

Minimization Problems Based on Relative α -Entropy II: Reverse Projection

M. Ashok Kumar and Rajesh Sundaresan

Abstract

In part I of this two-part work, certain minimization problems based on a parametric family of relative entropies (denoted \mathcal{I}_α) were studied. Such minimizers were called forward \mathcal{I}_α -projections. Here, a complementary class of minimization problems leading to the so-called reverse \mathcal{I}_α -projections are studied. Reverse \mathcal{I}_α -projections, particularly on log-convex or power-law families, are of interest in robust estimation problems ($\alpha > 1$) and in constrained compression settings ($\alpha < 1$). Orthogonality of the power-law family with an associated linear family is first established and is then exploited to turn a reverse \mathcal{I}_α -projection into a forward \mathcal{I}_α -projection. The transformed problem is a simpler quasiconvex minimization subject to linear constraints.

Index Terms

Best approximant; exponential family; information geometry; Kullback-Leibler divergence; linear family; power-law family; projection; Pythagorean property; relative entropy; Rényi entropy; robust estimation; Tsallis entropy.

I. INTRODUCTION

This paper is a continuation of our study of minimization problems based on a parametric generalization of relative entropies, denoted \mathcal{I}_α . See (12) for the definition of $\mathcal{I}_\alpha(P, Q)$, where P and Q are probability measures on an alphabet set \mathbb{X} . We say “parametric generalization of relative entropy” because $\lim_{\alpha \rightarrow 1} \mathcal{I}_\alpha(P, Q) = \mathcal{I}(P||Q)$, the usual relative entropy or Kullback-Leibler divergence. In part I [2], we showed how \mathcal{I}_α arises and studied the problem of a *forward \mathcal{I}_α -projection*, namely

$$\min_{P \in \mathbb{E}} \mathcal{I}_\alpha(P, R),$$

where R is a fixed probability measure on \mathbb{X} and \mathbb{E} is a convex set of probability measures on \mathbb{X} . In this paper, we shall study *reverse \mathcal{I}_α -projection*, namely

$$\min_{P \in \mathbb{E}} \mathcal{I}_\alpha(R, P).$$

The minimization now is with respect to the second argument of \mathcal{I}_α . Such problems arise in robust parameter estimation and constrained compression settings. The family \mathbb{E} is usually a parametric family such as the exponential family, or its generalization, called the *α -power-law family*.

We shall bring to light the geometric relation between the *α -power-law family* and a *linear family*¹ of probability measures. We shall turn the reverse \mathcal{I}_α -projection problem on an α -power-law family into a forward \mathcal{I}_α -projection problem on a linear family. The latter turns out to be a minimization of a quasiconvex objective function subject to linear constraints.

The outline of the paper is as follows. In Section II, we motivate reverse \mathcal{I}_α -projections for the cases $\alpha > 1$ and $\alpha < 1$. In Section III, we define the required terminologies and highlight the contributions of the paper. In Section IV, we study the existence of a reverse \mathcal{I}_α -projection on general log-convex sets. In Section V, we provide simplified proofs of some essential results from [2] on the forward \mathcal{I}_α -projection. Our simplified proofs also serve the purpose of keeping this paper self-contained. In Section VI, we explore the geometric relation between the α -power-law and the linear families, and then exploit it to study reverse \mathcal{I}_α -projection on α -power-law families. The paper ends with some concluding remarks in Section VII.

II. MOTIVATIONS

The purpose of this section is to motivate reverse \mathcal{I}_α -projections. The motivation for $\alpha > 1$ comes from robust statistics. The motivation for $\alpha < 1$ comes from information theory as well as from a strong similarity of the outcomes with the $\alpha = 1$ (relative entropy) setting.

M. Ashok Kumar was supported by a Council for Scientific and Industrial Research (CSIR) fellowship and by the Department of Science and Technology. R. Sundaresan was supported in part by the University Grants Commission by Grant Part (2B) UGC-CAS-(Ph.IV) and in part by the Department of Science and Technology. A part of the material in this paper (Section V alone) was presented at the National Conference on Communication (NCC 2015), Mumbai, India, held during February 2015 [1].

M. Ashok Kumar and R. Sundaresan are with the ECE Department, Indian Institute of Science, Bangalore 560012, India.

¹Example linear families are (1) the set of probability measures P on \mathbb{X} such that $\sum_x P(x)f(x) = 0$ for some $f: \mathbb{X} \rightarrow \mathbb{R}$, and (2) finite intersections of such sets. If there is an additive structure on \mathbb{X} , a concrete example is the set of all probability measures with a fixed mean.

A. Reverse \mathcal{I} -projection

Let \mathbb{X} be a finite alphabet set and let $\mathbb{E} = \{P_\theta : \theta \in \Theta\}$ denote a family of probability measures on \mathbb{X} indexed by the elements of the index set $\Theta \subset \mathbb{R}^k$ for some k . Let x_1, x_2, \dots, x_n be n samples drawn independently and with replacement from \mathbb{X} according to an unknown probability measure P_θ belonging to \mathbb{E} . The maximum likelihood estimate (MLE) of θ , denoted $\hat{\theta}$, is the element of the index set Θ that maximizes the likelihood (if it exists), i.e.,

$$\hat{\theta} = \arg \max_{\theta \in \Theta} \prod_{i=1}^n P_\theta(x_i). \quad (1)$$

Let \hat{P} denote the empirical measure of the n samples x_1, \dots, x_n , i.e.,

$$\hat{P} := \frac{1}{n} \sum_{i=1}^n \delta_{x_i},$$

where δ_a denotes the Dirac mass at a . One may then write

$$\begin{aligned} \frac{\prod_{i=1}^n P_\theta(x_i)}{\prod_{i=1}^n \hat{P}(x_i)} &= \prod_{i=1}^n \frac{P_\theta(x_i)}{\hat{P}(x_i)} \\ &= \prod_{x \in \mathbb{X}} \left(\frac{P_\theta(x)}{\hat{P}(x)} \right)^{n\hat{P}(x)} \\ &= \exp\{-n\mathcal{I}(\hat{P} \| P_\theta)\}, \end{aligned}$$

where

$$\mathcal{I}(P \| Q) := \sum_{x \in \mathbb{X}} P(x) \log \frac{P(x)}{Q(x)}$$

is the relative entropy² of P with respect to Q . Hence the MLE is the minimizer (if it exists)

$$\hat{\theta} = \arg \min_{\theta \in \Theta} \mathcal{I}(\hat{P} \| P_\theta), \quad (2)$$

and the corresponding probability measure $P_{\hat{\theta}}$ is known as the *reverse \mathcal{I} -projection of \hat{P} on the family \mathbb{E}* . Such reverse projections, particularly those related to robustifications of the MLE, are the subject matter of this paper.

Observe that the MLE depends on the samples only through their empirical measure. Let us write the MLE as a function of the empirical measure in a different way. Assume that the family \mathbb{E} is sufficiently smooth in the parameter θ on account of which we can define the *score function* as $s(\cdot; \theta) := \nabla_\theta \log P_\theta(\cdot)$, the gradient of $\log P_\theta(\cdot)$ with respect to θ . The first order optimality criterion applied to (1) after taking logarithms yields the so-called estimating equation for the MLE:

$$\frac{1}{n} \sum_{i=1}^n s(x_i; \theta) = 0;$$

the MLE $\hat{\theta}$ solves this equation. Write $E_P[\cdot \cdot \cdot]$ for expectation with respect to P . Noting that the score function satisfies

$$E_{P_\theta}[s(X; \theta)] = 0 \quad \forall P_\theta,$$

the estimating equation for the MLE can be rewritten as

$$\frac{1}{n} \sum_{i=1}^n s(x_i; \theta) = E_{P_\theta}[s(X; \theta)], \quad (3)$$

which is the same as

$$E_{\hat{P}}[s(X; \theta)] = E_{P_\theta}[s(X; \theta)]. \quad (4)$$

If we write $T(\hat{P})$ for the θ that solves (4), we then have $\hat{\theta} = T(\hat{P})$. The estimator $T(\hat{P})$ is *Fisher consistent*³, a fact that can be easily checked using (4).

²The usual convention is $p \log \frac{p}{q} = 0$ if $p = 0$ and $+\infty$ if $p > q = 0$.

³An estimator that maps an empirical measure to an element in Θ is *Fisher consistent* if it is continuous and maps P_θ to the true parameter θ . See [3, Sec. 5c.1]

B. Reverse \mathcal{I}_α -projection: $\alpha > 1$

Though the MLE is known to possess many good properties, asymptotic efficiency being an example, it is not appropriate when some of the data entries (x_i) are contaminated by outliers. To achieve robustness, one may consider scaling the scores $s(x_i; \theta)$ in the left-hand side of (3) by weights $w(x; \theta)$ that weigh down outlying observations “relative to the model” (see for example Basu et al. [4]). This type of robustification, along with the requirement of Fisher consistency, is accomplished by the estimator that maps the empirical measure \hat{P} to the θ that solves the equation

$$E_{\hat{P}}[w(X; \theta)s(X; \theta)] = E_{P_\theta}[w(X; \theta)s(X; \theta)]. \quad (5)$$

Basu et al. [4] proposed the natural weighting $w(x; \theta) = P_\theta(x)^c$ where $c > 0$. As another robustification procedure, Basu et al. [4] proposed a weighting of the model *by itself*, motivated by the works of Field and Smith [5] and Windham [6], prior to solving the estimating equation. Their procedure is as follows. Given a measure Q , its weighting with respect to a parameter $c > 0$ and a model $\theta \in \Theta$, denoted $Q^{(c, \theta)}$, is given by

$$Q^{(c, \theta)}(x) = \frac{w(x; \theta)Q(x)}{\sum_{y \in \mathbb{X}} w(y; \theta)Q(y)}, \quad x \in \mathbb{X},$$

where the dependence on c is through the weighting $w(x; \theta) = P_\theta(x)^c$ as before. Observe that $(P_\theta)^{(c, \theta)}$ weighs P_θ by itself, namely the weighting parameters are c and θ , and $(P_\theta)^{(c, \theta)}$ is the probability measure proportional to P_θ^{c+1} . The Basu et al. procedure⁴ [4] is to find the θ that solves the equation

$$E_{(\hat{P})^{(c, \theta)}}[s(X; \theta)] = E_{(P_\theta)^{(c, \theta)}}[s(X; \theta)]; \quad (6)$$

the \hat{P} and P_θ of (4) are replaced by the model reweighted $(\hat{P})^{(c, \theta)}$ and $(P_\theta)^{(c, \theta)}$, respectively. It is clear that the corresponding estimator is Fisher consistent. Now (6) can be rewritten as

$$\frac{\frac{1}{n} \sum_{i=1}^n w(x_i; \theta)s(x_i; \theta)}{\frac{1}{n} \sum_{i=1}^n w(x_i; \theta)} = \frac{\mathbb{E}_{P_\theta}[w(X; \theta)s(X; \theta)]}{\mathbb{E}_{P_\theta}[w(X; \theta)]},$$

which expands to

$$\frac{\sum_{i=1}^n P_\theta(x_i)^c s(x_i; \theta)}{\sum_{i=1}^n P_\theta(x_i)^c} = \frac{\sum_{x \in \mathbb{X}} P_\theta(x)^{c+1} s(x; \theta)}{\sum_{x \in \mathbb{X}} P_\theta(x)^{c+1}}. \quad (7)$$

Jones et al. [7] compare the robustness properties of estimators arising from (5) and (7). According to Jones et al. [7, p. 866], the former is more efficient, but the latter has better robustness with respect to a mixture model of contamination with outliers.

Equation (7) can be recognized as an estimating equation arising from the first order optimality criterion for the maximization

$$\max_{\theta \in \Theta} \left[\frac{1}{c} \log \left(\frac{1}{n} \sum_{i=1}^n P_\theta(x_i)^c \right) - \frac{1}{1+c} \log \sum_{x \in \mathbb{X}} P_\theta(x)^{1+c} \right]. \quad (8)$$

We shall soon see why it ought to be a maximization. The objective function in (8) is called *mean power likelihood*⁵. The corresponding estimator is called the maximum mean power likelihood estimate (MMPLE) by Eguchi and Kato [8]; we shall denote it $\hat{\theta}_{c+1}$. (The appearance of 1 in the subscript $\hat{\theta}_{c+1}$ will soon become clear.) When $c = 0$, we see that $\hat{\theta}_1$ becomes the MLE $\hat{\theta}$. The parameter c in (8) can thus be used to trade-off robustness for asymptotic efficiency as observed in [6], [7].

⁴This procedure may be viewed as a generalization of the self-weighting procedure suggested by Windham [6, p. 604].

⁵To see why the objective function in (8) is called *mean power likelihood*, verify that (7) is equivalent to

$$\frac{1}{n} \sum_{i=1}^n s_c(x_i; \theta) = 0$$

where

$$s_c(x; \theta) := P_\theta(x)^c \left[s(x; \theta) - \frac{1}{1+c} \nabla_\theta \left(\log \sum_{x \in \mathbb{X}} P_\theta(x)^{c+1} \right) \right].$$

The quantity $s_c(x_i; \theta)$ is a generalization of the power-weighted and centered score function. The centering ensures Fisher consistency. As $c \downarrow 0$, we have $s_c(x; \theta) \rightarrow s(x; \theta)$.

Let us now bring in the connection to a parametric family of relative entropies. Recall that \hat{P} is the empirical measure of the data. The argument $\theta \in \Theta$ that maximizes the objective in (8) is the same as minimizing

$$\begin{aligned} & -\frac{c+1}{c} \log \left(\frac{1}{n} \sum_{i=1}^n P_\theta(x_i)^c \right) + \frac{1}{c} \log \sum_{x \in \mathbb{X}} \hat{P}(x)^{c+1} + \log \sum_{x \in \mathbb{X}} P_\theta(x)^{c+1} \\ & = -\frac{c+1}{c} \log \sum_{x \in \mathbb{X}} \hat{P}(x) P_\theta(x)^c + \frac{1}{c} \log \sum_{x \in \mathbb{X}} \hat{P}(x)^{c+1} + \log \sum_{x \in \mathbb{X}} P_\theta(x)^{c+1} \\ & =: \mathcal{J}_{c+1}(\hat{P}, P_\theta), \end{aligned} \quad (9)$$

where \mathcal{J}_{c+1} in (9) is a parametric extension of relative entropies already studied in our companion paper [2]. We thus have

$$\hat{\theta}_{c+1} = \arg \min_{\theta \in \Theta} \mathcal{J}_{c+1}(\hat{P}, P_\theta), \quad (10)$$

and the probability measure $P_{\hat{\theta}_{c+1}}$ corresponding to the MMPLE $\hat{\theta}_{c+1}$ is called the *reverse \mathcal{J}_{c+1} -projection of the empirical measure \hat{P} on the family \mathbb{E}* . It is known (see for example [2, Lemma 1-b]) that $\lim_{c \downarrow 0} \mathcal{J}_{c+1}(P, Q) = \mathcal{J}(P \| Q)$, as it should be, for we already saw that $c = 0$ yields $\hat{\theta}_1 = \hat{\theta}$, the MLE, which is also the reverse \mathcal{J} -projection of the empirical measure \hat{P} on \mathbb{E} . This operational continuity intuitively suggests that we must have minimization in (10) and maximization in (8).

Let us now use large sample asymptotics to justify the minimization in (10) (and maximization in (8)). Let θ^* be the true parameter and let x_1, \dots, x_n be drawn independently and according to P_{θ^*} . As the number of samples n goes to infinity, almost surely, the empirical measure⁶ \hat{P} converges (point-wise) to the true probability measure P_{θ^*} . For a fixed candidate estimate θ , by virtue of the continuity of $\mathcal{J}_{c+1}(\cdot, P_\theta)$ in the first argument when $c > 0$, see [2, Prop. 2], we have (almost surely)

$$\mathcal{J}_{c+1}(\hat{P}, P_\theta) \xrightarrow{n \rightarrow \infty} \mathcal{J}_{c+1}(P_{\theta^*}, P_\theta) \geq \mathcal{J}_{c+1}(P_{\theta^*}, P_{\theta^*}),$$

where the last inequality follows from the fact that $\mathcal{J}_\alpha(P_{\theta^*}, P_\theta) \geq 0$ with equality if and only if $\theta = \theta^*$ [2, Lem. 1-a)]. From this, it is clear that one must minimize over $\theta \in \Theta$ (and not maximize) in (10) in order to identify the true parameter θ^* .

Some historical remarks are now called for. Basu et al. [4] studied a nonnormalized version of the estimating equation (7), namely (5) with $w(x; \theta) = P_\theta(x)^c$. They also identified an associated divergence which is now called β -divergence [9], [10]. The β -divergences belong to the class of Bregman divergences [11]. Jones et al. [7] proposed the normalized estimating equation (7) and identified a divergence associated with (7), see [7, Eq. (2.8)]. Fujisawa and Eguchi [9] found that \mathcal{J}_{c+1} is another divergence associated with the estimating equation (7) and termed it γ -divergence. They also established an approximate Pythagorean relation for \mathcal{J}_{c+1} (which is quite different from what we shall discuss in Section V) and used it to bound the error between estimates arising with and without contamination by outliers⁷. Recently, Cichocki and Amari [10] surveyed the properties of the β - and the \mathcal{J}_α -divergences and their connection to other divergences.

Earlier Sundaresan [12] and [13] arrived at \mathcal{J}_α -divergences in the context of redundancy in compression and guessing problems (for $\alpha < 1$). Let us now turn to this.

C. Reverse \mathcal{J}_α -projection: $\alpha < 1$

We now motivate reverse \mathcal{J}_α -projection for $\alpha < 1$. Rényi entropies play a role similar to Shannon entropy when one wishes to minimize the normalized cumulant of compressed lengths as opposed to expected compressed lengths. More precisely, with $\rho = \alpha^{-1} - 1 > 0$, Campbell [14] showed that

$$\min \frac{1}{n\rho} \log \mathbb{E}[\exp\{\rho L_n(X^n)\}] \rightarrow H_\alpha(\hat{P}) \quad (\text{as } n \rightarrow \infty)$$

for an i.i.d. source with marginal \hat{P} . The minimization is taken over all length functions L_n that satisfy the Kraft inequality. ρ is the cumulant parameter. As $\alpha \uparrow 1$, we have $\rho \downarrow 0$, and it is well known that $\lim_{\alpha \uparrow 1} H_\alpha(\hat{P}) = H(\hat{P})$, the Shannon entropy, so that Rényi entropy can be viewed as an operational generalization of Shannon entropy.

Suppose now that the compressor is forced to use for compression, not the true probability measure \hat{P} , but a probability measure P_θ from a family parameterized by $\theta \in \Theta$. Let us denote, as before, $\mathbb{E} = \{P_\theta : \theta \in \Theta\}$. As an example, \hat{P} may be a generic measure on $\mathbb{X} = \{0, 1, \dots, L\}$, but the compressor may wish to pick the best representation of \hat{P} among binomial distributions P_θ having $\theta \in (0, 1)$ as parameter⁸. If the compressor picks P_θ instead of the true \hat{P} , then the gap in the resulting normalized cumulant from the optimal value is $\mathcal{J}_\alpha(\hat{P}, P_\theta)$ [13]. It follows that the best compressor from within \mathbb{E} has parameter

$$\hat{\theta}_\alpha = \arg \min_{\theta \in \Theta} \mathcal{J}_\alpha(\hat{P}, P_\theta) \quad (11)$$

⁶The dependence of \hat{P} on n is understood and suppressed.

⁷The outliers are generated using a mixture model.

⁸More sophisticated examples are possible. Take $\mathbb{X} = \{0, 1\}^Z$, \hat{P} any fixed, stationary, and ergodic probability measure on \mathbb{X} , and \mathbb{E} the class of stationary Markov measures on \mathbb{X} of fixed Markov order. Since this \mathbb{X} is not finite, such examples are beyond the scope of this paper.

and the probability measure $P_{\hat{\theta}_\alpha}$ is the reverse \mathcal{I}_α -projection of \hat{P} on the family \mathbb{E} . While (10) defines reverse \mathcal{I}_α -projection for $\alpha > 1$, (11) defines such a projection for $\alpha < 1$. As one expects, $\lim_{\alpha \uparrow 1} \mathcal{I}_\alpha(\hat{P}, P_\theta) = \mathcal{I}(\hat{P} \| P_\theta)$, the penalty for mismatch in compression when *expected lengths* are considered, and one has the operational continuity that $\mathcal{I}(\hat{P} \| P_\theta)$ is the usual limiting penalty for mismatch as $\alpha \uparrow 1$.

\mathcal{I}_α also arises as the gap from optimality due to mismatch in performance of guessing schemes (Arikan [15], Hanawal and Sundaresan [16], Sundaresan [13]) and more recently in the performance of coding for tasks (Bunte and Lapidoth [17]).

III. THE SETTING AND CONTRIBUTIONS

In this section, we formalize the notions of projections and the families of interest. We then highlight our contributions. We begin by recalling the definition of \mathcal{I}_α and its alternate expressions.

Definition 1: The *relative α -entropy* of P with respect to Q is defined as

$$\mathcal{I}_\alpha(P, Q) := \frac{\alpha}{1-\alpha} \log \left[\sum_x P(x) Q(x)^{\alpha-1} \right] - \frac{1}{1-\alpha} \log \sum_x P(x)^\alpha + \log \sum_x Q(x)^\alpha \quad (12)$$

$$= \frac{\alpha}{1-\alpha} \log \left[\sum_x \frac{P(x)}{\|P\|} \left(\frac{Q(x)}{\|Q\|} \right)^{\alpha-1} \right], \quad (13)$$

where

$$\|Q\| = \left[\sum_x Q(x)^\alpha \right]^{1/\alpha}.$$

Equation (12) is the same as (9) but with the parameter space extended to $\alpha > 0, \alpha \neq 1$. Equation (13) follows after regrouping of terms using the definition of $\|P\|$ and $\|Q\|$. For any $\tau > 0$, since $Q/\|Q\| = \tau Q/\|\tau Q\|$, it follows that (13) can be extended to any pair of positive measures P and Q on \mathbb{X} , and not just probability measures on \mathbb{X} .

For each $\alpha > 0, \alpha \neq 1$, $\mathcal{I}_\alpha(P, Q) \geq 0$ with equality iff $P = Q$.

Note that $\mathcal{I}_\alpha(P, Q) = \infty$ if and only if either

- $\alpha < 1$ and P is not absolutely continuous with respect to Q (notation $P \not\ll Q$), or
- $\alpha > 1$ and P and Q are singular, i.e., the supports of P and Q are disjoint.

Let $\mathcal{P}(\mathbb{X})$ be the set of all probability measures on \mathbb{X} . For a probability measure P on \mathbb{X} , let $\text{Supp}(P) = \{x : P(x) > 0\}$ denote the support of P . For a set \mathbb{E} of probability measures, write $\text{Supp}(\mathbb{E})$ for the union of the supports of the members of \mathbb{E} .

Let us now formally define what we mean by a reverse \mathcal{I}_α -projection for $\alpha > 0, \alpha \neq 1$.

Definition 2 (Reverse \mathcal{I}_α -projection): Let R be a probability measure on \mathbb{X} . Let \mathbb{E} be a set of probability measures on \mathbb{X} such that $\mathcal{I}_\alpha(R, P) < \infty$ for some $P \in \mathbb{E}$. A probability measure $Q \in \mathbb{E}$ satisfying

$$\mathcal{I}_\alpha(R, Q) = \inf_{P \in \mathbb{E}} \mathcal{I}_\alpha(R, P) =: \mathcal{I}_\alpha(R, \mathbb{E}) \quad (14)$$

is called a *reverse \mathcal{I}_α -projection* of R on \mathbb{E} . If there is no such $Q \in \mathbb{E}$, a probability measure Q in the closure of \mathbb{E} satisfying (14) is called a *generalized reverse \mathcal{I}_α -projection* of R on \mathbb{E} .

In a previous paper [2], we studied the *forward \mathcal{I}_α -projection* of a probability measure R on a family. We reproduce [2, Defn. 6] here for it plays a crucial role in this paper.

Definition 3 (Forward \mathcal{I}_α -projection): Let R be a probability measure on \mathbb{X} . Let \mathbb{E} be a set of probability measures on \mathbb{X} such that $\mathcal{I}_\alpha(P, R) < \infty$ for some $P \in \mathbb{E}$. A probability measure $Q \in \mathbb{E}$ satisfying

$$\mathcal{I}_\alpha(Q, R) = \inf_{P \in \mathbb{E}} \mathcal{I}_\alpha(P, R) =: \mathcal{I}_\alpha(\mathbb{E}, R) \quad (15)$$

is called a *forward \mathcal{I}_α -projection* of R on \mathbb{E} .

In Definition 2, the minimization is with respect to the second argument, while in Definition 3 the minimization is with respect to the first argument. The focus in [2] was on forward projection on convex families and general alphabet spaces. We provided sufficient conditions for existence of the forward projection and argued that if the forward projection exists then it is unique. Convex families arise naturally from constraints placed by measurements of linear statistics. Examples of such families are linear families which we now define.

Definition 4 (Linear family): A linear family characterized by k functions $f_i : \mathbb{X} \rightarrow \mathbb{R}$, $1 \leq i \leq k$, is the set of probability measures given by

$$\mathbb{L} := \left\{ P \in \mathcal{P}(\mathbb{X}) : \sum_x P(x) f_i(x) = 0, i = 1, \dots, k \right\}. \quad (16)$$

Reverse \mathcal{S}_α -projections, however, correspond to maximum likelihood or robust estimations, and are often on exponential families which we now define.

Definition 5 (Exponential family): An exponential family characterized by a probability measure R and k functions $f_i: \mathbb{X} \rightarrow \mathbb{R}$, $1 \leq i \leq k$, is the set of probability measures given by

$$\mathbb{M} := \{P_\theta: \theta \in \Theta \subset \mathbb{R}^k\},$$

where

$$\begin{aligned} P_\theta(x)^{-1} &:= Z(\theta) \exp \left[\log (R(x)^{-1}) + \sum_{i=1}^k \theta_i f_i(x) \right] \\ &= Z(\theta) R(x)^{-1} \exp \left[\sum_{i=1}^k \theta_i f_i(x) \right] \quad \forall x \in \mathbb{X} \end{aligned}$$

with $Z(\theta)$ being the normalization constant and Θ being the subset of \mathbb{R}^k for which P_θ is a valid probability measure⁹.

Examples of exponential families include

- Bernoulli distribution ($\mathbb{X} = \{0, 1\}$, $\Theta = (0, 1)$),
- Binomial distribution ($\mathbb{X} = \{0, 1, \dots, L\}$, $\Theta = (0, 1)$),
- Poisson distribution ($\mathbb{X} = \{0, 1, \dots\}$, $\Theta = (0, \infty)$), and
- Gaussian distribution ($\mathbb{X} = \mathbb{R}^d$, the parameter θ denotes the pair of mean and covariance).

The last two are given only as illustrative examples for they do not satisfy the finite \mathbb{X} assumption of this paper. We will take up the study of reverse \mathcal{S}_α -projection on the more general *log-convex* families which we now define.

Definition 6 (Log-convex family): A set \mathbb{E} of probability measures on a finite alphabet set \mathbb{X} is said to be *log-convex* if for any two probability measures P and Q in \mathbb{E} that are not singular, and any $t \in [0, 1]$, the probability measure $\overline{P^t Q^{1-t}}$ defined by

$$\overline{P^t Q^{1-t}}(x) := \frac{P(x)^t Q(x)^{1-t}}{\sum_y P(y)^t Q(y)^{1-t}} \quad (17)$$

also belongs to \mathbb{E} .

Exponential families are log-convex, a fact that is easily checked.

We will also take up reverse projections on analogs of exponential families. To define these analogs, let us first define the generalized logarithm and the generalized exponential functions [18]. Let $\overline{\mathbb{R}}_+ = \mathbb{R} \cup \{+\infty\}$ and let $\overline{\mathbb{R}} = \mathbb{R} \cup \{+\infty, -\infty\}$.

Definition 7: For $\alpha > 0$, the α -logarithm function, denoted $\ln_\alpha: \overline{\mathbb{R}}_+ \rightarrow \overline{\mathbb{R}}$, is defined to be

$$\ln_\alpha(u) := \begin{cases} \frac{u^{1-\alpha}-1}{1-\alpha} & \alpha \neq 1 \\ \log(u) & \alpha = 1 \end{cases}$$

where the log function is the natural logarithm. Its functional inverse, the α -exponential function, denoted $e_\alpha: \overline{\mathbb{R}} \rightarrow \overline{\mathbb{R}}_+$, is defined to be

$$e_\alpha(u) := \begin{cases} (\max\{1 + (1 - \alpha)u, 0\})^{1/(1-\alpha)} & \alpha \neq 1 \\ \exp(u) & \alpha = 1. \end{cases}$$

It is easy to check that $e_\alpha(\ln_\alpha(u)) = u$ for $u > 0$ and that $\ln_\alpha(e_\alpha(u)) = u$ whenever $0 < e_\alpha(u) < \infty$.

The analogs of exponential families are the so-called α -power-law families which we now define. (Compare Definitions 5 and 8.)

Definition 8 (α -power-law family): Let R be a probability measure such that if $\alpha > 1$ then $\text{Supp}(R) = \mathbb{X}$. An α -power-law family characterized by the probability measure R and k functions $f_i: \mathbb{X} \rightarrow \mathbb{R}$, $1 \leq i \leq k$, is the set of probability measures given by

$$\mathbb{M}^{(\alpha)} := \{P_\theta: \theta \in \Theta \subset \mathbb{R}^k\},$$

where

$$P_\theta(x)^{-1} := Z(\theta) e_\alpha \left[\ln_\alpha (R(x)^{-1}) + \sum_{i=1}^k \theta_i f_i(x) \right] \quad \forall x \in \mathbb{X}, \quad (18)$$

⁹If $R(x)$ equals 0, then so does $P_\theta(x)$.

provided

$$1 + (1 - \alpha) \left[\ln_\alpha (R(x)^{-1}) + \sum_{i=1}^k \theta_i f_i(x) \right] > 0 \quad \forall x \in \mathbb{X},$$

with $Z(\theta)$ being the normalization constant and Θ being the subset of \mathbb{R}^k for which P_θ is a valid probability measure. Equivalently¹⁰,

$$P_\theta(x)^{\alpha-1} = Z(\theta)^{1-\alpha} \left[R(x)^{\alpha-1} + (1 - \alpha) \sum_{i=1}^k \theta_i f_i(x) \right] > 0 \quad \forall x \in \mathbb{X}. \quad (19)$$

When we wish to be explicit about the characterizing entities, we shall write $\mathbb{M}^{(\alpha)}(R, f_1, \dots, f_k)$ for the family. In Appendix 22, we show that $\mathbb{M}^{(\alpha)}$ depends on R in only a weak manner. Any member $P_{\theta^*} \in \mathbb{M}^{(\alpha)}$ may equally well play the role of R and this merely corresponds to translation and scaling of the parameter space.

$\mathbb{M}^{(\alpha)}$ is not closed. Sometimes it will be required to consider its closure $\text{cl}(\mathbb{M}^{(\alpha)})$.

One has the more general notion of \ln_α -convex family as well (see van Erven and Harremoës [19]¹¹).

Definition 9 (\ln_α -convex family): A set \mathbb{E} of probability measures is said to be \ln_α -convex if for any two probability measures P and Q in \mathbb{E} (that are not singular when $\alpha \leq 1$), and any $t \in [0, 1]$, the probability measure R defined by

$$R^{-1} := Z e_\alpha (t \ln_\alpha(P^{-1}) + (1 - t) \ln_\alpha(Q^{-1})) \quad (20)$$

also belongs to \mathbb{E} . The quantity Z is the normalization constant that makes R a probability measure.

Substitution of the definitions of e_α and \ln_α indicate that the probability measure R defined in (20) can be rewritten as

$$Z^{-1} [tP^{\alpha-1} + (1 - t)Q^{\alpha-1}]^{\frac{1}{\alpha-1}}. \quad (21)$$

When $\alpha = 1$, \ln_α -convexity is just log-convexity, thereby justifying that \ln_α -convexity is an extension of log-convexity. Just as exponential families are log-convex, α -power-law families are \ln_α -convex, a fact that can be easily checked using (21).

While forward projections of interest are on convex families, reverse projections of interest, particularly those arising in estimation problems, are on log-convex, and by analogy, on \ln_α -convex families. Log-convex or \ln_α -convex families are not necessarily convex in the usual sense.

Definition 9 is given only to complete the picture. We shall restrict attention in this paper to the α -power-law family.

A. A closer look at our contributions.

For a given R and a given \mathbb{E} with some P such that $\mathcal{J}_\alpha(R, P) < \infty$, we obviously have $\mathcal{J}_\alpha(R, \mathbb{E}) < \infty$. If we consider a sequence $(P_n) \subset \mathbb{E}$ such that $\lim_{n \rightarrow \infty} \mathcal{J}_\alpha(R, P_n) = \mathcal{J}_\alpha(R, \mathbb{E})$, by virtue of the continuity of $\mathcal{J}_\alpha(P, \cdot)$ in the second argument (see [2, Rem. 5]), all subsequential limits of (P_n) are generalized reverse \mathcal{J}_α -projections. In this paper, we study example settings when the generalized reverse \mathcal{J}_α -projection is unique, when it is not, and how one may characterize it, sometimes, as a forward \mathcal{J}_α -projection. Specifically, we do the following.

- In Section IV, we study reverse \mathcal{J}_α -projections on log-convex families. We show an example of nonuniqueness of generalized reverse \mathcal{J}_α -projections on an exponential family when $\alpha > 1$. However uniqueness holds for $\alpha < 1$.
- In Section V, our focus will be on the forward \mathcal{J}_α -projection on certain convex families, in particular, linear families. We identify the form of the forward \mathcal{J}_α -projection on a linear family \mathbb{L} and prove a necessary and sufficient condition for a $Q \in \mathbb{L}$ to be the forward \mathcal{J}_α -projection on \mathbb{L} . We consider the cases $\alpha > 1$ and $\alpha < 1$ separately in two subsections. The proof for the $\alpha < 1$ case is similar to Csiszár and Shields' proof for $\alpha = 1$ case [20]. For the proof of the $\alpha > 1$ case, we resort to the Lagrange multiplier technique. The structure of the forward \mathcal{J}_α -projection naturally suggests a statistical model, namely the α -power-law family $\mathbb{M}^{(\alpha)}$.
- In Section VI, we study reverse \mathcal{J}_α -projections on $\mathbb{M}^{(\alpha)}$, and show uniqueness of the generalized reverse projection for all $\alpha > 0, \alpha \neq 1$. To show this, we establish an orthogonality relationship between $\mathbb{M}^{(\alpha)}$ and an associated linear family. We then use this geometric property to turn a reverse \mathcal{J}_α -projection on $\mathbb{M}^{(\alpha)}$ into a forward \mathcal{J}_α -projection on the linear family. It will turn out that, sometimes, we may need to consider a larger family than just $\text{cl}(\mathbb{M}^{(\alpha)})$.

¹⁰A definition such as (18) is fraught with pesky issues of well-definedness. We have verified the equivalence of (19). But a skeptical reader may simply take (19) as the starting point to define $\mathbb{M}^{(\alpha)}$. The definition in (18) is given only to highlight its similarity with Definition 5. Observe that, from (19), if $\alpha < 1$, $R(x) = 0$ implies $P_\theta(x) = 0$.

¹¹van Erven and Harremoës [19] gave a different name to what we call \ln_α -convex family; they called this $(\alpha - 1)$ -convex family. Our convention follows the notation for and parametrization of the generalized logarithm.

IV. REVERSE PROJECTION ONTO LOG-CONVEX SETS

We consider the cases $\alpha > 1$ and $\alpha < 1$ separately in the next two subsections. Before that, we present a lemma of some independent interest. This is an extension of a result for relative entropy ($\alpha = 1$); see Csiszár and Matúš [21, Eq. (3)], where (22) below is an equality.

Lemma 10: Let P and Q be probability measures on \mathbb{X} that are mutually absolutely continuous. Let R be any probability measure on \mathbb{X} that is not singular with respect to P or Q . Let $t \in [0, 1]$.

(a) If $\alpha < 1$, then

$$t\mathcal{J}_\alpha(R, P) + (1-t)\mathcal{J}_\alpha(R, Q) \geq \mathcal{J}_\alpha(R, \overline{P^t Q^{1-t}}) - \log \sum_x P'(x)^t Q'(x)^{1-t}, \quad (22)$$

where P' is the escort probability measure associated with P given by

$$P'(x) := \frac{P(x)^\alpha}{\sum_y P(y)^\alpha}$$

and Q' is the escort probability measure associated with Q .

(b) If $\alpha > 1$, the inequality in (22) is reversed.

Proof: Let us first observe that if $\alpha < 1$ and $R \not\ll \overline{P^t Q^{1-t}}$, then, by the assumption that P and Q are mutually absolutely continuous, both sides of (22) are $+\infty$, and so (22) holds. We may thus assume that $R \ll \overline{P^t Q^{1-t}}$ when $\alpha < 1$. Also, notice that the hypotheses imply that R is not singular with respect to $\overline{P^t Q^{1-t}}$. Hence, for both $\alpha < 1$ and $\alpha > 1$, we may take all the terms in (22) to be finite.

Let us write

$$\frac{P(x)^t Q(x)^{1-t}}{\sum_y P(y)^t Q(y)^{1-t}} = \frac{\left(\frac{P(x)}{\|P\|}\right)^t \left(\frac{Q(x)}{\|Q\|}\right)^{1-t}}{\sum_y \left(\frac{P(y)}{\|P\|}\right)^t \left(\frac{Q(y)}{\|Q\|}\right)^{1-t}}.$$

Using this in (13) we get

$$\begin{aligned} \mathcal{J}_\alpha(R, \overline{P^t Q^{1-t}}) &= \frac{\alpha}{1-\alpha} \log \sum_x \frac{R(x)}{\|R\|} \left(\frac{\left(\frac{P(x)}{\|P\|}\right)^t \left(\frac{Q(x)}{\|Q\|}\right)^{1-t}}{\left(\sum_y \left(\frac{P(y)}{\|P\|}\right)^\alpha \left(\frac{Q(y)}{\|Q\|}\right)^{\alpha(1-t)}\right)^{\frac{1}{\alpha}}} \right)^{\alpha-1} \\ &= \frac{\alpha}{1-\alpha} \log \sum_x \frac{R(x)}{\|R\|} \left[\left(\frac{P(x)}{\|P\|}\right)^t \left(\frac{Q(x)}{\|Q\|}\right)^{1-t} \right]^{\alpha-1} + \log \sum_x \left(\frac{P(x)}{\|P\|}\right)^{\alpha t} \left(\frac{Q(x)}{\|Q\|}\right)^{\alpha(1-t)} \\ &= \frac{\alpha}{1-\alpha} \log \sum_x \left[\frac{R(x)}{\|R\|} \left(\frac{P(x)}{\|P\|}\right)^{\alpha-1} \right]^t \left[\frac{R(x)}{\|R\|} \left(\frac{Q(x)}{\|Q\|}\right)^{\alpha-1} \right]^{1-t} + \log \sum_x P'(x)^t Q'(x)^{1-t} \\ &\leq \frac{\alpha}{1-\alpha} \log \left[\sum_x \frac{R(x)}{\|R\|} \left(\frac{P(x)}{\|P\|}\right)^{\alpha-1} \right]^t \left[\sum_x \frac{R(x)}{\|R\|} \left(\frac{Q(x)}{\|Q\|}\right)^{\alpha-1} \right]^{1-t} + \log \sum_x P'(x)^t Q'(x)^{1-t} \\ &= t\mathcal{J}_\alpha(R, P) + (1-t)\mathcal{J}_\alpha(R, Q) + \log \sum_x P'(x)^t Q'(x)^{1-t}, \end{aligned}$$

for $\alpha < 1$, where the penultimate inequality follows by applying Hölder's inequality to the inner-product within the first logarithm term, with exponents $1/t$ and $1/(1-t)$. For $\alpha > 1$, the inequality is obviously reversed because the multiplication factor $\alpha/(1-\alpha)$ is negative. \blacksquare

A. Reverse \mathcal{J}_α -projection for $\alpha > 1$

Recall that the MMPLE on a log-convex family is the reverse \mathcal{J}_α -projection of the empirical measure on the family for the case when $\alpha > 1$. For log-convex families, it is possible that multiple reverse \mathcal{J}_α -projections may exist, and we provide an explicit example.

Example 1: Let $\mathbb{X} = \{0, 1, 2\}$, let R be the uniform probability measure on \mathbb{X} , and let \mathbb{E} be the log-convex family of binomial distributions on \mathbb{X} with parameter $\theta \in (0, 1)$. A member P_θ of the family is given by

$$P_\theta(0) = (1-\theta)^2, \quad P_\theta(1) = 2\theta(1-\theta), \quad P_\theta(2) = \theta^2.$$

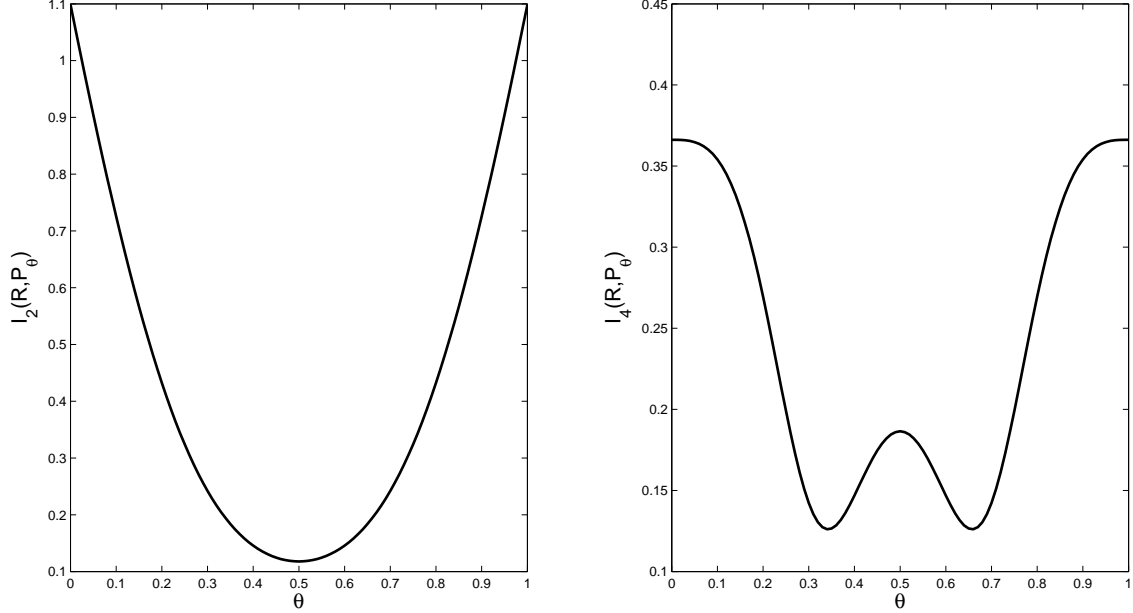


Fig. 1. Multiple reverse \mathcal{I}_α -projections are possible when $\alpha > 1$.

Figure 1 plots $\mathcal{I}_\alpha(R, P_\theta)$ as a function of θ for $\alpha = 2$ (plot on the left-hand side) and $\alpha = 4$ (plot on the right-hand side). Since $\mathcal{I}_\alpha(R, P_\theta)$ has mirror-symmetry around the point $\theta = 1/2$, a fact that can be easily checked, if there is a global minimum at $\theta^* \in (0, \frac{1}{2})$, then we have another global minimum at $1 - \theta^* \in (\frac{1}{2}, 1)$. This is the situation with the plot on the right-hand side.

Eguchi and Kato [8] consider the problem of *spontaneous* clustering for a Gaussian mixture model with an unknown number of components, and put the possibility of multiple minima to good use. Very briefly, their procedure operates on the data as follows, and we refer the interested reader to [8] for further details. They first choose the parameter α with some care using either the maximum range of the data or the Akaike information criterion. They then identify the resulting minima of $\mathcal{I}_\alpha(R, P_\theta)$ over the parameters $\theta \in \Theta$. Here R is the empirical measure¹² of the data and α is as chosen. They interpret each minimum point as the parameter of a “discovered” component of the mixture. Finally, they associate each data point to a nearby component, among those discovered, thereby arriving at a clustering. If the number of components is unknown, the number of minima is a *spontaneous* choice for the number of components of the mixture.

Example 1 suggests a sequence $(P_n) \subset \mathbb{E}$ that satisfies $\mathcal{I}_\alpha(R, P_n) \rightarrow \mathcal{I}_\alpha(R, \mathbb{E})$, and yet P_n does not converge: take $\alpha = 4$, $P_n = P_{\theta^*}$ for odd n , and $P_n = P_{1-\theta^*}$ for even n . All subsequential limits are of course generalized reverse \mathcal{I}_α -projections.

B. Reverse \mathcal{I}_α -projection for $\alpha < 1$

For $\alpha < 1$, the generalized reverse \mathcal{I}_α -projection is unique, unlike the situation in the previous subsection.

Theorem 11: Let $\alpha < 1$. Let \mathbb{E} be a log-convex set of mutually absolutely continuous probability measures on \mathbb{X} . Let R be a probability measure on \mathbb{X} such that $\mathcal{I}_\alpha(R, \mathbb{E}) < \infty$. Under these conditions, there exists a unique probability measure Q such that, for every sequence (P_n) in \mathbb{E} satisfying $\mathcal{I}_\alpha(R, P_n) \rightarrow \mathcal{I}_\alpha(R, \mathbb{E})$, we have $P_n \rightarrow Q$ and $\mathcal{I}_\alpha(R, Q) = \mathcal{I}_\alpha(R, \mathbb{E})$.

Proof: The proof broadly follows the proof of Csiszár’s [21, Th. 1].

Consider a sequence $(P_n) \subset \mathbb{E}$ such that $\lim_n \mathcal{I}_\alpha(R, P_n) = \mathcal{I}_\alpha(R, \mathbb{E})$. Since $\mathcal{I}_\alpha(R, \mathbb{E})$ is finite, we may assume without loss of generality that $\mathcal{I}_\alpha(R, P_n)$ is finite for all n . Hence, for all n , R is not singular with respect to P_n ; indeed, $R \ll P_n$

¹²The empirical measure R and the Gaussian P_θ are singular. Following the formal definition in [2, Sec. II], strictly speaking, we have the relative α -entropy $\mathcal{I}_\alpha(R, P_\theta) = \infty$. The expansion however does provide a valid expression for optimization although one cannot interpret it as the relative α -entropy, and Eguchi and Kato [8] minimize the expression to get the MMPLE.

for all n . Apply Lemma 10 with $P = P_m$, $Q = P_n$ to get

$$t\mathcal{J}_\alpha(R, P_m) + (1-t)\mathcal{J}_\alpha(R, P_n) \geq \mathcal{J}_\alpha(R, \overline{P_m^t P_n^{1-t}}) - \log \sum_x P_m'(x)^t P_n'(x)^{1-t} \quad (23)$$

$$\geq \mathcal{J}_\alpha(R, \mathbb{E}) - \log \sum_x P_m'(x)^t P_n'(x)^{1-t}, \quad (24)$$

where last inequality follows from the hypothesis that $\overline{P_m^t P_n^{1-t}} \in \mathbb{E}$. Also observe that, by Hölder's inequality,

$$\sum_x P_m'(x)^t P_n'(x)^{1-t} \leq \left(\sum_x P_m'(x) \right)^t \left(\sum_x P_n'(x) \right)^{1-t} = 1. \quad (25)$$

Let $m, n \rightarrow \infty$ in (24) and use (25) to get

$$\lim_{m, n \rightarrow \infty} \log \sum_x P_m'(x)^t P_n'(x)^{1-t} = 0.$$

Set $t = 1/2$ in this limit and undo the logarithm to get

$$\lim_{m, n \rightarrow \infty} \sum_x \sqrt{P_m'(x) P_n'(x)} = 1$$

so that

$$\begin{aligned} \sum_x \left(\sqrt{P_m'(x)} - \sqrt{P_n'(x)} \right)^2 &= 2 - 2 \cdot \sum_x \sqrt{P_m'(x) P_n'(x)} \\ &\rightarrow 0 \quad \text{as } m, n \rightarrow \infty. \end{aligned}$$

Thus (P_n') is a Cauchy sequence. It must converge to some Q' , an escort of some probability measure Q . Given our finite alphabet assumption, we must then have $P_n \rightarrow Q$.

If $(Q_n) \subset \mathbb{E}$ is another sequence such that $\mathcal{J}_\alpha(R, Q_n) \rightarrow \mathcal{J}_\alpha(R, \mathbb{E})$, then since P_n and Q_n can be merged together, (Q_n) must also converge to the same Q . The generalized reverse \mathcal{J}_α -projection is therefore unique.

By continuity of $\mathcal{J}_\alpha(R, \cdot)$, see [2, Rem. 5], we also have $\mathcal{J}_\alpha(R, Q) = \mathcal{J}_\alpha(R, \mathbb{E})$. ■

The proof fails for $\alpha > 1$ because the inequality in (24) is in the opposite direction, and one cannot conclude that (P_n') is a Cauchy sequence. Indeed, the previous subsection provides a counterexample for lack of convergence and nonuniqueness of reverse \mathcal{J}_α -projection on a log-convex family, when $\alpha > 1$.

V. FORWARD \mathcal{J}_α -PROJECTION

In this section, we will recall some results on forward \mathcal{J}_α -projection from [2] along with some refinements for our restricted finite alphabet setting. The proofs here use elementary tools and exploit the finite alphabet assumption. The results will then be used to turn a reverse \mathcal{J}_α -projection on an α -power-law family into a forward \mathcal{J}_α -projection on a linear family.

A. $\alpha < 1$:

The result for $\alpha < 1$ is the following. It establishes the form of the forward \mathcal{J}_α -projection on a linear family.

Theorem 12: Let $\alpha < 1$. Let \mathbb{L} be a linear family characterized by $f_i, i = 1, \dots, k$. Let R be a probability measure with full support. Then the following hold.

(a) R has a forward \mathcal{J}_α -projection on \mathbb{L} . Call it Q .

(b) $\text{Supp}(Q) = \text{Supp}(\mathbb{L})$ and the Pythagorean equality holds (see Figure 2):

$$\mathcal{J}_\alpha(P, R) = \mathcal{J}_\alpha(P, Q) + \mathcal{J}_\alpha(Q, R) \quad \forall P \in \mathbb{L}. \quad (26)$$

(c) The forward \mathcal{J}_α -projection Q satisfies

$$Z^{\alpha-1} Q(x)^{\alpha-1} = R(x)^{\alpha-1} + (1-\alpha) \sum_{i=1}^k \theta_i^* f_i(x) \quad \forall x \in \text{Supp}(\mathbb{L}), \quad (27)$$

where $\theta_1^*, \dots, \theta_k^*$ are scalars and Z is the normalization constant that makes Q a probability measure.

(d) The forward \mathcal{J}_α -projection is unique.

Proof: (a) The mapping $P \mapsto \mathcal{J}_\alpha(P, R)$ is continuous [2, Rem. 5] and \mathbb{L} is compact. Hence the forward \mathcal{J}_α -projection exists.

(b) This follows from [2, Props. 14-15, Th. 10-a].

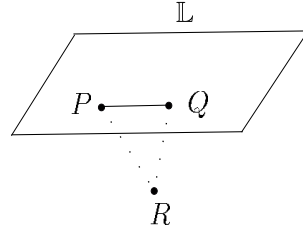


Fig. 2. Pythagorean property

(c) Our proof follows the proof of Csiszár and Shields proof for the case $\alpha = 1$ [20, Th. 3.2].

From (16), it is clear that the probability measures $P \in \mathbb{L}$, when considered as $|\text{Supp}(\mathbb{L})|$ -dimensional vectors, belong to the orthogonal complement \mathcal{F}^\perp of the subspace \mathcal{F} of $\mathbb{R}^{|\text{Supp}(\mathbb{L})|}$ spanned by the vectors $f_i(\cdot), i = 1, \dots, k$, restricted to $\text{Supp}(\mathbb{L})$. These $P \in \mathbb{L}$ actually span \mathcal{F}^\perp . (This follows from the fact that if a subspace of $\mathbb{R}^{|\text{Supp}(\mathbb{L})|}$ contains a vector all of whose components are strictly positive, here Q , then it is spanned by the probability vectors of that space.) Using (13), one can see (26) same as

$$\sum_x P(x) \left(\frac{R(x)^{\alpha-1}}{\sum_a Q(a)R(a)^{\alpha-1}} - \frac{Q(x)^{\alpha-1}}{\sum_a Q(a)^\alpha} \right) = 0 \quad \forall P \in \mathbb{L}.$$

Consequently, the vector

$$\frac{R(\cdot)^{\alpha-1}}{\sum_a Q(a)R(a)^{\alpha-1}} - \frac{Q(\cdot)^{\alpha-1}}{\sum_a Q(a)^\alpha}$$

belongs to $(\mathcal{F}^\perp)^\perp = \mathcal{F}$, that is,

$$\frac{R(x)^{\alpha-1}}{\sum_a Q(a)R(a)^{\alpha-1}} - \frac{Q(x)^{\alpha-1}}{\sum_a Q(a)^\alpha} = \sum_{i=1}^k \lambda_i f_i(x) \quad \forall x \in \text{Supp}(\mathbb{L})$$

for some scalars $\lambda_i, i = 1, \dots, k$. This verifies (27) for obvious choices of Z and θ_i^* .

(d) This follows from [2, Th. 8]. ■

One can also state a converse.

Theorem 13: Let $\alpha < 1$. Let $Q \in \mathbb{L}$ be a probability measure of the form (27). Then Q satisfies (26) and is the forward \mathcal{I}_α -projection of R on \mathbb{L} .

Proof: This follows from [2, Th. 11-b]. ■

B. $\alpha > 1$:

We now establish the form of the forward \mathcal{I}_α -projection on a linear family when $\alpha > 1$. The following result may be seen as a refinement of [2, Th. 10(a)].

Theorem 14: Let $\alpha > 1$. Let \mathbb{L} be a linear family characterized by $f_i, i = 1, \dots, k$. Let R be a probability measure with full support. Then the following hold.

- (a) R has a forward \mathcal{I}_α -projection on \mathbb{L} . Call it Q .
- (b) The forward \mathcal{I}_α -projection Q satisfies

$$Z^{\alpha-1}Q(x)^{\alpha-1} = \left[R(x)^{\alpha-1} + (1-\alpha) \sum_{i=1}^k \theta_i^* f_i(x) \right]_+ \quad \forall x \in \mathbb{X}, \quad (28)$$

where $\theta_1^*, \dots, \theta_k^*$ are scalars, Z is the normalization constant that makes Q a probability measure, and $[u]_+ = \max\{u, 0\}$.

- (c) The Pythagorean inequality holds:

$$\mathcal{I}_\alpha(P, R) \geq \mathcal{I}_\alpha(P, Q) + \mathcal{I}_\alpha(Q, R) \quad \forall P \in \mathbb{L}. \quad (29)$$

- (d) The forward \mathcal{I}_α -projection is unique.
- (e) If $\text{Supp}(Q) = \text{Supp}(\mathbb{L})$, then (29) holds with equality.

Proof: (a) The mapping $P \mapsto \mathcal{I}_\alpha(P, R)$ is continuous [2, Prop. 2] and \mathbb{L} is compact. Hence the forward \mathcal{I}_α -projection exists.

(b) The optimization problem for the forward \mathcal{J}_α -projection is

$$\min_P \mathcal{J}_\alpha(P, R) \quad (30)$$

$$\text{subject to } \sum_x P(x) f_i(x) = 0, \quad i = 1, \dots, k \quad (31)$$

$$\sum_x P(x) = 1 \quad (32)$$

$$P(x) \geq 0 \quad \forall x \in \mathbb{X}. \quad (33)$$

We will proceed in a sequence of steps.

(i) Observe that $\mathcal{J}_\alpha(\cdot, R)$, in addition to being continuous, is also continuously differentiable. Indeed, we have

$$\frac{\partial}{\partial P(x)} \mathcal{J}_\alpha(P, R) = \frac{\alpha}{1-\alpha} \left[\frac{R(x)^{\alpha-1}}{\sum_a P(a) R(a)^{\alpha-1}} - \frac{P(x)^{\alpha-1}}{\sum_a P(a)^\alpha} \right]. \quad (34)$$

Both denominators are bounded away from zero because for any $P \in \mathbb{L}$, we have $\max_x P(x) \geq 1/|\mathbb{X}|$, and therefore

$$\sum_a P(a) R(a)^{\alpha-1} \geq \frac{1}{|\mathbb{X}|} \cdot \min_a R(a)^{\alpha-1} > 0,$$

and

$$\sum_a P(a)^\alpha \geq \frac{1}{|\mathbb{X}|^\alpha} > 0.$$

Consequently, the partial derivative (34) exists everywhere on $\mathbb{R}_+^{|\mathbb{X}|}$, and is continuous because the terms involved are continuous. (The numerator of the second term in (34) is continuous because $\alpha > 1$).

(ii) Since the equality constraints in (31) and (32) arise from affine functions, and the inequality constraints in (33) arise from linear functions, we may apply [22, Prop. 3.3.7] to conclude that there exist Lagrange multipliers $(\lambda_i, i = 1, \dots, k)$, ν , and $(\mu(x), x \in \mathbb{X})$ associated with the constraints (31), (32), and (33), respectively, that satisfy:

$$\frac{\alpha}{1-\alpha} \left[\frac{Q(x)^{\alpha-1}}{\sum_a Q(a)^\alpha} - \frac{R(x)^{\alpha-1}}{\sum_a Q(a) R(a)^{\alpha-1}} \right] = \sum_{i=1}^k \lambda_i f_i(x) - \mu(x) + \nu \quad \forall x \quad (35)$$

$$\mu(x) \geq 0 \quad \forall x \quad (36)$$

$$\mu(x) Q(x) = 0 \quad \forall x. \quad (37)$$

In writing (35), we have substituted (34) for $\frac{\partial}{\partial P(x)} \mathcal{J}_\alpha(P, R)$.

(iii) Multiplying (35) by $Q(x)$, summing over all $x \in \mathbb{X}$, using $Q \in \mathbb{L}$, and using (37), we see that $\nu = 0$.

(iv) If $Q(x) > 0$, we must have $\mu(x) = 0$ from (37), and its substitution in (35) yields, for all such x ,

$$\frac{Q(x)^{\alpha-1}}{\sum_a Q(a)^\alpha} = \frac{R(x)^{\alpha-1}}{\sum_a Q(a) R(a)^{\alpha-1}} + \frac{1-\alpha}{\alpha} \sum_{i=1}^k \lambda_i f_i(x). \quad (38)$$

If $Q(x) = 0$, (35) implies that

$$\frac{R(x)^{\alpha-1}}{\sum_a Q(a) R(a)^{\alpha-1}} + \frac{1-\alpha}{\alpha} \sum_{i=1}^k \lambda_i f_i(x) = \frac{(1-\alpha)}{\alpha} \mu(x) \leq 0, \quad (39)$$

where the last inequality holds because of (36) and $\alpha > 1$. Therefore, (38) and (39) may be combined as

$$Z^{\alpha-1} Q(x)^{\alpha-1} = \left[R(x)^{\alpha-1} + (1-\alpha) \sum_{i=1}^k \theta_i^* f_i(x) \right]_+ \quad \forall x \in \mathbb{X},$$

where the choices of Z and θ_i^* are obvious. This verifies (28) and completes the proof of (b).

(c) This follows from [2, Th. 10-a].

(d) Follows from [2, Th. 8].

(e) This can be shown using the proof of [2, Prop. 15] and using [2, Th. 10-a]. ■

As in the $\alpha < 1$ case, one has a converse.

Theorem 15: Let $\alpha > 1$. Let $Q \in \mathbb{L}$ be a probability measure of the form (28). Then Q satisfies (29) for every $P \in \mathbb{L}$, and Q is the forward \mathcal{J}_α -projection of R on \mathbb{L} .

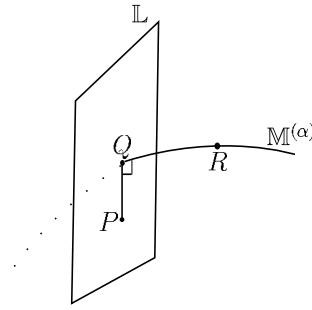


Fig. 3. Orthogonal intersection of an α -power-law family and a linear family

Proof: Follows from [2, Th. 11-b]. ■

When $\alpha > 1$, in general, $\text{Supp}(Q) \neq \text{Supp}(\mathbb{L})$ as shown by the following counterexample, and the Pythagorean inequality (29) may be strict.

Example 2: Let $\alpha = 2$. Let $\mathbb{X} = \{1, 2, 3, 4\}$. Write $P = (p_1, p_2, p_3, p_4)$ for a probability measure on \mathbb{X} . Define the linear family \mathbb{L} to be

$$\mathbb{L} = \{P \in \mathcal{P}(\mathbb{X}) : 8p_1 + 4p_2 + 2p_3 + p_4 = 7\}.$$

Let R be the uniform probability measure on \mathbb{X} . We claim that the forward \mathcal{S}_α -projection of R on \mathbb{L} is $Q = (3/4, 1/4, 0, 0)$. First, $Q \in \mathbb{L}$ because $8q_1 + 4q_2 + 2q_3 + q_4 = 8 \times 3/4 + 4 \times 1/4 + 0 + 0 = 7$. Second, Q is of the form (28). To see this, let us note that $f_1(\cdot) = (1, -3, -5, -6)$. Take $\theta_1^* = -1/20$ and $Z = 2/5$. Then

$$\begin{aligned} [R(\cdot)^{\alpha-1} + (1-\alpha)\theta_1^* f_1(\cdot)]_+ &= [R(\cdot) - \theta_1^* f_1(\cdot)]_+ \\ &= ([1/4 + 1/20]_+, [1/4 - 3/20]_+, [1/4 - 5/20]_+, [1/4 - 6/20]_+) \\ &= (6/20, 2/20, 0, 0) \\ &= Z \cdot Q(\cdot). \end{aligned}$$

That Q is the forward \mathcal{S}_α -projection now follows from Theorem 15.

Clearly $\text{Supp}(Q) \subsetneq \text{Supp}(\mathbb{L})$. Also for $P = (0.8227, 0.0625, 0.0536, 0.0612) \in \mathbb{L}$, numerical calculations yield a strict inequality in (29) since the left-hand side and the right-hand side of (29) evaluate to 1.0114 and 0.9871, respectively. See also [2, Rem. 13] where this counterexample showed that transitivity of projections does not hold for $\alpha > 1$. In both situations, the issue is that $\text{Supp}(Q) \neq \text{Supp}(\mathbb{L})$.

VI. ORTHOGONALITY BETWEEN THE α -POWER-LAW FAMILY AND THE LINEAR FAMILY

The focus of this section is on the geometry of the α -power-law family with respect to its associated linear family, and its exploitation. See Figure 3. We treat the cases $\alpha < 1$ and $\alpha > 1$ separately. Theorems 18 and 21 are the main contributions.

A. $\alpha < 1$:

This case is the simpler of the two. The core result of this section, one on which the main result Theorem 18 hinges, is the following that shows that the case $\alpha < 1$ is similar to $\alpha = 1$ [20, Th. 3.2].

Theorem 16: Let $\alpha < 1$. Let \mathbb{L} be a linear family characterized by $f_i, i = 1, \dots, k$, as in (16). Let R be a probability measure with full support. Let $\mathbb{M}^{(\alpha)}$ be the α -power-law family, as in Definition 8, characterized by R and the same k functions $f_i, i = 1, \dots, k$. Let Q be the forward \mathcal{S}_α -projection of R on \mathbb{L} . Then the following hold.

- (a) $\mathbb{L} \cap \text{cl}(\mathbb{M}^{(\alpha)}) = \{Q\}$.
- (b) For every $P \in \mathbb{L}$, we have

$$\mathcal{S}_\alpha(P, R) = \mathcal{S}_\alpha(P, Q) + \mathcal{S}_\alpha(Q, R). \quad (40)$$

- (c) If $\text{Supp}(\mathbb{L}) = \mathbb{X}$, then $\mathbb{L} \cap \mathbb{M}^{(\alpha)} = \{Q\}$.

Proof: Statement (b) is the same as Theorem 12-(c). Let us observe from Theorem 12 that when $\text{Supp}(\mathbb{L}) = \mathbb{X}$, the forward \mathcal{S}_α -projection Q of R on \mathbb{L} satisfies

$$Z^{\alpha-1} Q(x)^{\alpha-1} = R(x)^{\alpha-1} + (1-\alpha) \sum_{i=1}^k \theta_i^* f_i(x) \quad \forall x \in \mathbb{X}$$

for some scalars $Z, \theta_1^*, \dots, \theta_k^*$. Hence $Q \in \mathbb{M}^{(\alpha)}$. Since Q is also in \mathbb{L} , we have $Q \in \mathbb{L} \cap \mathbb{M}^{(\alpha)}$.

Thus, in general, when $\text{Supp}(\mathbb{L})$ is not necessarily \mathbb{X} , if we can show that (i) every member of $\mathbb{L} \cap \text{cl}(\mathbb{M}^{(\alpha)})$ satisfies (40), and (ii) $\mathbb{L} \cap \text{cl}(\mathbb{M}^{(\alpha)})$ is nonempty, then, since any member satisfying (40) is also forward \mathcal{S}_α -projection and since the forward \mathcal{S}_α -projection is unique, the theorem will be established. We now proceed to show (i) and (ii).

(i) Every $\tilde{Q} \in \mathbb{L} \cap \text{cl}(\mathbb{M}^{(\alpha)})$ satisfies (40). (41)

Let $(Q_n) \subset \mathbb{M}^{(\alpha)}$ be such that $Q_n \rightarrow \tilde{Q}$. Then, for each n , there exist $\theta^{(n)} = (\theta_1^{(n)}, \dots, \theta_k^{(n)}) \in \mathbb{R}^k$ and a constant Z_n such that

$$Z_n^{\alpha-1} Q_n(x)^{\alpha-1} = R(x)^{\alpha-1} + (1-\alpha) \sum_{i=1}^k \theta_i^{(n)} f_i(x) \quad \forall x \in \mathbb{X}. \quad (42)$$

Since, for any $P \in \mathbb{L}$, we have

$$\sum_x P(x) f_i(x) = \sum_x \tilde{Q}(x) f_i(x) = 0, \quad i = 1, \dots, k,$$

by taking expectation with respect to P and \tilde{Q} on both sides of (42), we get

$$Z_n^{\alpha-1} \sum_x P(x) Q_n(x)^{\alpha-1} = \sum_x P(x) R(x)^{\alpha-1}$$

and

$$Z_n^{\alpha-1} \sum_x \tilde{Q}(x) Q_n(x)^{\alpha-1} = \sum_x \tilde{Q}(x) R(x)^{\alpha-1},$$

respectively. Using the above two equations to eliminate $Z_n^{\alpha-1}$, we get

$$\sum_x P(x) R(x)^{\alpha-1} = \frac{\sum_x \tilde{Q}(x) R(x)^{\alpha-1}}{\sum_x \tilde{Q}(x) Q_n(x)^{\alpha-1}} \sum_x P(x) Q_n(x)^{\alpha-1}.$$

Letting $n \rightarrow \infty$, and then by using (12), we get (40) with Q replaced by \tilde{Q} . This proves (i).

(ii) $\mathbb{L} \cap \text{cl}(\mathbb{M}^{(\alpha)})$ is nonempty.

Let

$$\begin{aligned} \tau_i^{(n)} &:= \frac{\frac{1}{n} \sum_x R(x) f_i(x)}{\left(1 - \frac{1}{n}\right) \sum_x Q(x) R(x)^{\alpha-1} + \frac{1}{n} \sum_x R(x)^\alpha}, \\ \tilde{f}_i(\cdot) &:= f_i(\cdot) - \tau_i^{(n)} R(\cdot)^{\alpha-1}, \quad i = 1, \dots, k, \end{aligned}$$

and define the sequence of linear families

$$\mathbb{L}_n := \left\{ P \in \mathcal{P}(\mathbb{X}) : \sum_x P(x) \tilde{f}_i(x) = 0, \quad i = 1, \dots, k \right\}.$$

The $\tau_i^{(n)}$'s are chosen so that $(1 - \frac{1}{n})Q + \frac{1}{n}R \in \mathbb{L}_n$, and so $\text{Supp}(\mathbb{L}_n) = \mathbb{X}$. Let Q_n be the forward \mathcal{S}_α -projection of R on \mathbb{L}_n . Then, by virtue of Theorem 12-(b), we have $\text{Supp}(Q_n) = \mathbb{X}$, and by virtue of Theorem 12-(c), we have

$$\begin{aligned} Z_n^{\alpha-1} Q_n(x)^{\alpha-1} &= R(x)^{\alpha-1} + (1-\alpha) \sum_{i=1}^k \theta_i^{(n)} \tilde{f}_i(x) \\ &= \left[1 - (1-\alpha) \sum_{i=1}^k \theta_i^{(n)} \tau_i^{(n)} \right] R(x)^{\alpha-1} + (1-\alpha) \sum_{i=1}^k \theta_i^{(n)} f_i(x) \quad \forall x \in \mathbb{X}. \end{aligned} \quad (43)$$

Taking expectation with respect to Q on both sides, and using $\sum_x Q(x) f_i(x) = 0, i = 1, \dots, k$, we get

$$Z_n^{\alpha-1} \sum_x Q(x) Q_n(x)^{\alpha-1} = \left[1 - (1-\alpha) \sum_{i=1}^k \theta_i^{(n)} \tau_i^{(n)} \right] \cdot \sum_x Q(x) R(x)^{\alpha-1}.$$

As the summations on either side are finite and strictly positive for each n , the term within square brackets in the above equation is also strictly positive for each n . Rescaling (43) appropriately, we see that $Q_n \in \mathbb{M}^{(\alpha)}$. Note also that $\tau_i^{(n)} \rightarrow 0$ as $n \rightarrow \infty$ for $i = 1, \dots, k$. Hence the limit of any convergent subsequence of (Q_n) belongs to $\mathbb{L} \cap \text{cl}(\mathbb{M}^{(\alpha)})$. This verifies (ii) and concludes the proof of the theorem. ■

We now argue that the family $\text{cl}(\mathbb{M}^{(\alpha)})$ and \mathbb{L} are “orthogonal” to each other, in a sense made precise in the statement of the next result.

Corollary 17: Under the hypotheses of Theorem 16, the following additional statements hold.

(a) For every $P \in \mathbb{L}$ and every $S \in \text{cl}(\mathbb{M}^{(\alpha)})$, we have

$$\mathcal{J}_\alpha(P, S) = \mathcal{J}_\alpha(P, Q) + \mathcal{J}_\alpha(Q, S). \quad (44)$$

(b) For any $S \in \text{cl}(\mathbb{M}^{(\alpha)})$, the forward \mathcal{J}_α -projection of S on \mathbb{L} is Q .

Proof: Since any member of $\mathbb{M}^{(\alpha)}$ can play the role of R by Prop. 22 (in the Appendix), and since, by Theorem 16, $\mathbb{L} \cap \text{cl}(\mathbb{M}^{(\alpha)}) = \{Q\}$, Q is the forward \mathcal{J}_α -projection of any member of $\mathbb{M}^{(\alpha)}$ on \mathbb{L} . Therefore (44) holds for every $P \in \mathbb{L}$ and every $S \in \mathbb{M}^{(\alpha)}$. Furthermore, (44) holds for the limit of any sequence of members of $\mathbb{M}^{(\alpha)}$, and hence (a) and (b) hold for members of $\text{cl}(\mathbb{M}^{(\alpha)}) \setminus \mathbb{M}^{(\alpha)}$ as well. ■

Let us now return to the compression problem discussed in Section II-C and show the connection between the reverse \mathcal{J}_α -projection on an α -power-law family and a forward \mathcal{J}_α -projection on a linear family.

Theorem 18: Let $\alpha < 1$. Let \hat{P} be a probability measure on \mathbb{X} . Let $\mathbb{M}^{(\alpha)}$ be characterized by the probability measure R and the functions $f_i, i = 1, \dots, k$. Let \mathbb{L} be the associated linear family characterized by $f_i, i = 1, \dots, k$, and assume that it is nonempty. Let R have full support.

Define $\tilde{\mathbb{L}}$ as

$$\tilde{\mathbb{L}} := \left\{ P \in \mathcal{P}(\mathbb{X}) : \sum_x P(x) \tilde{f}_i(x) = 0 \right\}, \quad (45)$$

where

$$\tilde{f}_i(\cdot) = f_i(\cdot) - \tau_i^R R(\cdot)^{\alpha-1} \quad (46)$$

with

$$\tau_i^R = \frac{\sum_x \hat{P}(x) f_i(x)}{\sum_x \hat{P}(x) R(x)^{\alpha-1}}, i = 1, \dots, k. \quad (47)$$

Let Q be the forward \mathcal{J}_α -projection of R on $\tilde{\mathbb{L}}$.

(a) If $\text{Supp}(Q) = \mathbb{X}$, then Q is the unique reverse \mathcal{J}_α -projection of \hat{P} on $\mathbb{M}^{(\alpha)}$.

(b) If $\text{Supp}(Q) \neq \mathbb{X}$, then \hat{P} does not have a reverse \mathcal{J}_α -projection on $\mathbb{M}^{(\alpha)}$. However, Q is the unique reverse \mathcal{J}_α -projection of \hat{P} on $\text{cl}(\mathbb{M}^{(\alpha)})$.

Proof: $\tilde{\mathbb{L}}$ is constructed so that $\hat{P} \in \tilde{\mathbb{L}}$ (which is easy to check) and, further, $\tilde{\mathbb{L}}$ is orthogonal to $\mathbb{M}^{(\alpha)}$ in the sense of Corollary 17. We now verify the latter statement. For concreteness, we will index the the α -power-law family by its characterizing entities. By Corollary 17, $\tilde{\mathbb{L}}$ is orthogonal to $\mathbb{M}^{(\alpha)}(R, \tilde{f}_1, \dots, \tilde{f}_k)$. It therefore suffices to show that $\mathbb{M}^{(\alpha)}(R, \tilde{f}_1, \dots, \tilde{f}_k) = \mathbb{M}^{(\alpha)}(R, f_1, \dots, f_k)$. Take any $P_\theta \in \mathbb{M}^{(\alpha)}(R, f_1, \dots, f_k)$. Then, for each $x \in \mathbb{X}$, we have

$$\begin{aligned} Z(\theta)^{\alpha-1} P_\theta(x)^{\alpha-1} &= R(x)^{\alpha-1} + (1-\alpha) \sum_{i=1}^k \theta_i f_i(x) \\ &= (1 + (1-\alpha)\theta_i \tau_i^R) R(x)^{\alpha-1} + (1-\alpha) \sum_{i=1}^k \theta_i \tilde{f}_i(x). \end{aligned}$$

Taking expectation with respect to \hat{P} on both sides, and using $\sum_x \hat{P}(x) \tilde{f}_i(x) = 0, i = 1, \dots, k$, we get

$$Z(\theta)^{\alpha-1} \sum_x \hat{P}(x) P_\theta(x)^{\alpha-1} = [1 + (1-\alpha)\theta_i \tau_i^R] \cdot \sum_x \hat{P}(x) R(x)^{\alpha-1}.$$

Since P_θ and R have full support, it follows that $[1 + (1-\alpha)\theta_i \tau_i^R] > 0$, and hence $P_\theta \in \mathbb{M}^{(\alpha)}(R, \tilde{f}_1, \dots, \tilde{f}_k)$. This shows $\mathbb{M}^{(\alpha)}(R, f_1, \dots, f_k) \subset \mathbb{M}^{(\alpha)}(R, \tilde{f}_1, \dots, \tilde{f}_k)$. Similarly, using the assumption that \mathbb{L} is nonempty, one can show that $\mathbb{M}^{(\alpha)}(R, \tilde{f}_1, \dots, \tilde{f}_k) \subset \mathbb{M}^{(\alpha)}(R, f_1, \dots, f_k)$.

By Corollary 17, we have

$$\mathcal{J}_\alpha(\hat{P}, S) = \mathcal{J}_\alpha(\hat{P}, Q) + \mathcal{J}_\alpha(Q, S) \quad \forall S \in \text{cl}(\mathbb{M}^{(\alpha)}). \quad (48)$$

(a) If $\text{Supp}(Q) = \mathbb{X}$, then by Th. 16(c), $Q \in \mathbb{M}^{(\alpha)}$, and from (48), the minimum of $\mathcal{J}_\alpha(\hat{P}, S)$ over $S \in \mathbb{M}^{(\alpha)}$ is attained at $S = Q$. To prove the uniqueness, let $P_{\theta^*} \in \mathbb{M}^{(\alpha)}$ also attain the minimum. Then, from (48), we have

$$\mathcal{J}_\alpha(\hat{P}, P_{\theta^*}) = \mathcal{J}_\alpha(\hat{P}, Q) + \mathcal{J}_\alpha(Q, P_{\theta^*}). \quad (49)$$

Since $\mathcal{I}_\alpha(\hat{P}, P_{\theta^*}) = \mathcal{I}_\alpha(\hat{P}, Q)$, we have $\mathcal{I}_\alpha(Q, P_{\theta^*}) = 0$, and so $P_{\theta^*} = Q$.

(b) Let $\text{Supp}(Q) \neq \mathbb{X}$. Then, by Th. 16(a), $Q \in \text{cl}(\mathbb{M}^{(\alpha)}) \setminus \mathbb{M}^{(\alpha)}$. Uniqueness on the closure follows just as in (a) immediately above. If \hat{P} has a reverse \mathcal{I}_α -projection on $\mathbb{M}^{(\alpha)}$, say P_{θ^*} , then by continuity of $\mathcal{I}_\alpha(\hat{P}, \cdot)$ ([2, Rem. 5]), we have $\mathcal{I}_\alpha(\hat{P}, Q) = \mathcal{I}_\alpha(\hat{P}, P_{\theta^*})$. This contradicts the uniqueness. ■

B. $\alpha > 1$:

Let us begin with a counterexample that shows that Theorem 16 does not hold when $\alpha > 1$; $\text{cl}(\mathbb{M}^{(\alpha)})$ need not intersect the associated \mathbb{L} .

Example 3: Let $\alpha, \mathbb{X}, \mathbb{L}$, and R be as in Example 2. The associated α -power-law family and its closure are

$$\mathbb{M}^{(\alpha)} = \left\{ P_\theta : \theta \in (-1/24, 1/4) \right\},$$

and

$$\text{cl}(\mathbb{M}^{(\alpha)}) = \left\{ P_\theta : \theta \in [-1/24, 1/4] \right\},$$

where

$$P_\theta = \frac{1}{1 + 13\theta} \left(1/4 - \theta, 1/4 + 3\theta, 1/4 + 5\theta, 1/4 + 6\theta \right).$$

We assert that no such P_θ , either of $\mathbb{M}^{(\alpha)}$ or $\text{cl}(\mathbb{M}^{(\alpha)})$, is in \mathbb{L} . Furthermore, the forward \mathcal{I}_α -projection of every member in $\text{cl}(\mathbb{M}^{(\alpha)})$ on \mathbb{L} is $Q = (3/4, 1/4, 0, 0)$ which, of course, is not in $\text{cl}(\mathbb{M}^{(\alpha)})$.

One must therefore extend $\mathbb{M}^{(\alpha)}$ beyond its closure to identify the family that is orthogonal to \mathbb{L} and intersects \mathbb{L} at Q . An appropriate extension of $\mathbb{M}^{(\alpha)}$ that intersects \mathbb{L} turns out to be the following.

Definition 19: The family $\hat{\mathbb{M}}_+^{(\alpha)}$ characterized by a probability measure R and k functions $f_i : \mathbb{X} \rightarrow \mathbb{R}, i = 1, \dots, k$, is defined as follows. Let $Q = P_{\theta^*}$ be the forward \mathcal{I}_α -projection¹³ of R on \mathbb{L} . Define $\hat{\mathbb{M}}_+^{(\alpha)}$ to be the set of all probability measures P_θ satisfying (a), (b), and (c) below.

(a)

$$Z(\theta)^{\alpha-1} P_\theta^{\alpha-1}(x) = \left[R(x)^{\alpha-1} + (1-\alpha) \sum_{i=1}^k \theta_i f_i(x) \right]_+ \quad \forall x \in \mathbb{X},$$

where $Z(\theta)$ is the normalization constant that makes P_θ a valid probability measure on \mathbb{X} .

(b) $\text{Supp}(P_{\theta^*}) \subseteq \text{Supp}(P_\theta)$;

(c) $\sum_{i=1}^k \theta_i f_i(x) \leq \sum_{i=1}^k \theta_i^* f_i(x) \quad \forall x \notin \text{Supp}(P_\theta)$.

The following is the analog of the combined Theorem 16 and Corollary 17.

Theorem 20: Let $\alpha > 1$. Let \mathbb{L} be a linear family characterized by $f_i, i = 1, \dots, k$ as in (16). Let $\mathbb{M}^{(\alpha)}$ be as in Definition 8, characterized by R and the k functions $f_i, i = 1, \dots, k$. Let Q be the forward \mathcal{I}_α -projection of R on \mathbb{L} . Let $\hat{\mathbb{M}}_+^{(\alpha)}$ be the extension of $\mathbb{M}^{(\alpha)}$ as in Definition 19. We then have the following.

(a) $\mathbb{L} \cap \hat{\mathbb{M}}_+^{(\alpha)} = \{Q\}$ and

$$\mathcal{I}_\alpha(P, P_\theta) \geq \mathcal{I}_\alpha(P, Q) + \mathcal{I}_\alpha(Q, P_\theta) \quad (50)$$

for every $P \in \mathbb{L}$ and every $P_\theta \in \hat{\mathbb{M}}_+^{(\alpha)}$.

(b) If $Q \in \text{cl}(\mathbb{M}^{(\alpha)})$, then $\mathbb{L} \cap \text{cl}(\mathbb{M}^{(\alpha)}) = \{Q\}$ and (50) holds with equality for every $P \in \mathbb{L}$ and every $P_\theta \in \text{cl}(\mathbb{M}^{(\alpha)})$.

(c) If $Q \in \mathbb{M}^{(\alpha)}$, then $\mathbb{L} \cap \mathbb{M}^{(\alpha)} = \{Q\}$ and (50) holds with equality for every $P \in \mathbb{L}$ and every $P_\theta \in \mathbb{M}^{(\alpha)}$.

Proof: (a) By virtue of Theorem 14-(b), we have $Q \in \mathbb{L} \cap \hat{\mathbb{M}}_+^{(\alpha)}$. Furthermore, by Theorem 15, any member of $\mathbb{L} \cap \hat{\mathbb{M}}_+^{(\alpha)}$ is a forward \mathcal{I}_α -projection of R on \mathbb{L} . Since the forward projection is unique, $\mathbb{L} \cap \hat{\mathbb{M}}_+^{(\alpha)}$ must be the singleton $\{Q\}$.

Let $P_\theta \in \hat{\mathbb{M}}_+^{(\alpha)}$. We claim that P_θ has $P_{\theta^*} = Q$ as its forward projection on \mathbb{L} . Assuming the claim, by Theorem 14-(c), inequality (50) holds.

Let us now proceed to show the claim. By Theorem 15, it suffices to verify that P_{θ^*} can be written as

$$\tilde{Z}(\tilde{\theta})^{\alpha-1} P_{\theta^*}^{\alpha-1}(x) = \left[P_\theta(x)^{\alpha-1} + (1-\alpha) \sum_{i=1}^k \tilde{\theta}_i f_i(x) \right]_+ \quad \forall x \quad (51)$$

¹³By virtue of Th. 14(b), Q is of the form (28) for some θ^* and hence may be written as $Q = P_{\theta^*}$.

for some $\tilde{Z}(\tilde{\theta})$ and $\tilde{\theta} = (\tilde{\theta}_1, \dots, \tilde{\theta}_k)$. To see this, by definition of P_θ , we have

$$Z(\theta)^{\alpha-1} P_\theta(x)^{\alpha-1} = \left[R(x)^{\alpha-1} + (1-\alpha) \sum_{i=1}^k \theta_i f_i(x) \right]_+ \quad \forall x, \quad (52)$$

and, by Theorem 14-(b), we have

$$Z(\theta^*)^{\alpha-1} P_{\theta^*}(x)^{\alpha-1} = \left[R(x)^{\alpha-1} + (1-\alpha) \sum_{i=1}^k \theta_i^* f_i(x) \right]_+ \quad \forall x. \quad (53)$$

Let $x \in \text{Supp}(P_{\theta^*})$. By Definition 19-(a), $x \in \text{Supp}(P_\theta)$ as well. Hence, we can remove the $[\cdot]_+$ operation in (52) and (53) to get

$$Z(\theta)^{\alpha-1} P_\theta(x)^{\alpha-1} = R(x)^{\alpha-1} + (1-\alpha) \sum_{i=1}^k \theta_i f_i(x).$$

$$Z(\theta^*)^{\alpha-1} P_{\theta^*}(x)^{\alpha-1} = R(x)^{\alpha-1} + (1-\alpha) \sum_{i=1}^k \theta_i^* f_i(x),$$

Eliminating $R(x)^{\alpha-1}$ from the preceding equations, we get

$$Z(\theta^*)^{\alpha-1} P_{\theta^*}(x)^{\alpha-1} = Z(\theta)^{\alpha-1} P_\theta(x)^{\alpha-1} + (1-\alpha) \sum_{i=1}^k (\theta_i^* - \theta_i) f_i(x),$$

equivalently,

$$\left(\frac{Z(\theta^*)}{Z(\theta)} \right)^{\alpha-1} P_{\theta^*}(x)^{\alpha-1} = P_\theta(x)^{\alpha-1} + (1-\alpha) \sum_{i=1}^k \frac{(\theta_i^* - \theta_i)}{Z(\theta)^{\alpha-1}} f_i(x). \quad (54)$$

This suggests that $\tilde{Z}(\tilde{\theta}) = Z(\theta^*)/Z(\theta)$ and $\tilde{\theta}_i = (\theta_i^* - \theta_i)/Z(\theta)^{\alpha-1}$ should work. Let us now verify that they do, that is, that (51) holds for all x with these choices of \tilde{Z} and $\tilde{\theta}$.

The foregoing shows (51) holds for all $x \in \text{Supp}(P_{\theta^*})$. Next, let $x \in \text{Supp}(P_\theta) \setminus \text{Supp}(P_{\theta^*})$. The right-hand side of (54), upon substitution of (52) without the $[\cdot]_+$ operation, becomes

$$\frac{R(x)^{\alpha-1} + (1-\alpha) \sum_{i=1}^k \theta_i f_i(x)}{Z(\theta)^{\alpha-1}} + (1-\alpha) \sum_{i=1}^k \frac{(\theta_i^* - \theta_i)}{Z(\theta)^{\alpha-1}} f_i(x) = \frac{R(x)^{\alpha-1} + (1-\alpha) \sum_{i=1}^k \theta_i^* f_i(x)}{Z(\theta)^{\alpha-1}} \leq 0,$$

as is required for $x \notin \text{Supp}(P_{\theta^*})$. Hence (51) holds for $x \in \text{Supp}(P_\theta) \setminus \text{Supp}(P_{\theta^*})$ as well, and therefore for all $x \in \text{Supp}(P_\theta)$.

Finally, when $x \notin \text{Supp}(P_\theta)$,

$$R(x)^{\alpha-1} + (1-\alpha) \sum_{i=1}^k \theta_i f_i(x) \leq 0.$$

The right-hand side of (54) then satisfies

$$(1-\alpha) \sum_{i=1}^k \frac{(\theta_i^* - \theta_i)}{Z(\theta)^{\alpha-1}} f_i(x) \leq 0$$

because of condition (b) in Definition 19 and $\alpha > 1$. This establishes that P_{θ^*} is of the form (51), and is therefore the forward \mathcal{S}_α -projection of P_θ on \mathbb{L} .

Proofs of (b) and (c) are the same as in $\alpha < 1$ case considered in Theorem 16. ■

Having established the orthogonality between a linear family and its associated α -power-law family, let us now return to the problem of robust estimation discussed in section II-B. As in the case of $\alpha < 1$, we show a connection between the MMPLE on the extended α -power-law family $\hat{\mathbb{M}}_+^{(\alpha)}$, which is a reverse \mathcal{S}_α -projection on $\hat{\mathbb{M}}_+^{(\alpha)}$, and the forward \mathcal{S}_α -projection on the related linear family.

Theorem 21: Let $\alpha > 1$. Let \hat{P} be a probability measure on \mathbb{X} . Let $\mathbb{M}^{(\alpha)}$ be characterized by the probability measure R and the functions $f_i, i = 1, \dots, k$. Let R have full support. Let \mathbb{L} be the associated linear family characterized by $f_i, i = 1, \dots, k$, and assume that it is nonempty. Define $\tilde{\mathbb{L}}$ as in (45) using \tilde{f}_i and τ_i^R as defined in (46) and (47), respectively. Let Q be the forward \mathcal{S}_α -projection of R on $\tilde{\mathbb{L}}$. Then the following hold.

(a) If $Q \in \mathbb{M}^{(\alpha)}$, then Q is the unique reverse \mathcal{S}_α -projection of \hat{P} on $\mathbb{M}^{(\alpha)}$.

- (b) If $Q \in \text{cl}(\mathbb{M}^{(\alpha)}) \setminus \mathbb{M}^{(\alpha)}$, then \hat{P} does not have a reverse \mathcal{J}_α -projection on $\mathbb{M}^{(\alpha)}$. However, Q is the unique reverse \mathcal{J}_α -projection of \hat{P} on $\text{cl}(\mathbb{M}^{(\alpha)})$.
- (c) If $Q \notin \text{cl}(\mathbb{M}^{(\alpha)})$, then
- \hat{P} does not have a reverse \mathcal{J}_α -projection on $\mathbb{M}^{(\alpha)}$.
 - $\mathbb{M}^{(\alpha)}$ can be extended to $\hat{\mathbb{M}}_+^{(\alpha)}(R, \tilde{f}_1, \dots, \tilde{f}_k)$, and Q is the unique reverse \mathcal{J}_α -projection of \hat{P} on $\hat{\mathbb{M}}_+^{(\alpha)}(R, \tilde{f}_1, \dots, \tilde{f}_k)$.

Proof: Only (c)-(i) needs a proof. Proofs of all others follow the same arguments in the proof of Theorem 18, but now one uses Theorem 20 instead of Corollary 17.

Let us now prove (c)-(i) by contradiction. Suppose \hat{P} has a reverse \mathcal{J}_α -projection on $\mathbb{M}^{(\alpha)}$. Call it P_{θ^*} . Since P_{θ^*} has full support, there is a neighborhood N of θ^* such that $\theta \in N$ implies $P_\theta \in \mathbb{M}^{(\alpha)}$. The first order optimality condition applies, namely

$$\left. \frac{\partial}{\partial \theta_i} \mathcal{J}_\alpha(\hat{P}, P_\theta) \right|_{\theta=\theta^*} = 0, \quad i = 1, \dots, k.$$

We claim that this implies

$$\sum_x P_{\theta^*}(x) \tilde{f}_i(x) = 0, \quad i = 1, \dots, k. \quad (55)$$

But then $P_{\theta^*} \in \tilde{\mathbb{L}}$ and so $P_{\theta^*} = Q$, a contradiction to $Q \notin \text{cl}(\mathbb{M}^{(\alpha)})$.

We now proceed to prove the claim (55). Observe that, since $P_\theta \in \mathbb{M}^{(\alpha)}$, by Definition 8, we have

$$Z(\theta)^{\alpha-1} P_\theta(x)^{\alpha-1} = R(x)^{\alpha-1} + \sum_j \theta_j \tilde{f}_j(x), \quad (56)$$

and so

$$\begin{aligned} Z(\theta)^{\alpha-1} \sum_x \hat{P}(x) P_\theta(x)^{\alpha-1} &= \sum_x \hat{P}(x) R(x)^{\alpha-1} + \sum_j \theta_j \left(\sum_x \hat{P}(x) \tilde{f}_j(x) \right) \\ &= \sum_x \hat{P}(x) R(x)^{\alpha-1}, \end{aligned} \quad (57)$$

where the last equality holds because $\hat{P} \in \tilde{\mathbb{L}}$. Also,

$$\begin{aligned} \sum_x P_\theta(x)^\alpha &= \sum_x \left[P_\theta(x)^{\alpha-1} \right]^{\frac{\alpha}{\alpha-1}} \\ &= Z(\theta)^{-\alpha} \sum_x \left[R(x)^{\alpha-1} + \sum_j \theta_j \tilde{f}_j(x) \right]^{\frac{\alpha}{\alpha-1}}. \end{aligned} \quad (58)$$

Substituting (57) and (58) into (12) and taking the partial derivative, we get

$$\begin{aligned} \frac{\partial}{\partial \theta_i} \mathcal{J}_\alpha(\hat{P}, P_\theta) &= \frac{\alpha}{1-\alpha} \frac{\partial}{\partial \theta_i} \log Z(\theta)^{1-\alpha} + \frac{\partial}{\partial \theta_i} \log Z(\theta)^{-\alpha} + \frac{\partial}{\partial \theta_i} \log \sum_x \left[R(x)^{\alpha-1} + \sum_j \theta_j \tilde{f}_j(x) \right]^{\frac{\alpha}{\alpha-1}} \\ &= \frac{\partial}{\partial \theta_i} \log \sum_x \left[R(x)^{\alpha-1} + \sum_j \theta_j \tilde{f}_j(x) \right]^{\frac{\alpha}{\alpha-1}} \\ &= \frac{1}{A} \cdot \frac{\alpha}{1-\alpha} \sum_x \left[R(x)^{\alpha-1} + \sum_j \theta_j \tilde{f}_j(x) \right]^{\frac{1}{\alpha-1}} \tilde{f}_i(x) \\ &= \frac{1}{A} \cdot \frac{\alpha}{1-\alpha} Z(\theta) \sum_x P_\theta(x) \tilde{f}_i(x), \end{aligned}$$

where $A = \sum_x \left[R(x)^{\alpha-1} + \sum_j \theta_j \tilde{f}_j(x) \right]^{\frac{\alpha}{\alpha-1}}$, and the last equality follows from (56). Thus,

$$\left. \frac{\partial}{\partial \theta_i} \mathcal{J}_\alpha(\hat{P}, P_\theta) \right|_{\theta=\theta^*} = 0 \implies \sum_x P_{\theta^*}(x) \tilde{f}_i(x) = 0,$$

thereby proving the claim. ■

VII. EPILOGUE

We now provide some concluding remarks. Our focus has primarily been on the geometric relation between the α -power-law and the linear families. This geometric relation enabled us to characterize the reverse \mathcal{S}_α -projection on an α -power-law family $\mathbb{M}^{(\alpha)} := \mathbb{M}^{(\alpha)}(R, f_1, \dots, f_k)$ as a forward \mathcal{S}_α -projection on a linear family. The procedure is as follows.

“Given the family $\mathbb{M}^{(\alpha)}$, sweep through a collection of linear families (45)-(47) orthogonal to $\mathbb{M}^{(\alpha)}$ by varying $\tau_i^R, i = 1, \dots, k$, and find the linear family $\tilde{\mathbb{L}}$ that contains \hat{P} . Then find the forward \mathcal{S}_α -projection of R on $\tilde{\mathbb{L}}$; call it Q . If $Q \in \mathbb{M}^{(\alpha)}$, then Q is the reverse \mathcal{S}_α -projection of \hat{P} on the $\mathbb{M}^{(\alpha)}$. If $Q \in \text{cl}(\mathbb{M}^{(\alpha)}) \setminus \mathbb{M}^{(\alpha)}$, then \hat{P} does not have a reverse \mathcal{S}_α -projection on $\mathbb{M}^{(\alpha)}$. But Q attains the minimum in the closure.”

The cases $\alpha < 1$ and $\alpha > 1$ have different characteristics. The $\alpha < 1$ case is similar to $\alpha = 1$ and one always has $\tilde{\mathbb{L}} \cap \text{cl}(\mathbb{M}^{(\alpha)}) = \{Q\}$. On the other hand, when $\alpha > 1$, it is possible that $\tilde{\mathbb{L}} \cap \text{cl}(\mathbb{M}^{(\alpha)}) = \emptyset$, and $Q \notin \text{cl}(\mathbb{M}^{(\alpha)})$. Then \hat{P} does not have a reverse \mathcal{S}_α -projection on $\mathbb{M}^{(\alpha)}$. One then needs to extend $\mathbb{M}^{(\alpha)}$ to make it intersect $\tilde{\mathbb{L}}$. We showed that the extension $\tilde{\mathbb{M}}_+^{(\alpha)}$ is just right and satisfies $\tilde{\mathbb{L}} \cap \tilde{\mathbb{M}}_+^{(\alpha)} = \{Q\}$. However, Q , in the intersection $\tilde{\mathbb{L}} \cap \tilde{\mathbb{M}}_+^{(\alpha)}$, is no longer the reverse \mathcal{S}_α -projection of \hat{P} on $\text{cl}(\mathbb{M}^{(\alpha)})$. It would be interesting to see if Q can be used to simplify the computation of the true reverse \mathcal{S}_α -projection of \hat{P} on $\text{cl}(\mathbb{M}^{(\alpha)})$.

Our characterization has algorithmic benefits since the forward \mathcal{S}_α -projection is a minimization of a quasiconvex function subject to linear constraints. Standard techniques are available to solve such problems, for example, via a sequence of convex feasibility problems [23, Sec. 4.2.5], or via a sequence of simpler forward projections on single-constraint linear families [2, Th. 16, Rem. 13].

APPENDIX A

WEAK DEPENDENCE OF THE α -POWER-LAW FAMILY ON R

The following result shows that the α -power-law family depends on R only in a weak manner, and that any member of $\mathbb{M}^{(\alpha)}$ could equally well play the role of R . The same result is well-known for an exponential family.

Proposition 22: If $\alpha > 1$, let R have full support. Consider the $\mathbb{M}^{(\alpha)}(R, f_1, \dots, f_k)$ as in Definition 8. Fix $P_{\theta^*} \in \mathbb{M}^{(\alpha)}(R, f_1, \dots, f_k)$. Then $\mathbb{M}^{(\alpha)}(P_{\theta^*}, f_1, \dots, f_k) = \mathbb{M}^{(\alpha)}(R, f_1, \dots, f_k)$.

Proof: Write $\mathbb{M}^{(\alpha)}$ for $\mathbb{M}^{(\alpha)}(R, f_1, \dots, f_k)$ and $\tilde{\mathbb{M}}^{(\alpha)}$ for $\mathbb{M}^{(\alpha)}(P_{\theta^*}, f_1, \dots, f_k)$. We will check that an arbitrary element $P_\theta \in \mathbb{M}^{(\alpha)}$ is an element of $\tilde{\mathbb{M}}^{(\alpha)}$. This will establish $\mathbb{M}^{(\alpha)} \subset \tilde{\mathbb{M}}^{(\alpha)}$. The converse holds by symmetry.

From the formula for P_{θ^*} , observe that

$$P_{\theta^*}(x)^{\alpha-1} = Z(\theta^*)^{1-\alpha} \left[R(x)^{\alpha-1} + (1-\alpha) \sum_{i=1}^k \theta_i^* f_i(x) \right] \quad \forall x,$$

and so

$$R(x)^{\alpha-1} = Z(\theta^*)^{\alpha-1} P_{\theta^*}(x)^{\alpha-1} - (1-\alpha) \sum_{i=1}^k \theta_i^* f_i(x) \quad \forall x. \quad (59)$$

Substitute this into the formula for P_θ in (19) to get

$$\begin{aligned} P_\theta(x)^{\alpha-1} &= Z(\theta)^{1-\alpha} \left[Z(\theta^*)^{\alpha-1} P_{\theta^*}(x)^{\alpha-1} - (1-\alpha) \sum_{i=1}^k \theta_i^* f_i(x) + (1-\alpha) \sum_{i=1}^k \theta_i f_i(x) \right] \\ &= \left(\frac{Z(\theta^*)}{Z(\theta)} \right)^{\alpha-1} \left[P_{\theta^*}(x)^{\alpha-1} + (1-\alpha) \sum_{i=1}^k \frac{\theta_i - \theta_i^*}{Z(\theta^*)^{\alpha-1}} f_i(x) \right] \\ &= \tilde{Z}(\xi)^{1-\alpha} \left[P_{\theta^*}(x)^{\alpha-1} + (1-\alpha) \sum_{i=1}^k \xi_i f_i(x) \right], \end{aligned}$$

where $\xi = (\theta - \theta^*)/Z(\theta^*)^{\alpha-1}$, and $\tilde{Z}(\xi) = Z(\theta)/Z(\theta^*)$. Thus, $P_\theta \in \tilde{\mathbb{M}}^{(\alpha)}$. ■

Change of reference from R to P_{θ^*} merely amounts to a translation and rescaling of the parameter space.

ACKNOWLEDGEMENTS

We thank the reviewers whose comments/suggestions helped improve this manuscript enormously.

REFERENCES

- [1] M. Ashok Kumar and R. Sundaresan, "Relative α -entropy minimizers subject to linear statistical constraints," *arXiv:1410.4931*, October 2014.
- [2] —, "Minimization problems based on a parametric family of relative entropies I: Forward projection," *arXiv:1410.2346*, October 2014.
- [3] C. R. Rao, *Linear Statistical Inference and its Applications*, 2nd ed. New Delhi, India: Wiley Eastern Limited, 1973, 6th Wiley Eastern Reprint, March 1991.
- [4] A. Basu, I. R. Harris, N. L. Hjort, and M. C. Jones, "Robust and efficient estimation by minimising a density power divergence," *Biometrika*, vol. 85, pp. 549–559, 1998.
- [5] C. Field and B. Smith, "Robust estimation: A weighted maximum likelihood approach," *International Statistical Review*, vol. 62, no. 3, pp. 405–424, December 1994.
- [6] M. P. Windham, "Robustifying model fitting," *Journal of the Royal Statistical Society. Series B (Methodological)*, vol. 57, no. 3, pp. 599–609, 1995.
- [7] M. C. Jones, N. L. Hjort, I. R. Harris, and A. Basu, "A comparison of related density based minimum divergence estimators," *Biometrika*, vol. 88, no. 3, pp. 865–873, 2001.
- [8] S. Eguchi and S. Kato, "Entropy and divergence associated with power function and the statistical application," *Entropy*, vol. 12, no. 2, pp. 262–274, 2010.
- [9] H. Fujisawa and S. Eguchi, "Robust parameter estimation with a small bias against heavy contamination," *Journal of Multivariate Analysis*, vol. 99, pp. 2053–2081, 2008.
- [10] A. Cichocki and S. Amari, "Families of alpha- beta- and gamma- divergences: Flexible and robust measures of similarities," *Entropy*, vol. 12, pp. 1532–1568, 2010.
- [11] I. Csiszár, "Why least squares and maximum entropy? An axiomatic approach to inference for linear inverse problems," *The Annals of Statistics*, vol. 19, no. 4, pp. 2032–2066, 1991.
- [12] R. Sundaresan, "A measure of discrimination and its geometric properties," in *Proc. of the 2002 IEEE International Symposium on Information Theory*, Lausanne, Switzerland, June 2002, p. 264.
- [13] —, "Guessing under source uncertainty," *Information Theory, IEEE Transactions on*, vol. 53, no. 1, pp. 269–287, January 2007.
- [14] L. L. Campbell, "A coding theorem and Rényi's entropy," *Information and Control*, vol. 8, pp. 423–429, 1965.
- [15] E. Arikan, "An inequality on guessing and its application to sequential decoding," *Information Theory, IEEE Transactions on*, vol. 42, no. 1, pp. 99–105, January 1996.
- [16] M. K. Hanawal and R. Sundaresan, "Guessing revisited: A large deviations approach," *Information Theory, IEEE Transactions on*, vol. 57, no. 1, pp. 70–78, January 2011.
- [17] C. Bunte and A. Lapidoth, "Codes for tasks and Rényi entropy," *Information Theory, IEEE Transactions on*, vol. 60, no. 9, pp. 5065–5076, September 2014.
- [18] C. Tsallis, "What are the numbers that experiments provide," *Quimica Nova*, vol. 17, no. 6, pp. 468–471, 1994.
- [19] T. van Erven and P. Harremoës, "Rényi divergence and Kullback-Leibler divergence," *Information Theory, IEEE Transactions on*, vol. 60, no. 7, pp. 3797–3820, July 2014.
- [20] I. Csiszár and P. Shields, *Information Theory and Statistics: A Tutorial*, ser. Foundations and Trends in Communications and Information Theory. Hanover, USA: Now Publishers Inc, 2004, vol. 1, no. 4.
- [21] I. Csiszár and F. Matúš, "Information projections revisited," *Information Theory, IEEE Transactions on*, vol. 49, no. 6, pp. 1474–1490, June 2003.
- [22] D. P. Bertsekas, *Nonlinear Programming*, 2nd ed. Belmont, MA: Athena Scientific, 2003.
- [23] S. Boyd and L. Vandenberghe, *Convex Optimization*. The Edinburgh Building, Cambridge, CB2 8RU, UK: Cambridge University Press, 2004.