(r)}(x,y) := \left( 1 - r \frac{(xy+1)^2}{(x^2+1)(y^2+1)} \right)^{-s}.$$

Since $(x,y) \mapsto \frac{1}{(x^2+1)(y^2+1)}$ and $(x,y) \mapsto (xy)^2 + 2xy + 1$ are both positive definite kernels on $\mathbb{R}$, $(x,y) \mapsto \frac{(xy+1)^2}{(x^2+1)(y^2+1)}$ is also a positive definite kernel on $\mathbb{R}$.

By the Taylor expansion,

$$(1-x)^{-s} = \sum_{n=0}^{\infty} a_n x^n, \ |x| < 1,$$

for $a_n \geq 0, n = 0, 1, \cdots$. Hence $k_s^{(r)}(x,y)$ is a positive definite kernel on $\mathbb{R}$. Since $\lim_{r \to 1-0} k_s^{(r)}(x,y) = k_s(x,y)$, $k_s^{(r)}(x,y)$ is a positive definite kernel on $\mathbb{R}$. QED.

By this and the quadrature rule for the Riemannian integral for continuous functions, it holds that for every $a < b$ and $r \in (0,1)$,

$$\iint_{[a,b]^2} \frac{\varphi_s(x)\varphi_s(y)k_s^{(r)}(x,y)}{(x^2+1)(y^2+1)} dxdy \geq 0.$$

Since $0 \leq k_s^{(r)}(x,y) \leq k_s(x,y)$,

$$\iint_{\mathbb{R}^2} \frac{|\varphi_s(x)\varphi_s(y)|k_s^{(r)}(x,y)}{(x^2+1)(y^2+1)} dxdy \leq \iint_{\mathbb{R}^2} \frac{|\varphi_s(x)\varphi_s(y)|k_s(x,y)}{(x^2+1)(y^2+1)} dxdy \leq \sum_{i,j=1}^{n} |c_i||c_j|H_s(z_i,z_j) < +\infty.$$

By the Lebesgue convergence theorem, we see that for every $r \in (0,1)$,

$$\iint_{\mathbb{R}^2} \frac{\varphi_s(x)\varphi_s(y)k_s^{(r)}(x,y)}{(x^2+1)(y^2+1)} dxdy = \lim_{n \to \infty} \iint_{[-n,n]^2} \frac{\varphi_s(x)\varphi_s(y)k_s^{(r)}(x,y)}{(x^2+1)(y^2+1)} dxdy \geq 0.$$

and furthermore,

$$\iint_{\mathbb{R}^2} \frac{\varphi_s(x)\varphi_s(y)k_s(x,y)}{(x^2+1)(y^2+1)} dxdy = \lim_{r \to 1-0} \iint_{\mathbb{R}^2} \frac{\varphi_s(x)\varphi_s(y)k_s^{(r)}(x,y)}{(x^2+1)(y^2+1)} dxdy \geq 0.$$

This completes the proof.

# H Code snippet for Taylor expansions of $f$-divergences

We provide below a code using the MAXIMA[4] software to calculate the truncated Taylor series of $f$-divergences between two Cauchy distributions.

```
Cauchy(x,l,s) := (s/(%pi*((x-l)**2+s**2)));
KLCauchy(l1,s1,l2,s2) := log(((s1+s2)**2+(l1-l2)**2)/(4*s1*s2)) ;
l1:0;
s1:1;
```

---
[4] https://maxima.sourceforge.io/

```
l2:0.6;
s2:6/5;
k:40;
testcond: (9/16)-(l2**2+(s2-(4/5))**2);
print("Is condition>0 for Taylor expansion?:",testcond);
Cauchy1:Cauchy(x,l1,s1);
Cauchy2:Cauchy(x,l2,s2);
print("Exact KL");
KLCauchy(l1,s1,l2,s2);
ExactKL:float(%);
print("KL numerical integration:");
kla: quad_qagi( Cauchy1*log(Cauchy1/Cauchy2), x, minf, inf,'epsrel=1d-10);
NumKL:float(kla[1]);

for i:2 while (i<=k)
do( r[i]: quad_qagi( (Cauchy1-Cauchy2)**i/Cauchy2**(i-1), x, minf, inf,'epsrel=1d-10),
 print(i,r[i][1]));

print("KL Taylor truncated series:");
TaylorKL: sum( (((-1)**i)/i)*r[i][1], i, 2, k);
print("Exact:",ExactKL,"Numerical:",NumKL,"Trunc. Taylor", TaylorKL);
print("Error |Taylor-Exact|",abs(TaylorKL-ExactKL));
```