

PRIMAL AND DUAL FORMULATIONS RELEVANT FOR THE NUMERICAL ESTIMATION OF A PROBABILITY DENSITY VIA REGULARIZATION

ROGER KOENKER — IVAN MIZERA

ABSTRACT. General schemes relevant for the estimation of a probability density via regularization—primal and dual versions in the discretized setting—are investigated. Conditions for the dual solution to be a probability density are given, and a strong duality theorem is proved.

We study various instances of the problem

$$-w^T Lh + s^T \Psi(g) + J(-Ph) = \min_{g,h}!, \quad \text{subject to } h \preceq g, \quad (P)$$

where L and w are evaluation operator and averaging functional described later in the text; $\Psi(g)$ indicates the application of a real convex function ψ to the components of g , while $J(h)$ is a general convex function applied to the whole vector $-Ph$, the negative of the result of a linear operator P applied on h . We assume that vectors w and s have nonnegative elements; hereafter, \succeq and \preceq denote componentwise inequalities. If ψ is nondecreasing, the primal formulation (P) can be simplified—it is equivalent to the unconstrained problem

$$-w^T Lg + s^T \Psi(g) + J(-Pg) = \min_g!. \quad (U)$$

Convex functions are allowed to attain $+\infty$ as a value; the *domain*, $\text{dom } \Phi$, is the set where Φ is finite. We assume that all convex functions in (P) and (U) have domains with nonempty interiors. Concave functions are handled in an analogous manner, only the role of $+\infty$ is played by $-\infty$.

A *conjugate* of a convex function Φ is

$$\Phi^*(y) = \sup_x (y^T x - \Phi(x)) = \sup_{x \in \text{dom } \Phi} (y^T x - \Phi(x)),$$

2000 Mathematics Subject Classification: Primary 62G07; Secondary 94A17, 90C46, 65J20.

Keywords: density estimation, penalty methods, duality, Rényi entropies.

This research was partially supported by NSF grant SES-05-44673 and by the Natural Sciences and Engineering Research Council of Canada.

the latter formulation avoiding the need to compute with infinite values. The conjugate of the function $\lambda \|\cdot\|_p$, where $\|\cdot\|_p$ stands for the ℓ^p norm, is the indicator of the ball in the dual norm, $\{x: \|x\|_q \leq \lambda\}$, where q is the (Hölder) conjugate of p (that is, $q = \infty$ for $p = 1$, and $(1/p) + (1/q) = 1$ for $p > 1$) and the *indicator* of a convex set E is defined to be 0 for all $x \in E$ and $+\infty$ otherwise. The conjugate of the indicator of the cone $\{x: x \succeq 0\}$ is the indicator of the polar cone $\{x: x \preceq 0\}$. Finally, the function $(1/2)\|\cdot\|_2^2$ is conjugate to itself; and consequently, $\lambda\|\cdot\|_2^2$ to $1/(4\lambda)\|\cdot\|_2^2$. Our references for convex analysis are Rockafellar [10], Boyd and Vandenberghe [1].

We claim that the dual of (P), or when equivalent, (U) is the problem

$$\begin{aligned} -s^T \Psi^*(f) - J^*(e) &= \max_{f, e} !, \\ \text{subject to } Sf &= L^T w + P^T e \quad \text{and} \quad f \succeq 0, \end{aligned} \tag{D}$$

where $S = \text{diag}(s)$ and $\Psi^*(f)$ indicates the componentwise application of ψ^* .

Both (P)–(U) and (D) are relevant in the study of discretized, numerical formulations of regularized density estimation. The estimated density is represented by the vector f of its values on some collection of points, referred to as a *grid*. The evaluation operator L then expresses the position of n datapoints with respect to the grid via interpolation; for instance, if the datapoints are among gridpoints, then the i th row assigns 1 to a gridpoint equal to the i th datapoint and zero otherwise. The vector w assigns weights to the datapoints—as a rule, $1/n$ to each. Finally, s is the vector of integration weights attached to gridpoints: the identity $s^T f = 1$ expresses the fact that the estimated density integrates to 1. Estimated probability densities are approximated by the densities with respect to the dominating measure on the grid whose atoms are given by s .

As for the penalization term, a typical P is a discretized version of a differential operator appearing in the continuous formulation of the regularization proposal. Typical J involves an ℓ^p norm and a tuning constant, λ , customary in this context: say, $J(u) = \lambda\|u\|_1$ or $J(u) = \lambda\|u\|_2^2$.

Regularization may be also expressed in a constrained form, in which J is the indicator of a set $\{u: \|u\|_p \leq \Lambda\}$.

All these examples are symmetric: $J(-u) = J(u)$. An asymmetric example is provided by J equal to the indicator of $\{u: u \preceq 0\}$, the style of penalization used in density estimation under monotonicity or convexity constraints.

The fact that the estimated f is indeed a probability density can be most conveniently verified through the dual formulation (D).

THEOREM 1. *Suppose that $w^T L 1 = 1$ and $P 1 = 0$. Then the solution f of (D) satisfies $\sum_j s_j f_j = 1$ and $f_j \geq 0$ for every j .*

Proof. The nonnegativity constraint is directly included in the formulation of (D); it remains to verify that

$$\mathbf{s}^\top \mathbf{f} = \mathbf{1}^\top \mathbf{S} \mathbf{f} = \mathbf{1}^\top (\mathbf{L}^\top \mathbf{w} + \mathbf{P}^\top \mathbf{e}) = \mathbf{w}^\top \mathbf{L} \mathbf{1} + \mathbf{e}^\top \mathbf{P} \mathbf{1} = 1.$$

□

In the simplest case, when matrix \mathbf{L} is composed of zeros, except for a single 1 in each row corresponding to a datapoint, $\mathbf{w}^\top \mathbf{L}$ makes these 1's multiplied by $1/n$; further multiplying by 1 makes them sum to 1. More generally, common interpolation schemes yield evaluation operators satisfying the assumption of Theorem 1. As far as potential operators \mathbf{P} are concerned, they are discrete, difference versions of differential operators; as such, they annihilate constants—as can be directly verified for difference operators acting on sequences.

Compared to the dual (D), the relationship of the variables appearing in the primal formulations (P) or (U) to the estimated density is not explicit. However, once a strong duality of (P) and (D) is demonstrated true, then the relationship of \mathbf{g} to \mathbf{f} for qualified ψ is given by

$$\mathbf{f} = \Psi'(\mathbf{g}), \quad (\text{E})$$

where $\Psi'(\mathbf{g})$ indicates the componentwise application of ψ' , the derivative of ψ .

THEOREM 2. *Problem (D) is a strong dual of the problem (P). If ψ is differentiable on the interior I of its domain, then the corresponding solutions of (D) and (P) satisfy (E), whenever \mathbf{g} and \mathbf{f} are componentwise from I and the image of I under ψ' , respectively.*

Proof. We take a formulation equivalent to (P), obtained by rewriting it in terms of new variables \mathbf{u} and \mathbf{v} ,

$$\begin{aligned} -\mathbf{w}^\top \mathbf{L}(\mathbf{g} - \mathbf{v}) + \mathbf{s}^\top \Psi(\mathbf{g}) + J(\mathbf{u}) &= \min_{\mathbf{g}, \mathbf{u}, \mathbf{v}}!, \\ \text{subject to } \mathbf{v} &\succeq \mathbf{0} \quad \text{and} \quad -\mathbf{P}(\mathbf{g} - \mathbf{v}) = \mathbf{u}. \end{aligned} \quad (1)$$

The Lagrange dual of (1) is

$$\inf_{\mathbf{g}, \mathbf{u}, \mathbf{v}} \mathcal{L}(\mathbf{p}, \mathbf{e}; \mathbf{g}, \mathbf{u}, \mathbf{v}) = \max_{\mathbf{p}, \mathbf{e}}!, \quad \text{subject to } \mathbf{p} \succeq \mathbf{0}, \quad (2)$$

where (\mathbf{p} used here has no relationship to the parameter p used elsewhere)

$$\mathcal{L}(\mathbf{p}, \mathbf{e}; \mathbf{g}, \mathbf{u}, \mathbf{v}) = -\mathbf{w}^\top \mathbf{L}(\mathbf{g} - \mathbf{v}) + \mathbf{s}^\top \Psi(\mathbf{g}) + J(\mathbf{u}) + \mathbf{p}^\top (-\mathbf{v}) + \mathbf{e}^\top [-\mathbf{u} - \mathbf{P}(\mathbf{g} - \mathbf{v})]$$

is the Lagrangean of (1). The linear part of \mathcal{L} , in \mathbf{v} , leads to a feasibility constraint

$$\mathbf{L}^\top \mathbf{w} + \mathbf{P}^\top \mathbf{e} = \mathbf{p}, \quad (3)$$

preventing the objective function of (2) from becoming $-\infty$. Under (3), the minimization of the simplified Lagrangean can be done separately in \mathbf{g} and \mathbf{u} ,

$$\begin{aligned} \inf_{\mathbf{g}, \mathbf{u}} \mathcal{L}(\mathbf{p}, \mathbf{e}; \mathbf{g}, \mathbf{u}) &= \inf_{\mathbf{g}, \mathbf{u}} (-\mathbf{w}^\top \mathbf{L} \mathbf{g} - \mathbf{e}^\top \mathbf{P} \mathbf{g} + \mathbf{s}^\top \Psi(\mathbf{g}) - \mathbf{e}^\top \mathbf{u} + \mathbf{J}(\mathbf{u})) \\ &= \inf_{\mathbf{g}} (-(\mathbf{L}^\top \mathbf{w} + \mathbf{P}^\top \mathbf{e})^\top \mathbf{g} + \mathbf{s}^\top \Psi(\mathbf{g})) + \inf_{\mathbf{u}} (-\mathbf{e}^\top \mathbf{u} + \mathbf{J}(\mathbf{u})), \quad (4) \\ &= \inf_{\mathbf{g}} (-\mathbf{p}^\top \mathbf{g} + \mathbf{s}^\top \Psi(\mathbf{g})) - \mathbf{J}^*(\mathbf{e}). \end{aligned}$$

Minimizing in \mathbf{g} is done by expanding into components,

$$\begin{aligned} \inf_{\mathbf{g}} (-\mathbf{p}^\top \mathbf{g} + \mathbf{s}^\top \Psi(\mathbf{g})) &= \inf_{\mathbf{g}} \left(-\sum_j \mathbf{p}_j \mathbf{g}_j + \sum_j \mathbf{s}_j \psi(\mathbf{g}_j) \right) \\ &= \sum_j \mathbf{s}_j \inf_{\mathbf{g}_j} \left(-\frac{\mathbf{p}_j}{\mathbf{s}_j} \mathbf{g}_j + \psi(\mathbf{g}_j) \right) = -\sum_j \mathbf{s}_j \psi^* \left(\frac{\mathbf{p}_j}{\mathbf{s}_j} \right). \end{aligned} \quad (5)$$

The dual formulation (D) is obtained as the summary of (2)–(5), rewritten in terms of $\mathbf{f}_j = \mathbf{p}_j/\mathbf{s}_j$. Finally, (1) satisfies the Slater constraint qualification condition; therefore strong duality holds.

For fixed y , the domain of the concave function $\varphi(x) = yx - \psi(x)$ is the same as the domain of ψ . If ψ has a derivative on I , so does φ ; if y belongs to a range of I under ψ' , then there is x^* in I , depending on y , such that $y = \psi'(x^*)$. That is, $\varphi'(x^*) = 0$, and consequently φ attains its global maximum at x^* , because φ is concave. Hence, the conjugate is $\psi^*(y) = yx^* - \psi(x^*)$ and can be obtained via taking the derivative of ψ and setting it equal to zero. Applying this procedure componentwise in (5) yields $\psi'(\mathbf{g}_j) = \mathbf{p}_j/\mathbf{s}_j = \mathbf{f}_j$, whenever the additional assumptions of the theorem are satisfied. \square

EXAMPLE (Maximum Likelihood). In continuous version, this primal formulation can be traced back to Leonard [7] and Silverman [12]. The latter proposed

$$-\int g dP_n + \int e^g dx + \lambda \int (g^{(k)})^2 dx = \min_g! \quad (6)$$

using the third ($k = 3$) derivative to estimate the logarithm, g , of a density f , with the symbol P_n denoting the empirical probability supported by the datapoints; Gu [3] and others championed second ($k = 2$) derivative instead. The total variation penalty

$$\int |g^{(k)}| dx = \bigvee g^{(k-1)}$$

was considered by Koenker and Mizera [5, 6] for $k = 1, 2, 3$. Rufibach and Dümbgen [11] investigated maximum likelihood estimation of a log-concave density, which in our setting corresponds to $k = 2$ and the penalty

in the form of the non-positivity constraint on the second derivative (with no tuning parameter λ).

In the discrete setting, the k th derivative operator is replaced by an appropriate difference operator \mathbf{P} , and the evaluation operator \mathbf{L} and vector of weights \mathbf{w} by their typical instances described above. Since $\psi(x) = e^x$ is nondecreasing, (P) is equivalent to the unconstrained formulation (U), whose specific form is, for symmetric $J(\mathbf{u}) = \lambda \|\mathbf{u}\|_p^p$ and $p = 1, 2$,

$$-\mathbf{w}^\top \mathbf{L} \mathbf{g} + \mathbf{s}^\top \mathbf{e}^g + \lambda \|\mathbf{P} \mathbf{g}\|_p^p = \min_{\mathbf{g}} !, \quad (7)$$

where \mathbf{e}^g is understood componentwise. The additional assumptions of Theorem 2 are satisfied, so that indeed $\mathbf{f} = \mathbf{e}^g$, and

$$\psi^*(y) = \begin{cases} y \log y - y, & \text{for } y > 0, \\ 0, & \text{for } y = 0, \\ +\infty, & \text{otherwise.} \end{cases}$$

The feasibility requirement related to the fact that $\text{dom } \psi^* = [0, +\infty)$ independently enforces the nonnegativity constraint on \mathbf{f} . Silverman [12] showed, via an argument based on the specific properties of the exponential function, that the result of (6) is a probability density; the same conclusion follows, in the discrete setting, from our Theorems 1 and 2 for all formulations of the type (7). If the assumptions of Theorem 1 regarding \mathbf{P} , \mathbf{L} , and \mathbf{w} are satisfied, then the dual objective function

$$-\sum_j \mathbf{s}_j \mathbf{f}_j \log \mathbf{f}_j + \sum_j \mathbf{s}_j \mathbf{f}_j,$$

can be further simplified, because the second sum is equal to 1, a constant. The resulting dual of (7), cast in the minimization form, is, for $p = 1$,

$$\begin{aligned} \sum \mathbf{s}_j \mathbf{f}_j \log \mathbf{f}_j &= \min_{\mathbf{f}, \mathbf{e}} !, \\ \text{subject to } \quad \mathbf{S} \mathbf{f} &= \mathbf{L}^\top \mathbf{w} + \mathbf{P}^\top \mathbf{e}, \quad \mathbf{f} \succeq 0, \quad \text{and} \quad \|\mathbf{e}\|_\infty \leq \lambda, \end{aligned} \quad (8)$$

and for $p = 2$,

$$\begin{aligned} \sum \mathbf{s}_j \mathbf{f}_j \log \mathbf{f}_j + \frac{1}{4\lambda} \|\mathbf{e}\|_2^2 &= \min_{\mathbf{f}, \mathbf{e}} !, \\ \text{subject to } \quad \mathbf{S} \mathbf{f} &= \mathbf{L}^\top \mathbf{w} + \mathbf{P}^\top \mathbf{e}, \quad \text{and} \quad \mathbf{f} \succeq 0. \end{aligned} \quad (9)$$

The dual of the penalty-constrained version of the primal (7),

$$-\mathbf{w}^\top \mathbf{L} \mathbf{g} + \mathbf{s}^\top \mathbf{e}^g = \min_{\mathbf{g}} !, \quad \text{subject to} \quad \|\mathbf{P} \mathbf{g}\|_p \leq \Lambda, \quad (10)$$

is (p and q being conjugate)

$$\begin{aligned} \sum_j s_j f_j \log f_j + \Lambda \|e\|_q = \min_{f,e} !, \\ \text{subject to} \quad S f = L^T w + P^T e, \quad \text{and} \quad f \succeq 0. \end{aligned} \quad (11)$$

Finally, the dual of the shape-constrained formulation,

$$-w^T L g + s^T \Psi(g) = \min_g !, \quad \text{subject to} \quad P g \preceq 0 \quad (12)$$

(yielding log-concave f when P is a second-order difference operator), is

$$\begin{aligned} \sum_j s_j f_j \log f_j = \min_{f,e} !, \\ \text{subject to} \quad S f = L^T w + P^T e, \quad f \succeq 0, \quad \text{and} \quad e \preceq 0. \end{aligned} \quad (13)$$

The essence of all the dual variants is the maximization of the Shannon entropy of f , or, equivalently, the minimization of the Kullback-Leibler divergence

$$\mathcal{K}(f, \sigma^{-1}) = \sum_j s_j f_j \log \frac{f_j}{\sigma^{-1}} = \sum_j s_j f_j \log \frac{s_j f_j}{s_j \sigma^{-1}} = \sum_j s_j f_j \log f_j + \log \sigma,$$

where $\sigma^{-1} = (\sum_j s_j)^{-1}$ can be viewed as a discretization of the uniform density.

The dual formulation of the penalized likelihood problem as a maximum entropy problem can be generalized by replacing the Shannon entropy term by some of the Rényi entropies, indexed by a parameter $\alpha > 0$; similarly to the Kullback-Leibler case, the appropriate minimum divergence interpretations follow. Formally, Rényi's entropies include the Shannon one for $\alpha = 1$; the Rényi [9] entropy with exponent $\alpha \neq 1$ is defined as $(1 - \alpha)^{-1} \log(s^T f^\alpha)$, where f^α is interpreted componentwise. The maximization of this function is equivalent to the maximization of $-\text{sign}(\alpha - 1)s^T f^\alpha$ or, equivalently, $-\text{sign}(\alpha - 1)s^T f^\alpha / \alpha$.

Let ψ_p be a function equal to x^p/p for $x \geq 0$ and to 0 for $x < 0$. The conjugate, ψ_p^* , of ψ_p is for $p > 1$ equal to y^q/q for $y \geq 0$ (p and q conjugate), and to $+\infty$ otherwise. Note that ψ_p is nondecreasing, hence (P) is equivalent to (U) whenever $\psi = \psi_p$.

EXAMPLE (Minimum Pearson χ^2). The special case of the Rényi system for $\alpha = 2$ yields $\psi(x) = \psi_2$ and $\psi_2^* = y^2/2$ for $y \geq 0$. The dual is obtained by replacing the entropy term $\sum_j s_j f_j \log f_j$ in the objective function of (8), (9), (11), and (13) by $s^T f^2$, and eliminating the redundant constant in the objective.

The corresponding primal results from replacing $\sum_j s_j e^{g_j}$ in (7), (10), and (12) by $\sum_j s_j \psi_2(g_j)$. Minimizing the dual (and in this case also primal) objective is

equivalent to minimizing the χ^2 -divergence

$$\chi^2(\mathbf{f}, \sigma^{-1}) = \sum_j s_j \frac{(\mathbf{f}_j - \sigma^{-1})^2}{\sigma^{-1}} = \sum_j \frac{(\mathbf{s}_j \mathbf{f}_j - \mathbf{s}_j \sigma^{-1})^2}{\mathbf{s}_j \sigma^{-1}} = \sigma \left(\sum_j \mathbf{s}_j \mathbf{f}_j^2 \right) - 1.$$

If instead of ψ_2 we consider $\psi(x) = (1/2)x^2$ for all x , we can cast both primal and dual in a quadratic programming form. However, the correct primal has to be written in the constrained form (P) now, because ψ is no longer monotone. In particular, the correct formulation for the setting corresponding to (7) is

$$-\mathbf{w}^\top \mathbf{L} \mathbf{h} + \frac{1}{2} \mathbf{s}^\top \mathbf{g}^2 + \lambda \|\mathbf{P} \mathbf{h}\|_p^p = \min_{\mathbf{g}, \mathbf{h}}!, \quad \text{subject to } \mathbf{h} \preceq \mathbf{g}.$$

In all variants, both primal and dual yield directly $\mathbf{f} = \mathbf{g}$, because $\psi'(x) = x$.

EXAMPLE. Another special case of the Rényi scheme, with $\alpha = 3/2$, puts $\mathbf{s}^\top \mathbf{f}^{3/2}$ into the objective function of (8), (9), (11), and (13). For the primal, we may take either $\sum_j \mathbf{s}_j \psi^3(\mathbf{g}_j)$ in (7), (10), and (12); or we may use $\psi(x) = (1/3)|x|^3$ instead, leaving the dual unchanged, but making the primal constrained; for instance, the formulation (7) becomes

$$-\mathbf{w}^\top \mathbf{L} \mathbf{h} + \frac{1}{3} \mathbf{s}^\top \mathbf{g}^3 + \lambda \|\mathbf{P} \mathbf{h}\|_p^p = \min_{\mathbf{g}, \mathbf{h}}!, \quad \text{subject to } \mathbf{h} \preceq \mathbf{g}.$$

Due to the fact that $\mathbf{f} = \mathbf{g}^2$ in any of these variants, we could nickname this example “Silverman for Good”. Apart from the additional middle term, the objective function differs from the original proposal of Good [2] also in the first term; ours is not based on the logarithm of the square root of the estimated density, but on the square root itself. It would be interesting to know whether there is any Bayesian justification for such an approach, whether in “mufti” or “full regalia”. In any case, the primal formulation yields a square root of a probability density, a “rootogram” in Tukey’s terminology.

EXAMPLE (Minimum Hellinger). Another example from the Rényi system, with $\alpha = 1/2$, sets $\psi(x) = -1/x$, for $x < 0$ and $+\infty$ elsewhere. The conjugate is $\psi^*(y) = -2\sqrt{y}$, for $y \geq 0$, and ∞ elsewhere. The dual (for $p = 1$) has, in the minimization form and after the elimination of the redundant constant, $-\mathbf{s}^\top \sqrt{\mathbf{f}}$ in the objective of (8), (9), (11), and (13); $\sqrt{\mathbf{f}}$ is again applied componentwise. The dual objective minimizes the Hellinger distance

$$\begin{aligned} \mathcal{H}(\mathbf{f}, \sigma^{-1}) &= \sum_j s_j \left(\sqrt{\mathbf{f}_j} - \sqrt{\sigma^{-1}} \right)^2 \\ &= \sum_j \left(\sqrt{\mathbf{s}_j \mathbf{f}_j} - \sqrt{\mathbf{s}_j \sigma^{-1}} \right)^2 \\ &= 2 - 2\sqrt{\sigma^{-1}} \sum_j \mathbf{s}_j \sqrt{\mathbf{f}_j}. \end{aligned}$$

Since ψ is nondecreasing, the primal can be cast in its unconstrained version (U), just replacing the $\sum_j s_j e^{g_j}$ term in (7), (10), or (12) by $-\mathbf{s}^\top \mathbf{g}^{-1}$, where \mathbf{g}^{-1} is the componentwise reciprocal value of \mathbf{g} ; however, the domain restriction for ψ has to be included as a feasibility constraint. The resulting primal analog of (7) is

$$-\mathbf{w}^\top \mathbf{L} \mathbf{g} - \mathbf{s}^\top \mathbf{g}^{-1} + \lambda \|\mathbf{P} \mathbf{g}\|_1 = \min_{\mathbf{g}}!, \quad \text{subject to} \quad \mathbf{g} \preceq 0.$$

For symmetric penalties, it is more convenient to recast the primal in terms of $\mathbf{h} = -\mathbf{g}$:

$$\mathbf{w}^\top \mathbf{L} \mathbf{h} + \mathbf{s}^\top \mathbf{h}^{-1} + \lambda \|\mathbf{P} \mathbf{h}\|_1 = \min_{\mathbf{h}}!, \quad \text{subject to} \quad \mathbf{h} \succeq 0.$$

The estimated density $\mathbf{f} = 1/\mathbf{g}^2 = 1/\mathbf{h}^2$; hence \mathbf{h} could be called, in the Tukey spirit, a “rootosparsity”, and \mathbf{g} , being negative, a “hanging rootosparsity”.

In our implementations, we observed that numerical performance may be improved by adding the (theoretically redundant) nonnegativity constraint $\mathbf{f} \succeq 0$ also in the primal formulation. However, this is rather an unimportant detail, because dual formulations always ran significantly faster and were more numerically stable than their primal counterparts.

EXAMPLE (Maximum empirical likelihood). The limiting variant of the Rényi system for $\alpha = 0$ is $\psi(x) = -1/2 - \log(-x)$ for $x < 0$, and $+\infty$ otherwise. The dual puts $-\mathbf{s}^\top \log \mathbf{f}$ into the objective function of (8), (9), (11), and (13), while the primal (unconstrained, but with a feasibility constraint) puts $-\mathbf{s}^\top \log(-\mathbf{g})$ into in (7), (10), or (12). For instance, recasting (7) in terms of $\mathbf{h} = -\mathbf{g}$ gives

$$\mathbf{w}^\top \mathbf{L} \mathbf{h} - \mathbf{s}^\top \log \mathbf{h} + \lambda \|\mathbf{P} \mathbf{h}\|_1 = \min_{\mathbf{h}}!, \quad \text{subject to} \quad \mathbf{h} \succ 0.$$

The dual objective is equivalent to the reversed Kullback-Leibler divergence

$$\begin{aligned} \mathcal{K}(\sigma^{-1}, \mathbf{f}) &= \sum_j s_j \sigma^{-1} \log \frac{\sigma^{-1}}{f_j} \\ &= \sum_j s_j \sigma^{-1} \log \frac{s_j \sigma^{-1}}{s_j f_j} \\ &= -\sigma^{-1} \sum_j s_j \log f_j - \log \sigma, \end{aligned}$$

whose minimizing is known to be equivalent to maximizing the so-called empirical likelihood in the sense of Owen [8]; see Hall and Presnell [4]. The salient feature of this example is that the function \mathbf{h} penalized in the primal is “sparsity”, the reciprocal of the estimated density \mathbf{f} .

EXAMPLE. One can easily come to the idea to employ the popular and simple total variation distance in the minimum divergence formulation, choosing the dual objective to minimize

$$V(\mathbf{f}, \sigma^{-1}) = \sum_j s_j |\mathbf{f}_j - \sigma^{-1}|,$$

a formulation leading to a linear programming problem. However, the primal in this case would involve a function $\psi(x)$ equal to x/σ for x in the interval $[-1, 1]$, and $+\infty$ elsewhere. This indicates difficulties and likely explains the strange results we observed in our implementations.

Acknowledgement. We are indebted to Xuming He for valuable discussions, to Erling D. Andersen for creating MOSEK, a convex optimization toolbox for MATLAB, and to the anonymous referee for the careful reading of the manuscript.

REFERENCES

- [1] BOYD, S.—VANDENBERGHE, L.: *Convex Optimization*. Cambridge University Press, Cambridge, 2004.
- [2] GOOD, I. J.: *A nonparametric roughness penalty for probability densities*, *Nature* **229** (1971), 29–30.
- [3] GU, C.: *Smoothing Spline ANOVA Models*. Springer-Verlag, New York, 2002.
- [4] HALL, P.—PRESNELL, B.: *Density estimation under constraints*, *J. Comput. Graph. Statist.* **8** (1999), 259–277.
- [5] KOENKER, R.—MIZERA, I.: *The alter egos of the regularized maximum likelihood density estimators: deregularized maximum-entropy, Shannon, Renyi, Simpson, Gini, and stretched strings*. In: *Proceedings of 7th Prague Symposium on Asymptotic Statistics and 15th Prague Conference on Information Theory, Statistical Decision Functions and Random Processes—Prague Stochastics '06* (M. Hušková, M. Janžura, eds.), Prague, August 21–25, 2006, Matfyzpress, Prague, 2006, pp. 145–157.
- [6] KOENKER, R.—MIZERA, I.: *Density estimation by total variation regularization*. In: *Advances in statistical modeling and inference, Essays in honor of Kjell A. Doksum* (V. Nair, ed.), World Scientific, Singapore, 2006.
- [7] LEONARD, T.: *Density estimation, stochastic processes and prior information*, *J. Roy Stat. Soc.* **40** (1978), 113–132.
- [8] OWEN, A. B.: *Empirical Likelihood*. Chapman & Hall/CRC, Boca Raton, 2001.
- [9] RÉNYI, A.: *On measures of entropy and information*. In: *Proc. 4th Berkeley Symp. Math. Stat. Probab., Vol. 1: Contributions to the Theory of Statistics* (J. Neyman, ed.), University of California June 20–July 30, 1960, University of California Press, Berkeley, 1961, pp. 547–561.
- [10] ROCKAFELLAR, R. T.: *Convex Analysis*. Princeton University Press, Princeton, 1970.

- [11] RUFIBACH, K.—DÜMBGEN, L.: *Maximum likelihood estimation of a log-concave density: basic properties and uniform consistency*, preprint arXiv:0709.0334.
- [12] SILVERMAN, B. W.: *On the estimation of a probability density function by the maximum penalized likelihood method*, Ann. Statist. **10** (1982), 795–810.

Received October 25, 2006

Roger Koenker
University of Illinois at Urbana-Champaign
Departments of Economics and Statistics
Champaign, Illinois, 61620
USA
E-mail: rkoenker@uiuc.edu

Ivan Mizera
University of Alberta
Department of Mathematical and Statistical Sciences
CAB 632, Edmonton, Alberta, T6J 0Z2
CANADA
E-mail: mizera@stat.ualberta.ca