

Maximum Entropy on the Mean and the Cramér Rate Function in Statistical Estimation and Inverse Problems: Properties, Models, and Algorithms*

Yakov Vaisbourd[†], Rustum Choksi[†], Ariel Goodwin[†], Tim Hoheisel[†], and Carola-Bibiane Schönlieb[‡]

Abstract. We explore a method of statistical estimation called *Maximum Entropy on the Mean* (MEM) which is based on an information-driven criterion that quantifies the compliance of a given point with a reference prior probability measure. At the core of this approach lies the *MEM function* which is a partial minimization of the Kullback-Leibler divergence over a linear constraint. In many cases, it is known that this function admits a simpler representation (known as the *Cramér rate function*). Via the connection to exponential families of probability distributions, we study general conditions under which this representation holds. We then address how the associated *MEM estimator* gives rise to a wide class of MEM-based regularized linear models for solving inverse problems. Finally, we propose an algorithmic framework to solve these problems efficiently based on the Bregman proximal gradient method, alongside proximal operators for commonly used reference distributions. The article is complemented by a software package for experimentation and exploration of the MEM approach in applications.

Key words. Maximum Entropy on the Mean, Statistical Estimation, Cramér Rate Function, Kullback-Leibler Divergence, Prior Distribution, Regularization, Linear Inverse Problems, Bregman Proximal Gradient, Convex Duality, Large Deviations.

MSC codes. 49M27, 29M29, 60F10, 62B10, 62H12, 90C25, 90C46

1. Introduction. Many models for modern applications in various disciplines are based on some form of *statistical estimation*, for example the very common *maximum likelihood* (ML) principle. In this study, we consider an alternative approach known as the *maximum entropy on the mean* (MEM). At its core lies the MEM function κ_P induced by some *reference distribution* P and defined as

$$\kappa_P(y) := \inf \{ \text{KL}(Q|P) : \mathbb{E}_Q = y, Q \in \mathcal{P}(\Omega) \},$$

where $P(\Omega)$ stands for the set of probability measures on $\Omega \subseteq \mathbb{R}^d$, \mathbb{E}_Q is the expected value of $Q \in P(\Omega)$ and $\text{KL}(Q|P)$ stands for the Kullback-Leibler (KL) divergence of Q with respect to P [38] (see Section 2 for precise definitions). Thus, the MEM modeling paradigm stems from the principle of minimum discrimination information [37] which generalizes the well-known principle of maximum entropy [36]. In the context of information theory [24], the argmin of $\kappa_P(y)$ is often referred to as the *information projection* of P onto the set $\{Q \in P(\Omega) : \mathbb{E}_Q = y\}$, the *closest* member of the set to P .

Various forms and interpretations of MEM have been studied (see for example, [26, 30, 31, 32, 34, 39, 40]) and found applications in various disciplines, including earth sciences [29, 42, 43, 45, 52], and medical imaging [1, 19, 22, 33, 35]. A version of the MEM method

*Submitted to the editors DATE.

[†]Department of Mathematics and Statistics, McGill University

[‡]Department of Applied Mathematics and Theoretical Physics, University of Cambridge.

39 was recently explored for blind deblurring of images possessing some form of fixed symbology
 40 (for example, in barcodes) [47, 46]. There one exploited the ability of of the MEM framework to
 41 facilitate the incorporation of nonlinear constraints via the introduction of a prior distribution.

42 Despite its many interesting properties in both theory and applications, the MEM method-
 43 ology has yet to find its place as a mainstream tool for statistical estimation, particularly as it
 44 pertains to solving inverse problems. One factor that might have contributed to this centers
 45 on the practical issue that there are no dedicated optimization algorithms designed to tackle
 46 models based on the MEM methodology. Indeed, the MEM function is defined by means of
 47 an infinite-dimensional optimization problem. Previous attempts to solve models involving
 48 the MEM function relied on its finite-dimensional dual problem. To the best of the authors'
 49 knowledge, there are no dedicated optimization algorithms designed to tackle models based
 50 on the MEM methodology. Therefore, any researcher or practitioner wishing to employ the
 51 MEM framework must first overcome a notable barrier of deriving an appropriate optimization
 52 algorithm for its solution. In this work, our goal is to fill in this gap, providing an accessible
 53 gate to the MEM methodology.

54 Our approach is based on the fundamental work by Brown [18, Chapter 6] and comple-
 55 ments [39] by first proving the equivalence of the MEM function to the *Cramér's rate* func-
 56 tion, mostly known from its role in *large deviation theory*. Cramér's rate function is defined
 57 by means of a finite-dimensional optimization problem as it is simply the convex conjugate of
 58 the log-normalizer (aka the cumulant generating function) of the reference distribution P . In
 59 many cases (i.e., choices of P) it admits a closed form expression while in others it can still
 60 be evaluated efficiently. The connection between these seemingly different functions is well
 61 established in the large deviations [27], statistics [18], and information theory [39] literature.
 62 Nonetheless, various assumptions imposed in the aforementioned works limit the scope of ex-
 63 isting results. Employing the framework of exponential families of probability distributions
 64 [18], we establish the equivalence between the two functions under very mild and natural con-
 65 ditions, allowing us to cover many distributions of practical interest. Thus, models involving
 66 MEM functions can be explicitly stated using the corresponding Cramér functions.

67 Central to our study is *the MEM estimator* which is shown to be well defined under very
 68 mild conditions. We further recall an insightful connection between the MEM and ML esti-
 69 mators as presented in [18] for the case of a reference distribution from an exponential family.
 70 As with the ML counterpart, the MEM estimator has vast applications, and hence we restrict
 71 the remainder of the paper to a wide class of regularized linear models for solving inverse
 72 problems. Each model in this class involves two MEM functions, one in the role of a fidelity
 73 term and another as a regularizer (comparable to the *maximum a priori (MAP) estimation*
 74 framework which extends ML). Let us provide an example: given a measurement matrix
 75 $A \in \mathbb{R}^{m \times d}$, an observation vector $\hat{y} \in \mathbb{R}^m$ and an additional vector $p \in [0, 1]^d$ representing
 76 some prior knowledge, the following optimization problem

$$77 \quad \min \left\{ \underbrace{\frac{1}{2} \|Ax - \hat{y}\|_2^2}_{\text{Fidelity}} + \underbrace{\sum_{i=1}^d \left[x_i \log \left(\frac{x_i}{p_i} \right) + (1 - x_i) \log \left(\frac{1 - x_i}{1 - p_i} \right) \right]}_{\text{Regularization}} : x \in [0, 1]^d \right\},$$

78
79

80 fits the MEM framework with normal (Gaussian) and Bernoulli reference distributions of the

81 fidelity and regularization terms, respectively. Other choices of reference distributions will
 82 lead to additional models that admit similar additive composite structure. Moreover, the
 83 closed form expressions of the two functions in our example follow from the definition of
 84 Cramér’s rate function. In models of these forms, concrete expressions and structures with
 85 distinct geometry can be exploited to customize appropriate optimization strategies. Here we
 86 highlight the class of *Bregman proximal gradient* (BPG) methods as an especially suitable
 87 choice for this family of models. Nevertheless, other methods are also viable alternatives; for
 88 example, adaptive and scaled, accelerated variants and dual decomposition methods which
 89 are defined by means of the same operators developed here.

90 Our overall aim is to provide a self-contained, mathematically sound toolbox for working
 91 with the MEM methodology for a wide variety of models. For this reason, we provide a
 92 comprehensive list of Cramér functions and operators used in the algorithms, and complement
 93 it with a software package. We believe this sets the basis for (and hopefully triggers) further
 94 experimentation and exploration of the MEM approach in contemporary applications.

95 The paper is organized as follows. In [Section 2](#), we recall some concepts and preliminary
 96 results from convex analysis and probability theory which will be used in this work. In
 97 [Section 3](#), we study the MEM and Cramér rate functions and establish the equivalence between
 98 the two under very mild and natural conditions. This allows us to use the accessible definition
 99 of the Cramér function and derive tractable expressions for a wide class of possible reference
 100 distributions which closes this section (see [Table 1](#)). [Section 4](#) is devoted to the MEM models
 101 considered in this work, and in [Section 5](#), we present the algorithms for solving such models.
 102 We end with a few concrete examples of problems and corresponding algorithms crafted from
 103 the operators derived in this work. An appendix provides deferred proofs and the details of a
 104 variety of Cramér rate function computations.

105 2. Preliminaries.

106 **2.1. Convex Analysis.** We recall here some definitions and results from convex analysis.
 107 Further details and proofs can be found in various textbooks such as [\[9, 11, 48\]](#).

108 The *affine hull* of a set $S \subseteq \mathbb{R}^d$ is the smallest affine subspace containing S . For any point
 109 $y \in S$, we have the following relation

$$110 \quad (2.1) \quad \text{aff } S = y + \text{span}(S - y),$$

112 where $\text{span } S$ stands for the linear hull of S . The dimension of $\text{aff } S$ is defined as $\dim(\text{aff } S) :=$
 113 $\dim(\text{span}(S - y))$. The interior, closure and boundary of a set are denoted as $\text{int } S$, $\text{cl } S$ and
 114 $\text{bd } S$, respectively.

115 The (Fenchel) conjugate of $\psi : \mathbb{R}^d \rightarrow [-\infty, \infty]$ is defined as

$$116 \quad \psi^*(y) := \sup\{\langle y, x \rangle - \psi(x) : x \in \mathbb{R}^d\}.$$

118 The function ψ is proper if $\psi(x) > -\infty$ for all $x \in \mathbb{R}^d$ and $\text{dom } \psi := \{x \in \mathbb{R}^d : \psi(x) < \infty\} \neq \emptyset$.
 119 In addition, ψ is closed, if its epigraph $\{(x, \alpha) \in \mathbb{R}^d \times \mathbb{R} : \psi(x) \leq \alpha\}$ is a closed set.

120 If ψ is proper and convex then ψ^* is closed, proper and convex. For a proper function
 121 $\psi : \mathbb{R}^d \rightarrow (-\infty, +\infty]$, the *Fenchel-Young inequality* states that $\psi(x) + \psi^*(y) \geq \langle y, x \rangle$. If ψ is

122 proper, closed and convex then we obtain that [11, Theorem 4.20]

$$123 \quad (2.2) \quad \psi(x) + \psi^*(y) = \langle y, x \rangle \iff y \in \partial\psi(x) \iff x \in \partial\psi^*(y),$$

125 where $\partial\psi(x) := \{g \in \mathbb{R}^d : \psi(y) \geq \psi(x) + \langle g, y - x \rangle \ (y \in \mathbb{R}^d)\}$ is the *subdifferential* of ψ at
126 $x \in \mathbb{R}^d$.

127 The *indicator function* of a set $S \subseteq \mathbb{R}^d$ is denoted by δ_S and defined as $\delta_S(x) = 0$ if
128 $x \in S$ and $\delta_S(x) = +\infty$ otherwise. Its convex conjugate is known as the *support function*
129 $\sigma_S(y) := \delta_S^*(y) = \sup\{\langle y, x \rangle : x \in S\}$.

130 **Definition 2.1 (Essential smoothness and Legendre type).** *Let $\psi : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ be*
131 *proper and convex. Then, ψ is called essentially smooth if it satisfies the following conditions:*

- 132 1. $\text{int}(\text{dom } \psi) \neq \emptyset$;
- 133 2. ψ is differentiable on $\text{int}(\text{dom } \psi)$;
- 134 3. $\|\nabla\psi(x^k)\| \rightarrow \infty$ for any sequence $\{x^k \in \text{int}(\text{dom } \psi)\}_{k \in \mathbb{N}} \rightarrow \bar{x} \in \text{bd}(\text{dom } \psi)$.

135 *The last condition listed above is called steepness. An essentially smooth function ψ is said*
136 *to be of Legendre type if it is strictly convex on $\text{int}(\text{dom } \psi)$.*

137 For $\psi : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ closed and of Legendre type, the following hold [48, Theorem 26.5]:

- 138 1. ψ^* is of Legendre type.
- 139 2. $\nabla\psi : \text{int}(\text{dom } \psi) \rightarrow \text{int}(\text{dom } \psi^*)$ is a bijection with $(\nabla\psi)^{-1} = \nabla\psi^*$.

140 The *Bregman distance* induced by a function ψ of Legendre type is defined as [17]

$$141 \quad D_\psi(y, x) = \psi(y) - \psi(x) - \langle \nabla\psi(x), y - x \rangle \quad (x \in \text{int}(\text{dom } \psi), y \in \text{dom } \psi).$$

143 For any $(x, y) \in \text{int}(\text{dom } \psi) \times \text{dom } \psi$, the Bregman distance is nonnegative $D_\psi(y, x) \geq 0$, and
144 equality holds if and only if $x = y$ due to strict convexity of ψ [17]. However, in general, D_ψ
145 is not symmetric, unless $\psi = (1/2)\|\cdot\|^2$ [7, Lemma 3.16]. The Bregman distance induced by
146 a function ψ of Legendre type satisfies the following additional properties [8, Theorem 3.7]:
147 For any $x, y \in \text{int}(\text{dom } \psi)$ it holds that

$$148 \quad (2.3) \quad D_\psi(y, x) = D_{\psi^*}(\nabla\psi(x), \nabla\psi(y)).$$

150 The Bregman distance is strictly convex with respect to its first argument. Moreover, for two
151 functions ψ_1 and ψ_2 differentiable at $x \in \text{int}(\text{dom } \psi_1) \cap \text{int}(\text{dom } \psi_2)$

$$152 \quad (2.4) \quad D_{\alpha\psi_1 + \beta\psi_2}(y, x) = \alpha D_{\psi_1}(y, x) + \beta D_{\psi_2}(y, x) \quad (y \in \text{dom } \psi_1 \cap \text{dom } \psi_2, \alpha, \beta \in \mathbb{R}).$$

154 **2.2. Probability Theory and Exponential Families.** We recall some concepts from prob-
155 ability theory with an emphasis on exponential families. For further detail, see e.g. [4, 18].

156 Let $\mathcal{M}(\Omega)$ be the set of σ -finite measures defined over a measurable space (Ω, Σ) where
157 $\Omega \subseteq \mathbb{R}^d$ and Σ is a σ -algebra on Ω . The *support* of ρ , namely the minimal closed measurable
158 set $A \in \Sigma$ such that $\rho(\Omega \setminus A) = 0$, is denoted by Ω_ρ . We denote by $\Omega_\rho^{cc} := \text{cl}(\text{conv } \Omega_\rho)$
159 the closure of the convex hull of the support Ω_ρ , which is known as the *convex support* of ρ .
160 Recall further that, if μ is another measure defined over (Ω, Σ) , then μ is *absolutely continuous*
161 with respect to ρ (denoted by $\mu \ll \rho$) if for every $A \in \Sigma$ such that $\rho(A) = 0$ it holds that
162 $\mu(A) = 0$. In this case, the *Radon-Nikodym derivative* is the unique function $h = \frac{d\mu}{d\rho}$ such that

163 $\mu(A) = \int_A h d\rho$ for any $A \in \Sigma$. For a measurable space (Ω, Σ) we denote by $\nu \in \mathcal{M}(\Omega)$ the
 164 *dominating measure*. Throughout, we restrict ourselves to two scenarios: either $\Omega = \mathbb{R}^d$ and
 165 ν is the Lebesgue measure or Ω is a countable subset of \mathbb{R}^d and ν is the counting measure.
 166 Let $\mathcal{P}(\Omega)$ be the set of probability measures defined over Ω and absolutely continuous with
 167 respect to ν . We emphasize that for $P \in \mathcal{P}(\Omega)$ the support Ω_P might be a proper subset of
 168 Ω , and thus there is no loss of generality in our setting even when $\Omega = \mathbb{R}^d$. Furthermore,
 169 for any set $A \subseteq \mathbb{R}^d$ the expression $P(A)$ should be understood as $P(A \cap \Omega)$. For $P \in \mathcal{P}(\Omega)$,
 170 the Radon-Nikodym derivative $f_P := \frac{dP}{d\nu}$ is either a probability density or mass function,
 171 depending on the set Ω . In both cases, we will refer to f_P as the density of the distribution.¹
 172 The *expected value* (if it exists) and *moment generating function* of $P \in \mathcal{P}(\Omega)$ are given by

$$173 \quad \mathbb{E}_P := \int_{\Omega} y dP(y) \in \mathbb{R}^d \quad \text{and} \quad M_P[\theta] := \int_{\Omega} \exp(\langle \cdot, \theta \rangle) dP,$$

174
 175 respectively. For $P \in \mathcal{M}(\Omega)$ absolutely continuous with respect to ν , we define

$$176 \quad \Theta_P := \left\{ \theta \in \mathbb{R}^d : \int_{\Omega} \exp(\langle \cdot, \theta \rangle) dP < \infty \right\},$$

177
 178 and consider the function $\psi_P : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ given by

$$179 \quad (2.5) \quad \psi_P(\theta) := \begin{cases} \log \int_{\Omega} \exp(\langle \cdot, \theta \rangle) dP, & \theta \in \Theta_P, \\ +\infty, & \theta \notin \Theta_P. \end{cases}$$

180
 181 Then $\mathcal{F}_P := \{f_{P_\theta}(y) := \exp(\langle y, \theta \rangle - \psi_P(\theta)) : \theta \in \Theta_P\}$, is a *standard exponential family* gener-
 182 ated by P . Note that, the probability measure P_θ satisfying $dP_\theta = f_{P_\theta} dP$ is, by construction, a
 183 probability measure such that P_θ and P are mutually absolutely continuous, hence $\Omega_{P_\theta} = \Omega_P$
 184 for all $\theta \in \Theta_P$ [4, Section 8.1]. The function ψ_P is called the *log-normalizer* (also known as
 185 the *log-partition* or *log-Laplace transform* of P). The vector $\theta \in \mathbb{R}^d$ is known as the *natural*
 186 *parameter* and the set $\Theta_P = \text{dom } \psi_P$ is called the *natural parameter space*.²

187 The following results summarize some well-known properties of the log-normalizer ψ_P .

188 **Proposition 2.2 (Convexity, [18, Theorem 1.13]).** *Let \mathcal{F}_P be an exponential family generated*
 189 *by $P \in \mathcal{M}(\Omega)$. Then, the natural parameter space Θ_P is a convex set and the log-normalizer*
 190 *function $\psi_P : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ is closed, proper and convex.*

191 **Proposition 2.3 (Differentiability, [18, Theorem 2.2, Corollary 2.3]).** *Let \mathcal{F}_P be an exponential*
 192 *family generated by $P \in \mathcal{M}(\Omega)$ and let $\theta \in \text{int } \Theta_P$. Then, the log normalizer $\psi_P : \mathbb{R}^d \rightarrow$
 193 $(-\infty, +\infty]$ is infinitely differentiable at θ and it holds that $\nabla \psi_P(\theta) = \mathbb{E}_{P_\theta}$.*

194 The dimension of a convex set $S \subseteq \mathbb{R}^d$, denoted by $\dim S$, is equal to the affine dimension
 195 of $\text{aff } S$. We assume that the exponential family generated by $P \in \mathcal{M}(\Omega)$ is *minimal*, i.e.,
 196 $\dim \Theta_P = \dim \Omega_P^{\text{cc}} = d$ or, equivalently, $\text{int } \Theta_P \neq \emptyset$ and $\text{int } \Omega_P^{\text{cc}} \neq \emptyset$. This is not restrictive as a
 197 non-minimal exponential family can be always reduced to a minimal form [18, Theorem 1.9].
 198 The following result strengthens Proposition 2.2 for minimal exponential families.

¹We will interchangeably refer to $P \in \mathcal{P}(\Omega)$ as either a distribution or measure.

²It is possible to define the exponential family \mathcal{F}_P over a subset of the natural parameter space [18, Definition 1.1], but this is not needed for our study.

199 **Proposition 2.4** (Strict convexity, [18, Theorem 1.13]). Let \mathcal{F}_P be a minimal exponential
 200 family generated by $P \in \mathcal{M}(\Omega)$. Then, the log-normalizer function $\psi_P : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ is
 201 strictly convex over Θ_P .

202 If the log-normalizer ψ_P is essentially smooth (or 'steep' in the exponential family terminology,
 203 see, e.g., [4, Theorem 5.27] and [18, Definition 3.2]), we say that the exponential family \mathcal{F}_P is
 204 steep. This condition is automatically satisfied when Θ_P is open [4, Theorem 8.2]. While most
 205 exponential families encountered in practice have this property, there are relevant cases when
 206 this assumption is too restrictive (e.g., [18, Example 3.4]). Thus, in order to cover all examples
 207 provided in this work, we will assume that the exponential family is steep. Summarizing the
 208 above discussion and recalling Definition 2.1 we have the following corollary.

209 **Corollary 2.5.** Let \mathcal{F}_P be a minimal and steep exponential family generated by $P \in \mathcal{M}(\Omega)$.
 210 Then, the log normalizer function ψ_P is of Legendre type.

211 From the last corollary we can see that $\nabla\psi_P$ forms a bijection between $\text{int}(\text{dom } \psi_P) = \text{int } \Theta_P$
 212 and $\text{int}(\text{dom } \psi_P^*)$. This relation, provides a dual representation of the log-normalizer ψ_P
 213 and, consequently, the distribution in question. The so-called *mean value parametrization*
 214 is obtained by applying a change of variables where the natural parameter θ is replaced by
 215 $\mu \in \mathbb{R}^d$ such that $\mu = \mathbb{E}_{P_\theta} = \nabla\psi_P(\theta)$, i.e., $\theta = \nabla\psi_P^*(\mu)$.

216 The *Kullback-Leibler (KL) divergence* (also known as the relative entropy) of a probability
 217 measure $Q \in \mathcal{P}(\Omega)$ with respect to $P \in \mathcal{P}(\Omega)$ is given by (see [38])

$$218 \quad \text{KL}(Q|P) := \begin{cases} \int_{\Omega} \log \left(\frac{dQ}{dP} \right) dQ, & Q \ll P, \\ +\infty, & \text{otherwise.} \end{cases}$$

220 It holds that $\text{KL}(Q|P) \geq 0$ with equality if and only if $Q = P$ [38, Lemma 3.1]. Thus, the
 221 Kullback-Leibler information quantifies the dissimilarity between two probability measures.
 222 We note that, in general, $\text{KL}(Q|P)$ is not symmetric. Furthermore, $\text{KL}(Q|P)$ is jointly convex
 223 in $(Q|P)$. We record a special case for which the KL divergence is of particular interest.

224 **Remark 2.6** (Kullback-Leibler divergence for exponential family). Let \mathcal{F}_P be an exponential
 225 family generated by $P \in \mathcal{M}(\Omega)$. Let $\theta_1 \in \Theta_P$ and $\theta_2 \in \text{int } \Theta_P$, thus for $i = 1, 2$ we have that
 226 $f_{P_{\theta_i}} \in \mathcal{F}_P$. In this case, the KL divergence between the two measures $P_{\theta_i} \in \mathcal{P}(\Omega)$ such that
 227 $dP_{\theta_i} := f_{P_{\theta_i}} dP$ ($i = 1, 2$) satisfies $\text{KL}(P_{\theta_2}|P_{\theta_1}) = D_{\psi_P}(\theta_1, \theta_2)$ [18, Proposition 6.3]. \diamond

228 **3. Maximum entropy on the mean and Cramér's rate function.** For $y \in \mathbb{R}^d$, the density

$$229 \quad (3.1) \quad f_P(y) := \frac{dP}{d\nu}(y)$$

231 provides an indication of the likelihood of y under the distribution $P \in \mathcal{P}(\Omega)$. The method of
 232 *Maximum Entropy on the Mean* (MEM) suggests an alternative, information driven function
 233 $\kappa_P : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ given by

$$234 \quad (3.2) \quad \kappa_P(y) := \inf \{ \text{KL}(Q|P) : \mathbb{E}_Q = y, Q \in \mathcal{P}(\Omega) \}.$$

236 Here, κ_P measures how y complies with the distribution P , by seeking a distribution Q
 237 with expected value y that minimizes $\text{KL}(\cdot|P)$. The distance, in terms of the KL divergence

238 (the information gain) between the resulting and the original distributions quantifies the
 239 compliance of y with P . We will refer to κ_P as the *MEM function* and to P as the *reference*
 240 *distribution*. Since $\text{KL}(Q|P) \geq 0$ and $\text{KL}(Q|P) = 0$ if and only if $Q = P$, we find that the
 241 MEM function satisfies $\kappa_P(y) \geq 0$ for any $y \in \mathbb{R}^d$ and $\kappa_P(y) = 0$ if and only if $y = \mathbb{E}_P$.

242 In most cases of interest, the MEM function admits an alternative representation which
 243 sheds light on many of its additional properties (cf. [Theorem 3.10](#)). More precisely, under
 244 suitable conditions (cf. [Theorem 3.8](#)), the MEM function coincides with the *Cramér rate*
 245 *function* [25], to which we turn now. For a given reference distribution $P \in \mathcal{P}(\Omega)$, recall the
 246 log-normalizer previously defined for a general measure in (2.5):

$$247 \quad \psi_P(\theta) := \log M_P[\theta] = \log \int_{\Omega} \exp(\langle \cdot, \theta \rangle) dP.$$

249 In the context of probability measures P , ψ_P is often known as the *cumulant generating*
 250 *function*. The *Cramér rate function* ψ_P^* associated with P is the conjugate of ψ_P , that is,

$$251 \quad \psi_P^*(y) = \sup\{\langle y, \theta \rangle - \psi_P(\theta) : \theta \in \mathbb{R}^d\}.$$

252 Our central assumption (which is not too restrictive in view of our discussion above) on the
 253 prior P and its exponential family \mathcal{F}_P is provided below. The additional condition $0 \in \text{int } \Theta_P$
 254 insures the existence of \mathbb{E}_P .

255 **Assumption 3.1.** *The reference distribution $P \in \mathcal{P}(\Omega)$ generates a minimal and steep ex-*
 256 *ponential family \mathcal{F}_P such that $0 \in \text{int } \Theta_P$.*

257 The equivalence between the two seemingly different functions³ ψ_P^* and κ_P was previously
 258 established under various assumptions: the authors of [27, Theorem 5.2] (see also [28]) impose
 259 the (restrictive) assumption that ψ_P is finite. On the other hand, the results in [18, Theorem
 260 6.17] and [39, Proposition 1] (see also [13] and a closely related result in [54, Theorem 3.4]) do
 261 not address the challenging case when y resides on the boundary of the domain. This scenario
 262 turns out to be important if (and only if) the reference distribution is defined over a countable
 263 set. Here, we provide a complete proof that overcomes these assumptions previously imposed.
 264 Our approach emphasizes the role played by the convex support of the reference distribution
 265 and leads to natural and easy to verify conditions. To this end, we will first need to examine
 266 the domains $\text{dom } \kappa_P$ and $\text{dom } \psi_P^*$. For Cramér's rate function ψ_P^* , a characterization of the
 267 domain is summarized in the following proposition.

268 **Proposition 3.2 (Domain of the Cramér rate function ψ_P^* [4, Theorems 9.1, 9.4 and 9.5]).** *Let*
 269 *$P \in \mathcal{P}(\Omega)$ be a reference distribution satisfying [Assumption 3.1](#). Then, $\text{int } \Omega_P^{\text{cc}} \subseteq \text{dom } \psi_P^* \subseteq$*
 270 *Ω_P^{cc} . Moreover, the following hold:*

- 271 (a) *If Ω_P is finite, then $\text{dom } \psi_P^* = \Omega_P^{\text{cc}}$.*
- 272 (b) *If Ω_P is countable, then $\text{dom } \psi_P^* \supseteq \text{conv } \Omega_P$.*
- 273 (c) *If Ω_P is uncountable, then $\text{dom } \psi_P^* = \text{int } \Omega_P^{\text{cc}}$.*

³ ψ_P^* appears in *Cramér's Theorem* central in large deviations theory [28]. A more general form of κ_P appears in *Sanov's Theorem*.

274 In order to establish a similar characterization for the domain of the MEM function, we
 275 will need to make precise the relation between Ω_P and the expected value \mathbb{E}_P for a given
 276 probability measure $P \in \mathcal{P}(\Omega)$. To this end, we first recall some additional definitions and
 277 results (see, for example, [48, Section 6]). Consider two subsets $S, \hat{S} \subseteq \mathbb{R}^d$ and assume further
 278 that $S \subseteq \hat{S}$. Then $\text{cl } S \subseteq \text{cl } \hat{S}$, $\text{int } S \subseteq \text{int } \hat{S}$ and $\text{conv } S \subseteq \text{conv } \hat{S}$.

279 Denote the closed Euclidean unit ball in \mathbb{R}^d by \mathcal{B}_d . The *relative interior* [48, Section 6] of
 280 a convex set $S \subseteq \mathbb{R}^d$ is defined as

$$281 \quad \text{ri } S := \left\{ x \in \mathbb{R}^d : \exists \tau > 0 \text{ such that } (x + \tau \mathcal{B}_d) \cap \text{aff } S \subseteq S \right\}.$$

283 E.g., for the *unit simplex* $\Delta_d := \{y \in \mathbb{R}_+^d : \langle e, y \rangle = 1\}$ we have $\text{ri } \Delta_d := \{y \in \mathbb{R}_{++}^d : \langle e, y \rangle = 1\}$.
 284 Some facts which will be used in the sequel are summarized in the following lemma. Further
 285 details and proofs can be found in [48, Section 6, Theorem 13.1].

286 **Lemma 3.3 (On the relative interior).** *Let $S \subseteq \mathbb{R}^d$ be nonempty and convex. Then:*

- 287 (a) *It holds that $\text{ri}(\text{cl } S) = \text{ri } S$ and $\text{ri } S \subseteq S \subseteq \text{cl } S$.*
 288 (b) *If $\dim S = d$ then $\text{ri } S = \text{int } S$ and, in particular, $\text{int } S \neq \emptyset$.*
 289 (c) *It holds that $x \in \text{ri } S$ if and only if $\sigma_{S-x}(v) \geq 0$ where the last inequality is strict for*
 290 *every $v \in \mathbb{R}^d$ such that $-\sigma_S(-v) \neq \sigma_S(v)$.*

291 **Lemma 3.4 (Domain of expected value).** *Let $P \in \mathcal{P}(\Omega)$ and assume that \mathbb{E}_P exists. Then*
 292 *$\mathbb{E}_P \in \text{ri } \Omega_P^{\text{cc}} = \text{ri}(\text{conv } \Omega_P)$.*

293 *Proof.* By definition of σ_{Ω_P} , for any $v \in \mathbb{R}^d$, it holds that $-\sigma_{\Omega_P}(-v) \leq \langle v, y \rangle \leq \sigma_{\Omega_P}(v)$.
 294 As $P \in \mathcal{P}(\Omega)$, this implies, for all $v \in \mathbb{R}^d$, that

$$295 \quad (3.3) \quad \langle v, \mathbb{E}_P \rangle = \int_{\Omega_P} \langle v, y \rangle dP(y) \leq \sigma_{\Omega_P}(v) \int_{\Omega_P} dP(y) = \sigma_{\Omega_P}(v).$$

297 If there exists some subset $A \subseteq \Omega_P$ such that $P(\{y \in A : \langle v, y \rangle < \sigma_{\Omega_P}(v)\}) > 0$, then the
 298 inequality in (3.3) is strict. We will show that, for any $v \in \mathbb{R}^d$ such that $-\sigma_{\Omega_P}(-v) \neq \sigma_{\Omega_P}(v)$,
 299 such a subset exists; the desired result then follows from Lemma 3.3 (c) and the equivalence
 300 $\sigma_{\Omega_P^{\text{cc}}}(v) = \sigma_{\Omega_P}(v)$ [49, Theorem 8.24]. Indeed, let $v \in \mathbb{R}^d$ such that $-\sigma_{\Omega_P}(-v) \neq \sigma_{\Omega_P}(v)$, i.e.
 301 $-\sigma_{\Omega_P}(-v) < \sigma_{\Omega_P}(v)$. Pick $\tau \in (-\sigma_{\Omega_P}(-v), \sigma_{\Omega_P}(v))$ and consider $A = \{y \in \Omega_P : \langle v, y \rangle \leq \tau\}$.
 302 As $\tau < \sigma_{\Omega_P}(v)$, we have $A \subset \{y \in \Omega_P : \langle v, y \rangle < \sigma_{\Omega_P}(v)\}$, and

$$303 \quad P(A) = P(\{y \in \Omega_P : \langle -v, y \rangle \geq -\tau\}) = P(\{y \in \Omega_P : \sigma_{\Omega_P}(-v) \geq \langle -v, y \rangle \geq -\tau\}) > 0,$$

305 where the strict inequality follows from the definition of $\sigma_{\Omega_P}(-v)$ and $\sigma_{\Omega_P}(-v) > -\tau$. Hence,
 306 A satisfies the desired conditions, which establishes the result. ■

307 We are now in a position to present and prove a characterization for the domain of the MEM
 308 function, analogous to Proposition 3.2. We will use the following notation

$$309 \quad \mathcal{Q}_P(y) := \{Q \in \mathcal{P}(\Omega) : \mathbb{E}_Q = y, Q \ll P\}.$$

311 Observe that $y \in \text{dom } \kappa_P$ if and only if $\mathcal{Q}_P(y) \neq \emptyset$.

312 **Lemma 3.5 (Domain of the MEM function κ_P).** *Let $P \in \mathcal{P}(\Omega)$ be a reference distribution*
 313 *satisfying Assumption 3.1. Then:*

- 314 (a) If Ω_P is countable, then $\text{dom } \kappa_P = \text{conv } \Omega_P$. Hence, if Ω_P is finite, then $\text{dom } \kappa_P = \Omega_P^{\text{cc}}$.
 315 (b) If Ω_P is uncountable, then $\text{dom } \kappa_P = \text{int } \Omega_P^{\text{cc}}$.

316 *Proof.* (a) Let $y \in \text{dom } \kappa_P$, hence there exists $Q \in \mathcal{Q}_P(y)$. As $Q \ll P$, we obtain
 317 $\Omega_Q \subseteq \Omega_P$, thus $\text{conv } \Omega_Q \subseteq \text{conv } \Omega_P$. Hence, by [Lemma 3.3](#) (a) and [Lemma 3.4](#), we
 318 know that $y = \mathbb{E}_Q \in \text{ri } \Omega_Q^{\text{cc}} \subseteq \text{conv } \Omega_Q \subseteq \text{conv } \Omega_P$. Thus, $\text{dom } \kappa_P \subseteq \text{conv } \Omega_P$. For
 319 the converse inclusion, let $y \in \text{conv } \Omega_P$. By Carathéodory's theorem [20], there exist
 320 $n \leq d + 1$ points p_1, \dots, p_n in Ω_P such that $y = \sum_{i=1}^n \lambda_i p_i$ for some $\lambda \in \Delta_n$. Consider
 321 a distribution $Q \in \mathcal{P}(\Omega)$ satisfying $Q(\{p_i\}) = \lambda_i$ for all $i = 1, \dots, n$. Then, $Q \in \mathcal{Q}_P(y)$
 322 by construction. Thus, $y \in \text{dom } \kappa_P$, and we can conclude that $\text{conv } \Omega_P \subseteq \text{dom } \kappa_P$.
 323 (b) First, let $y \in \text{dom } \kappa_P$, then there exists $Q \in \mathcal{Q}_P(y)$. Since $Q \ll P$ which satisfies
 324 [Assumption 3.1](#), it holds that $\dim \Omega_Q^{\text{cc}} = \Omega_P^{\text{cc}} = d$. Otherwise, the probability measure
 325 Q ($Q(\Omega_Q) = 1$) is concentrated on a lower dimensional affine subspace in contradiction
 326 to the absolute continuity of Q with respect to P . Hence, using [Lemma 3.4](#) and
 327 [Lemma 3.3](#) (b), we obtain that $y = \mathbb{E}_Q \in \text{ri } \Omega_Q^{\text{cc}} = \text{int } \Omega_Q^{\text{cc}} \subseteq \text{int } \Omega_P^{\text{cc}}$. For the converse
 328 inclusion, by [Proposition 3.2](#), $y \in \text{int } \Omega_P^{\text{cc}} = \text{dom } \psi_P^* = \text{int } (\text{dom } \psi_P^*) = \text{dom } \nabla \psi_P^*$, and
 329 we conclude that $y = \mathbb{E}_{P_\theta}$ for $\theta = \nabla \psi_P^*(y)$. Since $P_\theta \ll P$ for P_θ from the exponential
 330 family generated by P , we find that $P_\theta \in \mathcal{Q}_P(y)$ and therefore $y \in \text{dom } \kappa_P$. ■

331 Combining [Lemma 3.5](#) with [Proposition 3.2](#) yields the following corollary.

- 332 [Corollary 3.6.](#) Let $P \in \mathcal{P}(\Omega)$ be a reference distribution satisfying [Assumption 3.1](#). Then,
 333 (a) If Ω_P is countable and $\text{conv } \Omega_P$ is closed (i.e., $\text{conv } \Omega_P = \Omega_P^{\text{cc}}$), then $\text{dom } \kappa_P =$
 334 $\text{dom } \psi_P^* = \Omega_P^{\text{cc}}$. In particular, $\text{dom } \kappa_P = \text{dom } \psi_P^* = \Omega_P^{\text{cc}}$ if Ω_P is finite.
 335 (b) If Ω_P is uncountable, then $\text{dom } \kappa_P = \text{dom } \psi_P^* = \text{int } \Omega_P^{\text{cc}}$.

336 The following lemma will be crucial for proving the equivalence between the MEM function
 337 κ_P and Cramér's rate function ψ_P^* . The proof of the lower bound follows similar arguments
 338 as in [18, Theorem 6.17] and [39, Proposition 1] and we include it here for completeness.

339 [Lemma 3.7.](#) Let $P \in \mathcal{P}(\Omega)$ be a reference distribution satisfying [Assumption 3.1](#). Then:

$$340 \quad \psi_P^*(y) \leq \kappa_P(y) \leq \psi_P^*(y) + \text{KL}(Q|P_\theta) - D_{\psi_P^*}(y, \nabla \psi_P(\theta)),$$

342 for any $y \in \text{dom } \kappa_P$, $Q \in \mathcal{Q}_P(y)$ and $\theta \in \text{int } \Theta_P$.

343 *Proof.* For any $\theta \in \text{int } \Theta_P$ and $Q \in \mathcal{Q}_P(y)$ we obtain that $Q \ll P_\theta$ due to the mutual
 344 absolute continuity between P_θ and P . Hence,
 345

$$346 \quad (3.4) \quad \text{KL}(Q|P) = \int_{\Omega} \log \left(\frac{dQ}{dP} \right) dQ = \int_{\Omega} \log \left(\frac{dQ}{dP_\theta} \right) dQ + \int_{\Omega} \log \left(\frac{dP_\theta}{dP} \right) dQ$$

$$347 \quad = \text{KL}(Q|P_\theta) + \int_{\Omega} [\langle z, \theta \rangle - \psi_P(\theta)] dQ(z) = \text{KL}(Q|P_\theta) + \langle y, \theta \rangle - \psi_P(\theta),$$

349 where the last identity uses $y = \mathbb{E}_Q$. Since (3.4) holds for all $\theta \in \text{int } \Theta_P$ and $\text{KL}(Q|P_\theta) \geq 0$,

$$350 \quad (3.5) \quad \text{KL}(Q|P) \geq \sup \{ \langle y, \theta \rangle - \psi_P(\theta) : \theta \in \text{int } \Theta_P \} = \psi_P^*(y),$$

352 due to the closedness of ψ_P , see [Proposition 2.2](#). The lower bound for κ_P follows immediately
 353 from its definition and the above inequality.

354 As for the upper bound: by (3.4) and (2.2), for any $Q \in \mathcal{Q}_P(y)$ and $\theta \in \text{int } \Theta_P$, we have

$$\begin{aligned}
\text{KL}(Q|P) &= \text{KL}(Q|P_\theta) + \langle y, \theta \rangle - \psi_P(\theta) \\
&= \text{KL}(Q|P_\theta) + \langle y - \nabla\psi_P(\theta), \theta \rangle + \langle \nabla\psi_P(\theta), \theta \rangle - \psi_P(\theta) \\
355 &= \text{KL}(Q|P_\theta) - [\psi_P^*(y) - \psi_P^*(\nabla\psi_P(\theta)) - \langle y - \nabla\psi_P(\theta), \theta \rangle] + \psi_P^*(y) \\
356 &= \text{KL}(Q|P_\theta) - D_{\psi_P^*}(y, \nabla\psi_P(\theta)) + \psi_P^*(y).
\end{aligned}$$

357 Then the result follows due to the fact that $\kappa_P(y) \leq \text{KL}(Q|P)$ for all $Q \in \mathcal{Q}_P(y)$. ■

358 **Theorem 3.8 (Equivalence between Cramér's rate function and the MEM function).** *Let*
359 *$P \in \mathcal{P}(\Omega)$ satisfy Assumption 3.1, and assume that one of the following two conditions holds:*

360 (i) Ω_P is uncountable.

361 (ii) Ω_P is countable and $\text{conv } \Omega_P$ is closed (as is the case when Ω_P is finite).

362 Then, $\kappa_P = \psi_P^*$. In particular, κ_P is closed, proper and convex.

363 *Proof.* First, let $y \in \text{int } \Omega_P^{\text{cc}}$. By Assumption 3.1, $\nabla\psi_P$ is a bijection between $\text{int}(\text{dom } \psi_P)$
364 $= \text{int } \Theta_P$ and $\text{int}(\text{dom } \psi_P^*) = \text{int } \Omega_P^{\text{cc}}$, where the latter uses Proposition 3.2. Thus, there exists
365 $\theta \in \text{int } \Theta_P$ such that $y = \nabla\psi_P(\theta) = \mathbb{E}_{P_\theta}$. Applying Lemma 3.7 with $Q = P_\theta$ yields

$$366 \quad (3.6) \quad \kappa_P(y) = \psi_P^*(y) \quad (y \in \text{int } \Omega_P^{\text{cc}}).$$

368 Due to Corollary 3.6, this establishes the result when Ω_P is uncountable. To complete the
369 proof, we only need to address the case when $y \in \text{bd } \Omega_P^{\text{cc}}$ under assumption (ii). By Corol-
370 lary 3.6, in this case $\text{dom } \kappa_P = \text{dom } \psi_P^* = \Omega_P^{\text{cc}}$ and $\mathcal{Q}_P(y) \neq \emptyset$ for $y \in \text{bd } \Omega_P^{\text{cc}}$. Consider any
371 $Q \in \mathcal{Q}_P(y)$, then, by definition of κ_P , we have that

$$372 \quad (3.7) \quad \kappa_P(y) \leq \text{KL}(Q|P) < +\infty.$$

374 Choose any $\hat{y} \in \text{int } \Omega_P^{\text{cc}}$ and set $\hat{\theta} = \nabla\psi_P^*(\hat{y})$ (i.e., $\hat{y} = \nabla\psi(\hat{\theta})$). For any $\lambda \in [0, 1)$ consider
375 $Q_\lambda = \lambda Q + (1 - \lambda)P_{\hat{\theta}}$. Then, by linearity of $Q \mapsto \mathbb{E}_Q$ [46, Lemma 2], we obtain

$$376 \quad y_\lambda := \mathbb{E}_{Q_\lambda} = \lambda\mathbb{E}_Q + (1 - \lambda)\mathbb{E}_{P_{\hat{\theta}}} = \lambda y + (1 - \lambda)\hat{y}.$$

378 By convexity of Ω_P^{cc} and the line segment principle [10, Lemma 6.28] we conclude that
379 $y_\lambda \in \text{int } \Omega_P^{\text{cc}}$. Set $\theta_\lambda := \nabla\psi_P^*(y_\lambda)$ and observe that, by Lemma 3.7 and the nonnegativity
380 of the Bregman distance, it holds that

$$381 \quad (3.8) \quad \psi_P^*(y) \leq \kappa_P(y) \leq \psi_P^*(y) + \text{KL}(Q|Q_\lambda).$$

383 In addition, due to (3.7) and the fact that $Q \ll P \ll P_{\hat{\theta}}$, we conclude that $\text{KL}(Q|P_{\hat{\theta}}) < \infty$.
384 Thus, by (3.8) and convexity of $\text{KL}(Q|\cdot)$, we obtain

$$385 \quad \text{KL}(Q|Q_\lambda) \leq \lambda\text{KL}(Q|Q) + (1 - \lambda)\text{KL}(Q|P_{\hat{\theta}}) \rightarrow 0 \quad \text{as } \lambda \rightarrow 1. \quad \blacksquare$$

387 We refer to a solution of the optimization problem (3.2) as the *MEM distribution* and denote
388 it as Q_{MEM} . By similar arguments to the ones used in order to establish the lower bound in

389 **Lemma 3.7**, one can show that, when $y \in \text{int}(\text{dom } \kappa_P) = \text{int}(\text{conv } \Omega_P)$, the MEM distribution
 390 is a particular member of the exponential family generated by the reference distribution P .
 391 More precisely, it holds that $Q_{MEM} = P_\theta$ where $\theta = \nabla \psi_P^*(y)$ and consequently

$$392 \quad f_{Q_{MEM}}(x) = \frac{dP_\theta}{dP}(x) = \exp\left(\langle x, \theta \rangle - \log \int_{\Omega} \exp(\langle \cdot, \theta \rangle) dP\right) = \frac{\exp(\langle x, \theta \rangle)}{\int_{\Omega} \exp(\langle \cdot, \theta \rangle) dP}.$$

394 This, again, highlights the intimate connection between the MEM function and exponential
 395 families. The case $y \in \text{bd}(\text{dom } \kappa_P)$ is more subtle and will be the topic of future research.

396 In what follows, we assume that the reference distribution of the MEM function satisfies
 397 the conditions stated in **Theorem 3.8**, that is:

398 **Assumption 3.9.** *The distribution $P \in \mathcal{P}(\Omega)$ satisfies one of the following conditions:*

- 399 (i) Ω_P is uncountable.
- 400 (ii) Ω_P is countable and $\text{conv } \Omega_P$ is closed (as is the case when Ω_P is finite).

401 Under **Assumptions 3.1** and **3.9**, the MEM function and the Cramér rate function coincide.
 402 As an immediate consequence, we obtain that the MEM function κ_P is of Legendre type.
 403 More importantly, we will see that the alternative representation by means of Cramér's rate
 404 function is more tractable compared to the original definition given in (3.2).

405 **Theorem 3.10 (Properties of the MEM function).** *Let $P \in \mathcal{P}(\Omega)$ satisfy **Assumptions 3.1**
 406 and **3.9**. Then the following hold:*

- 407 (a) $\kappa_P(y) \geq 0$ and equality holds if and only if $y = \mathbb{E}_P$.
- 408 (b) κ_P is of Legendre type.
- 409 (c) κ_P is coercive in the sense that $\lim_{\|y\| \rightarrow \infty} \kappa_P(y) = +\infty$ [**9**, Definition 11.10]. In
 410 particular, $\kappa_P(y)$ is level bounded.
- 411 (d) If M_P is finite (which holds, in particular, when Ω_P is bounded), then κ_P is superco-
 412 coercive in the sense that $\lim_{\|y\| \rightarrow \infty} \kappa_P(y)/\|y\| = +\infty$ [**9**, Definition 11.10].

413 *Proof.* Part (a) is evident from the definition of κ_P as given in (3.2) and [**18**, Proposition
 414 6.2]. Part (b) follows directly from the equivalence to the Cramér rate function ψ_P^* and
 415 **Corollary 2.5**. To see (c), observe that (a) implies that κ_P admits a unique minimizer \mathbb{E}_P
 416 which combined with the fact that κ_P is closed, proper and convex (since κ_P is of Legendre type
 417 due to (b)) establishes the result by [**2**, Proposition 3.1.3]. Lastly, if the moment generating
 418 function is finite, then so is ψ_P , and the supercoercivity of $\kappa_P = \psi_P^*$ follows from [**49**, Theorem
 419 11.8(d)].⁴ If Ω_P is bounded then $\text{dom } \kappa_P$ is bounded due to **Lemma 3.5**. In this case, $\kappa_P = \psi_P^*$
 420 is trivially supercoercive and the claim that ψ_P is finite follows from [**49**, Theorem 11.8(d)]. ■

421 The results presented in the remainder of this work are established under **Assumptions 3.1**
 422 and **3.9** which, in particular, ensure the equivalence between the MEM and Cramér rate
 423 functions. For this reason, we take this opportunity to standardize our nomenclature: between
 424 the two options (κ_P or ψ_P^*) we will opt for the one that corresponds to the Cramér rate function
 425 ψ_P^* . This choice is motivated by our intent to emphasize the more computationally appealing
 426 definition and the connection to the log-normalizer function ψ_P . Nevertheless, in the definition

⁴The definition of supercoercive convex functions we use here follows [**9**, Definition 11.10]. In [**49**] the authors refer to such functions as coercive (see [**49**, Definition 3.25]).

427 of some new concepts defined by means of Cramér's rate function, we will adopt the MEM
428 terminology in order to emphasize the motivation in the context of estimation.

429 If the reference distribution belongs to an exponential family generated by some measure
430 $P \in \mathcal{M}(\Omega)$, i.e., if for some $\hat{\theta} \in \Theta_P$ we consider a new exponential family generated by the
431 probability measure $P_{\hat{\theta}}$,⁵ then the corresponding moment generating function takes the form

$$432 \quad (3.9) \quad M_{P_{\hat{\theta}}}[\theta] = \exp\left(\psi_P(\hat{\theta} + \theta) - \psi_P(\hat{\theta})\right).$$

434 In this case, the Cramér rate functions that corresponds to $P_{\hat{\theta}}$ and P share a useful relation
435 summarized in the following lemma. We include the simple proof in [Appendix A](#).

436 **Lemma 3.11.** *Let \mathcal{F}_P be a minimal and steep exponential family generated by $P \in \mathcal{M}(\Omega)$
437 and assume further that, for any $\theta \in \text{int } \Theta_P$, [Assumption 3.9](#) holds for $P_{\theta} \in \mathcal{P}(\Omega)$. Then, for
438 any $\hat{\theta} \in \text{int } \Theta_P$ and $y \in \text{dom } \psi_P^*$, we have $\psi_{P_{\hat{\theta}}}^*(y) = D_{\psi_P^*}(y, \hat{y})$ where $\hat{y} := \nabla \psi_P(\hat{\theta}) \in \text{int } \Omega_P^c$.*

439 We list in [Table 1](#) below a number of examples of Cramér rate functions that correspond
440 to most of the popular distributions (i.e. choices of the reference distribution $P \in \mathcal{P}(\Omega)$).
441 Some of the functions admit a closed form expression while others are given implicitly.⁶ The
442 derivations and further details are included as a supplementary material. Observe that all
443 cases considered below satisfy [Assumptions 3.1](#) and [3.9](#) which guarantees the equivalence
444 established in [Theorem 3.8](#): indeed, with some exceptions, all the distributions in [Table 1](#) are
445 minimal with a natural parameter space Θ_P open which implies steepness. These exceptions
446 are: the multinomial distribution which is minimal under an appropriate reformulation, and
447 the multivariate normal-inverse Gaussian which is steep (see supplementary material). Here,
448 we provide the Cramér rate function of the multinomial distribution in minimal form. Thus,
449 [Assumption 3.1](#) holds for all the distributions given in [Table 1](#). This comprehensive list
450 complements and extends some previously established formulas [[39](#), [54](#)].

451 Many computations are facilitated in the presence of separability as described in the
452 following remark.

453 **Remark 3.12 (Separability of ψ_P^*).** In most examples, the reference distribution $P \in \mathcal{P}(\Omega)$
454 admits a separable structure of the form $P(y) = P_1(y_1)P_2(y_2) \cdots P_d(y_d)$ where $P_i \in \mathcal{P}(\Omega_i)$,
455 $\Omega_i \subset \mathbb{R}$, i.e., each component corresponds to an i.i.d. random variable. In this case, since
456 $\mathbb{M}_P[\theta] = \prod_{i=1}^d \mathbb{M}_{P_i}[\theta_i]$ [[50](#), Section 4.4], we have

$$457 \quad \psi_P^*(y) = \sup \left\{ \langle y, \theta \rangle - \log(\mathbb{M}_P[\theta]) : \theta \in \mathbb{R}^d \right\} = \sum_{i=1}^d \sup \{ y_i \theta_i - \log(\mathbb{M}_{P_i}[\theta_i]) : \theta_i \in \mathbb{R} \}.$$

458 Hence, in most of our examples below we will consider only the case $d = 1$. ◇

459 In [Table 1](#) we employ the convention that $0 \log(0) = 0$ and define

$$460 \quad \Delta_{(d)} := \left\{ y \in \mathbb{R}_+^d : \sum_{i=1}^d y_i \leq 1 \right\} \quad \text{and} \quad I(p) := \{ y \in \mathbb{R}^d : y_i = 0 \ (p_i = 0) \} \quad (p \in \mathbb{R}^d).$$

⁵Recall from the definition of \mathcal{F}_P that $P_{\hat{\theta}}$ is the probability measure with $\frac{dP_{\hat{\theta}}}{dP}(y) = \exp(\langle y, \hat{\theta} \rangle - \psi_P(\hat{\theta}))$.

⁶One can evaluate Cramér's rate function value at a point of interest by solving a nonlinear system.

| Reference Distribution (P) | Cramér Rate Function ($\psi_P^*(y)$) | dom ψ_P^* |
|--|--|----------------------------|
| Multivariate Normal ($\mu \in \mathbb{R}^d, \Sigma \in \mathbb{S}^d : \Sigma \succ 0$) | $\frac{1}{2}(y - \mu)^T \Sigma^{-1} (y - \mu)$ | \mathbb{R}^d |
| Multivar. Normal-inverse Gaussian ($\mu, \beta \in \mathbb{R}^d, \alpha, \delta \in \mathbb{R}, \Sigma \in \mathbb{R}^{d \times d}$: $\delta > 0, \Sigma \succ 0, \alpha \geq \sqrt{\beta^T \Sigma \beta}$ $\gamma := \sqrt{\alpha^2 - \beta^T \Sigma \beta}$) | $\alpha \sqrt{\delta^2 + (y - \mu)^T \Sigma^{-1} (y - \mu)} - \beta^T (y - \mu) - \delta \gamma$ | \mathbb{R}^d |
| Gamma ($\alpha, \beta \in \mathbb{R}_{++}$) | $\beta y - \alpha + \alpha \log\left(\frac{\alpha}{\beta y}\right)$ | \mathbb{R}_{++} |
| Laplace ($\mu \in \mathbb{R}, b \in \mathbb{R}_{++}$) | $\begin{cases} 0, & y = \mu, \\ \sqrt{1 + \rho(y)^2} - 1 + \log\left(\frac{\sqrt{1 + \rho(y)^2} - 1}{\rho(y)^2/2}\right), & y \neq \mu, \end{cases}$ ($\rho(y) := (y - \mu)/b$) | \mathbb{R} |
| Poisson ($\lambda \in \mathbb{R}_{++}$) | $y \log(y/\lambda) - y + \lambda$ | \mathbb{R}_+ |
| Multinomial ($n \in \mathbb{N}, p \in \Delta_{(d)}$: $\sum_{i=1}^d p_i < 1$) | $\sum_{i=1}^d y_i \log\left(\frac{y_i}{n p_i}\right) + \left(n - \sum_{i=1}^d y_i\right) \log\left(\frac{n - \sum_{i=1}^d y_i}{n(1 - \sum_{i=1}^d p_i)}\right)$ | $n \Delta_{(d)} \cap I(p)$ |
| Negative Multinomial ($p \in [0, 1]^d$, $y_0 \in \mathbb{R}_{++}, p_0 := 1 - \sum_{i=1}^d p_i > 0$) | $\sum_{i=0}^d y_i \log\left(\frac{y_i}{p_i \bar{y}}\right)$ ($\bar{y} := \sum_{i=0}^d y_i$) | $\mathbb{R}_+^d \cap I(p)$ |
| Discrete Uniform ($a, b \in \mathbb{Z} : a \leq b$, $\mu := (a + b)/2, n := b - a + 1$) | $\begin{cases} 0, & y = \mu, \\ (y - \mu)\theta - \log\left(\frac{e^{(b-\mu)\theta} - e^{(a-\mu)\theta}}{n(e^\theta - 1)}\right), & y \neq \mu, \end{cases}$ where $\theta \in \mathbb{R} : y + \frac{e^\theta}{e^\theta - 1} = \frac{(b+1)e^{(b+1)\theta} - a e^{a\theta}}{e^{(b+1)\theta} - e^{a\theta}}$ | $[a, b]$ |
| Continuous Uniform ($a, b \in \mathbb{R} : a < b, \mu := (a + b)/2$) | $\begin{cases} 0, & y = \mu, \\ (y - \mu)\theta - \log\left(\frac{e^{(b-\mu)\theta} - e^{(a-\mu)\theta}}{(b-a)\theta}\right), & y \neq \mu, \end{cases}$ where $\theta \in \mathbb{R} : y + \frac{1}{\theta} = \frac{b e^{b\theta} - a e^{a\theta}}{e^{b\theta} - e^{a\theta}}$ | (a, b) |
| Logistic ($\mu \in \mathbb{R}, s \in \mathbb{R}_{++}$) | $\begin{cases} 0, & y = \mu, \\ (y - \mu)\theta - \log(B(1 - s\theta, 1 + s\theta)), & y \neq \mu, \end{cases}$ where $\theta \in \mathbb{R}_+ : y - \mu = \frac{1}{\theta} + \frac{\pi s}{\tan(-\pi s\theta)}$ | \mathbb{R} |

Table 1: Cramér rate functions for popular distributions.

461 *Remark 3.13 (On Table 1).* We provide some additional comments on Table 1 here.

462 (a) (Special cases)

463 – As special cases of the Gamma distribution we obtain Chi-squared with pa-
 464 rameter k ($\alpha = k/2, \beta = 1/2$), Erlang (α positive integer) and exponential
 465 ($\alpha = 1$) distributions.

- 466 – As special cases of the multinomial distribution, we obtain binomial ($d = 1$,
 467 $n > 1$), Bernoulli ($d = 1, n = 1$) and categorical ($d > 1, n = 1$) distributions.
 468 – As special cases of the negative multinomial distribution we obtain the negative
 469 binomial ($d = 1$) and (shifted) geometric ($d = 1, y_0 = 1$) distributions.
 470 (b) (Statistical interpretation) For many reference distributions, ψ_P^* recovers well-known
 471 functions from information theory and related areas. Here, the MEM provides an in-
 472 formation driven, statistical interpretation for these functions. Examples include the
 473 squared Mahalanobis distance (multivariate normal), pseudo-Huber loss (multivariate
 474 normal-inverse Gaussian), Itakura-Saito distance (Gamma), Burg entropy (exponen-
 475 tial), Fermi-Dirac entropy (Bernoulli), and the generalized cross entropy (Poisson).
 476 \diamond

477 **4. The MEM Estimator and Models for Inverse Problems.** In this section we show how
 478 the MEM function can be used in various modeling paradigms. We start by presenting the
 479 MEM estimator and explore some of its properties. We then discuss its (primal and dual)
 480 analogy to the maximum likelihood (ML) estimator. Finally we will illustrate its efficacy by
 481 considering a class of linear models involving a regularization term.

482 **4.1. The Maximum Entropy on the Mean Estimator.** The maximum entropy on the
 483 mean (MEM) function gives rise to an information driven criterion for measuring the compli-
 484 ance of given data with a prior distribution. Based on this function, we can define the MEM es-
 485 timator as given in [Definition 4.1](#) below. First, we introduce some additional terminology and
 486 notation that will be used in the sequel. Let $\Omega \subseteq \mathbb{R}^d$ and let $F_\Lambda = \{P_\lambda : \lambda \in \Lambda \subseteq \mathbb{R}^d\} \subset \mathcal{P}(\Omega)$
 487 be a parameterized family of distributions indexed by $\lambda \in \Lambda$ such that $\mathbb{E}_{P_{\lambda_1}} = \mathbb{E}_{P_{\lambda_2}}$ if and
 488 only if $\lambda_1 = \lambda_2$. We call F_Λ as the *reference family* and say that it satisfies [Assumptions 3.1](#)
 489 and [3.9](#) if they hold for each $P_\lambda \in F_\Lambda$. When F_Λ is an exponential family (in this case Λ is
 490 the natural parameter space Θ_P for some $P \in \mathcal{M}(\Omega)$) the MEM estimator was studied in [[18](#),
 491 Chapter 6]. We stress that, in our presentation, F_Λ need *not* be an exponential family.

492 **Definition 4.1 (MEM estimator).** *Let $F_\Lambda \subset \mathcal{P}(\Omega)$ be a reference family satisfying [Assump-](#)
 493 [tions 3.1](#) and [3.9](#) and assume that $\mathbb{E}_{P_{\lambda_1}} = \mathbb{E}_{P_{\lambda_2}}$ if and only if $\lambda_1 = \lambda_2$. For an observation
 494 $\hat{y} \in \mathbb{R}^d$, let $P_{\hat{\lambda}} \in F_\Lambda$ be such that $\hat{y} = \mathbb{E}_{P_{\hat{\lambda}}}$, and let $S^* \subseteq \mathbb{R}^d$ be (nonempty) closed. The MEM
 495 estimator is defined as*

$$496 \quad y_{MEM}(\hat{y}, F_\Lambda, S^*) := \operatorname{argmin}\{\psi_{P_{\hat{\lambda}}}^*(y) : y \in S^*\}.$$

498 In order to simplify notation, in what follows, we will write $y_{MEM} := y_{MEM}(\hat{y}, F_\Lambda, S^*)$ when
 499 the dependence on the triple $(\hat{y}, F_\Lambda, S^*)$ is clear from the context.

500 **Remark 4.2 (The observation vector and its domain).** In [Definition 4.1](#), the condition that
 501 $P_{\hat{\lambda}} \in F_\Lambda$ is chosen such that $\hat{y} = \mathbb{E}_{P_{\hat{\lambda}}}$ implies that the reference distribution is indexed by the
 502 observation vector \hat{y} . This condition combined with [Assumption 3.1](#) entails that $\hat{y} \in \operatorname{int}\Omega_{P_{\hat{\lambda}}}^{cc}$
 503 must hold due to [Lemma 3.4](#). \diamond

504 In order to establish the well-definedness of the MEM estimator, we will use the following
 505 extension of [[18](#), Lemma 5.4]. The proof is included in [Appendix A](#).

506 **Lemma 4.3.** *Let $\phi : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ be closed and Legendre-type, let $\varphi : \mathbb{R}^d \rightarrow (-\infty, +\infty]$*
 507 *be proper, closed and convex such that $\text{int}(\text{dom } \phi) \cap \text{dom } \varphi \neq \emptyset$. Assume that one of the*
 508 *functions is coercive while the other is bounded from below. Then there exists a unique solution*
 509 *$y^* \in \mathbb{R}^d$ to $\min\{\phi(y) + \varphi(y) : y \in \mathbb{R}^d\}$, which also satisfies $y^* \in \text{int}(\text{dom } \phi) \cap \text{dom } \varphi$.*

510 **Theorem 4.4 (Well-definedness of the MEM estimator).** *Let $F_\Lambda \subset \mathcal{P}(\Omega)$ be a reference*
 511 *family satisfying [Assumptions 3.1](#) and [3.9](#). For $\hat{y} \in \mathbb{R}^d$, let $P_{\hat{\lambda}} \in F_\Lambda$ such that $\hat{y} = E_{P_{\hat{\lambda}}}$, and*
 512 *let $S^* \subseteq \mathbb{R}^d$ be closed with $S^* \cap \text{dom } \psi_{P_{\hat{\lambda}}}^* \neq \emptyset$. Then, the MEM estimator y_{MEM} exists. If, in*
 513 *addition, S^* is convex and $\text{int}(\text{dom } \psi_{P_{\hat{\lambda}}}^*) \cap S^* \neq \emptyset$, y_{MEM} is unique and in $\text{int}(\text{dom } \psi_{P_{\hat{\lambda}}}^*) \cap S^*$.*

514 *Proof.* Recall that, by [Theorem 3.10](#), $\psi_{P_{\hat{\lambda}}}^*$ is coercive and of Legendre type (proper, closed,
 515 steep and strictly convex on the interior of its domain). Observe that $S^* \subset \mathbb{R}^d$ is closed and
 516 $S^* \cap \text{dom } \psi_{P_{\hat{\lambda}}}^* \neq \emptyset$. Thus, the function $\psi_{P_{\hat{\lambda}}}^* + \delta_{S^*}$ is proper, closed and coercive. Hence, the
 517 existence of the MEM estimator follows from [[2](#), Remark 3.4.1, Theorem 3.4.1]. The case
 518 when S^* is convex and $\text{int}(\text{dom } \psi_{P_{\hat{\lambda}}}^*) \cap S^* \neq \emptyset$ follows from [Lemma 4.3](#) with $\phi = \psi_{P_{\hat{\lambda}}}^*$ and
 519 $\varphi = \delta_{S^*}$ due to the coercivity of $\psi_{P_{\hat{\lambda}}}^*$ and the fact that δ_{S^*} is bounded from below. ■

520 **4.1.1. Analogy Between MEM and ML (for Exponential Families).** *Maximum likelihood*
 521 *(ML) is arguably the most popular principle for statistical estimation. Here, the estimated*
 522 *parameters are chosen as the most likely to produce a given sample of observed data while*
 523 *satisfying model assumptions. More precisely, for some $\Omega \subseteq \mathbb{R}^d$, the model is defined by*
 524 *means of a nonempty, closed set $S \subseteq \mathbb{R}^d$ of admissible parameters and a parameterized family*
 525 *of distributions $F_\Lambda = \{P_\lambda : \lambda \in \Lambda \subseteq \mathbb{R}^m\} \subset \mathcal{P}(\Omega)$ with densities f_{P_λ} . Given a sample of*
 526 *observed data $\hat{y} \in \mathbb{R}^d$, the ML estimator $\lambda_{ML}(\hat{y}, F_\Lambda, S)$ is defined as*

$$527 \quad \lambda_{ML}(\hat{y}, F_\Lambda, S) := \operatorname{argmax}\{\log f_{P_\lambda}(\hat{y}) : \lambda \in S \cap \Lambda\}.$$

529 In order to simplify notation, we will write $\lambda_{ML} := \lambda_{ML}(\hat{y}, F_\Lambda, S)$ when the dependence on the
 530 triple (\hat{y}, F_Λ, S) is clear from the context.

531 An intriguing connection between the ML and MEM estimator comes to light when Λ is
 532 the natural parameter space Θ_P of an exponential family induced by $P \in \mathcal{M}(\Omega)$. The MEM
 533 estimator can then be retrieved by solving one of two alternative optimization problems each
 534 of which has a closely related problem that yields the ML estimator. One problem is driven
 535 by information theoretic arguments, while the other emphasizes a connection motivated by
 536 convex duality. These connections were previously observed in [[18](#), Chapter 6] (also [[14](#)]) and
 537 are summarized in the following theorem whose proof is in [Appendix A](#). For consistency, we
 538 denote the ML estimator as θ_{ML} .

539 **Theorem 4.5 (MEM and ML estimator analogy).** *Let \mathcal{F}_P be a minimal and steep exponential*
 540 *family generated by $P \in \mathcal{M}(\Omega)$ and assume that, for any $\theta \in \text{int } \Theta_P$, [Assumption 3.9](#) holds*
 541 *with respect to $P_\theta \in \mathcal{P}(\Omega)$. Let $S, S^* \subseteq \mathbb{R}^d$ such that $S \cap \text{dom } \psi_P \neq \emptyset$ and $S^* \cap \text{dom } \psi_P^* \neq \emptyset$.*
 542 *Finally, let $\hat{y} \in \text{int } \Omega_P^{cc}$ and set $\hat{\theta} := \nabla \psi_P^*(\hat{y})$. Then the following hold:*

543 (a) *(Primal analogy) If $S^* \cap \text{int}(\text{dom } \psi_P^*) \neq \emptyset$ and $\nabla \psi_P^*(S^* \cap \text{int}(\text{dom } \psi_P^*)) = S \cap \text{int}(\text{dom } \psi_P)$,*
 544 *then $y_{MEM} = \nabla \psi_P(\theta_{MEM})$ where*

$$545 \quad (4.1) \quad \theta_{MEM} \in \operatorname{argmin}\{KL(P_\theta | P_{\hat{\theta}}) : \theta \in S\} \quad \text{and} \quad \theta_{ML} \in \operatorname{argmin}\{KL(P_{\hat{\theta}} | P_\theta) : \theta \in S\}.$$

547 (b) (Dual analogy): We have

$$548 \quad (4.2) \quad y_{MEM} \in \operatorname{argmin}\{D_{\psi_P^*}(y, \hat{y}) : y \in S^*\} \quad \text{and} \quad \theta_{ML} \in \operatorname{argmin}\{D_{\psi_P}(\theta, \hat{\theta}) : \theta \in S\}.$$

550 The primal and dual analogy between the MEM and ML estimator for exponential families
551 clarifies that the two are symmetric principles.

552 **4.2. Examples - Linear Models.** To illustrate the versatility of the MEM estimation
553 framework, we will consider the broad class of linear models which are among the most
554 popular paradigms in statistical estimation with applications in numerous fields such as image
555 processing, bio-informatics, machine learning etc.

556 We assume that the set S^* of admissible mean value parameters is the image of a convex
557 set $X \subseteq \mathbb{R}^d$ under a linear mapping defined by a measurement matrix $A \in \mathbb{R}^{m \times d}$. In many
558 practical scenarios, this matrix satisfies some application-related properties, which in combi-
559 nation with the set X restricts the image space to a subset of \mathbb{R}^m . We will denote by \mathcal{C} the
560 set of all matrices that satisfy such a condition for the application in question. The second
561 component in the model is $F_\Lambda = \{P_\lambda : \lambda \in \Lambda \subseteq \mathbb{R}^m\} \subset \mathcal{P}(\Omega)$, a reference family indexed by
562 $\lambda \in \Lambda$ such that $\mathbb{E}_{P_{\lambda_1}} = \mathbb{E}_{P_{\lambda_2}}$ if and only if $\lambda_1 = \lambda_2$. The reference distribution is specified
563 from this family by means of the observation vector \hat{y} . From [Remark 4.2](#) it follows that such
564 a family of distributions must satisfy $\hat{y} \in \operatorname{int} \Omega_{P_\lambda}^{cc}$ for $\hat{\lambda}$ such that $\mathbb{E}_{P_{\hat{\lambda}}} = \hat{y}$. In some cases, this
565 condition imposes additional assumptions that must be satisfied by the measurement vector.
566 We will denote the set of measurement vectors that satisfy such an assumption with respect
567 to the family of distributions under consideration by $D := \{y \in \mathbb{R}^m : \mathbb{E}_{P_\lambda} = y (\lambda \in \Lambda)\}$. To
568 summarize, an MEM estimator of the linear model outlined above is obtained by solving

$$569 \quad (4.3) \quad \min \left\{ \psi_{P_{\hat{\lambda}}}^*(Ax) : x \in X \right\} \quad (\hat{\lambda} \in \Lambda : \mathbb{E}_{P_{\hat{\lambda}}} = \hat{y}),$$

571 under the following set of assumptions:

572 **Assumption 4.6 (MEM estimation for linear models).**

- 573 1. The reference family F_Λ satisfies [Assumptions 3.1 and 3.9](#).
- 574 2. The set $X \subseteq \mathbb{R}^d$ is nonempty and convex.
- 575 3. $A \in \mathcal{C}$ and for any $x \in X$ it holds that $Ax \in \operatorname{dom} \psi_P^*$.
- 576 4. The observation vector satisfies $\hat{y} \in D$.

577 In the following table, we present some examples of MEM linear models that correspond to
578 particular choices of a reference family. In all cases, we assume that the reference family
579 admits a separable structure as outlined in [Remark 3.12](#). The vectors a_i ($i = 1, \dots, m$) stand
580 for the i th row of the matrix A . We set

$$581 \quad \mathcal{C}_0 := \{A \in \mathbb{R}_+^{m \times d} : A \text{ has no zero rows or columns}\}.$$

| Reference family | Objective function ($\psi_{P_\lambda}^* \circ A$) | \mathcal{C} | X | D |
|-----------------------|--|---------------------------|---------------------|---------------------|
| Normal | $\frac{1}{2} \ Ax - \hat{y}\ _2^2$ | $\mathbb{R}^{m \times d}$ | \mathbb{R}^d | \mathbb{R}^m |
| Poisson | $\sum_{i=1}^m [\langle a_i, x \rangle \log(\langle a_i, x \rangle / \hat{y}_i) - \langle a_i, x \rangle + \hat{y}_i]$ | \mathcal{C}_0 | \mathbb{R}_+^d | \mathbb{R}_{++}^m |
| Gamma ($\beta = 1$) | $\sum_{i=1}^m [\langle a_i, x \rangle - \hat{y}_i \log(\langle a_i, x \rangle) - (\hat{y}_i - \hat{y}_i \log(\hat{y}_i))]$ | \mathcal{C}_0 | \mathbb{R}_{++}^d | \mathbb{R}_+^m |

Table 2: Linear models under the MEM estimation framework for various reference families.

582 *Remark 4.7.* Additional models are readily available by choosing any of the reference
583 distributions presented in Table 1. Alternatively, one may consider a family of linear models
584 where the natural parameters are the ones restricted to the image of a convex set under a
585 linear mapping. This class of models is commonly referred to as *generalized linear models*
586 with a *canonical link function* [44]. \diamond

587 The MEM linear model with reference family that corresponds to the normal distribution
588 coincides with its ML counterpart, resulting in the celebrated least-squares model [15]. This
589 phenomenon is unique for the normal distribution and is a direct consequence of the fact that
590 the squared Euclidean norm is the only self-conjugate function [48, Section 12].

591 Linear inverse models under the Poisson noise assumption have been successfully applied in
592 various disciplines including fluorescence microscopy, optical/infrared astronomy and medical
593 applications such as positron emission tomography (PET) (see, for example, [14, 53]). The
594 MEM linear model with Poisson reference distribution outlined in Table 2 was previously
595 suggested in [6, Subsection 5.3] as an example for the algorithmic setting considered in that
596 work (see further details in Section 5 where we expand on the framework considered in [6]).

597 If, for example, $X = \mathbb{R}^d$ and $\text{rge}A = \mathbb{R}^m$ with $m < d$, then $x \in \mathbb{R}^d$ such that $y_{ML} =$
598 $y_{MEM} = Ax = \hat{y}$. This outcome is not a result of a deep statistical characteristic but a simple
599 consequence of the model's ill-posedness, a situation when the desired solution is not uniquely
600 characterized by the model. Situations like this are among the reasons which motivate the use
601 of *regularizers* which allow to incorporate some additional (prior) knowledge of the solution.
602 This approach give rise to the following extended version of model (4.3)

$$603 \quad (4.4) \quad \min \left\{ \psi_{P_\lambda}^*(Ax) + \varphi(x) : x \in X \right\} \quad (\hat{\lambda} \in \Lambda : \mathbb{E}_{P_\lambda} = \hat{y}),$$

605 where, in our setting, $\varphi : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ stands for a proper, closed and convex function.
606 In (4.4), the optimization formulation is designed to find a solution (model estimator) that
607 balances between two criteria represented by the *fidelity* term $\psi_{P_\lambda}^* \circ A$ and the *regularization*
608 term φ . While the fidelity term penalizes the violation between the model and observations,
609 the regularization term incorporates prior information (belief) on the solution, and in many
610 cases, when the problem with the fidelity term alone is ill-posed, it also serves as a regularizer.

611 In the context of MEM, the Cramér rate function can be used to penalize violations of the
 612 solution vector $x \in \mathbb{R}^d$ with respect to some prior reference measure $R \in \mathcal{P}(\Omega)$ that satisfies
 613 [Assumptions 3.1](#) and [3.9](#). In other words, we can set $\varphi(x) = \psi_R^*(x)$.

614 In many applications, the desired reference distribution of the regularizer will admit a
 615 separable structure (à la [Remark 3.12](#)). While this is advantageous from an algorithmic per-
 616 spective (cf. [Remark 5.3](#)), other alternatives are viable. Non-separable priors can be consid-
 617 ered in order to promote desirable correlations between the entries of the solution to problem
 618 (4.4). E.g., by considering the multinomial, negative multinomial, multivariate normal in-
 619 verse Gaussian or multivariate normal (with non-diagonal correlation matrix in the latter)
 620 reference distributions intrinsically give rise to non-separable modeling. But there are other
 621 options which involve separable reference distributions with a composite structure such as

$$622 \quad (4.5) \quad \varphi(x) = \psi_R^*(Lx) \quad \text{or} \quad \varphi(x) = \sum_{i=1}^d \psi_R^*(L_i x),$$

623
 624 where $L \in \mathbb{R}^{r \times d}$, $L_i \in \mathbb{R}^{r \times d}$. For example, new variants of the well-known (discrete) *total*
 625 *variation* (TV) regularizer [[51](#)] can be considered by replacing the norm appearing in the
 626 original definition by a Cramér rate function while keeping the first-order finite difference ma-
 627 trix (further details are given in the end of [Section 5](#)). Different reference distributions might
 628 be used to promote desirable, application-specific, properties of the solution. Nevertheless, for
 629 all choices of reference distribution the resulting function will admit some desirable properties,
 630 including convexity, differentiability and coerciveness as established in [Theorem 3.10](#). As we
 631 will see in the following section, these properties allows us to consider a unified algorithmic
 632 approach for tackling problem (4.4).

633 **5. Algorithms.** The optimization formulations of statistical estimation problems as pre-
 634 sented in the previous section are solved by optimization algorithms. Customized methods,
 635 such as the ones we consider here, allow to leverage the structure of a given problem, thus
 636 resulting in a significant efficiency improvement compared to general purpose solvers. The
 637 structure of problems which are of interest for us is given by the *additive composite model*

$$638 \quad (5.1) \quad \min\{f(x) + g(x) : x \in \mathbb{R}^d\},$$

640 where $f, g : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ are proper, closed and convex.

641 We will assume that both the fidelity and regularization term, represented by f and g ,
 642 respectively, are continuously differentiable on the interior of their domain. This assumption
 643 holds for all the modeling paradigms discussed in the previous section. In particular, model
 644 (4.4) is recovered with $f = \psi_P^* \circ A$ and $g = \psi_R^*$. Our focus on this type of problem is for
 645 convenience only as our goal is merely to illustrate how modern first-order methods can be used
 646 for computing MEM estimators, much like their popular ML counterparts. We point out that
 647 we are not limited to this setting. Other models can be considered as well, e.g., by blending
 648 a fidelity term originating from an MEM modeling paradigm with a traditional regularizer or
 649 vice versa. In this case, similar algorithms are applicable under suitable adjustments.

650 The method we consider is the *Bregman proximal gradient* (BPG) method. This first-
 651 order iterative algorithm admits a comparably mild per-iteration complexity and as such it is

652 particularly suitable for contemporary large-scale applications. It is important to notice that
 653 many other methods, including second-order and primal-dual decomposition methods, can be
 654 also considered in some scenarios and can benefit from the operators derived in this work.
 655 Before we present the BPG method, we need to define its fundamental components [6, 16].

656 **Smooth adaptable kernel:** Let $f : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ be proper, closed and continuously
 657 differentiable on $\text{int}(\text{dom } f)$. Then $h : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ of Legendre type is a *smooth adaptable*
 658 *kernel* with respect to f if $\text{dom } h \subseteq \text{dom } f$ and there exists $L > 0$ such that $Lh - f$ is convex.

659 **Bregman proximal operator:** Let $g : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ be closed and proper and $h : \mathbb{R}^d \rightarrow$
 660 $(-\infty, +\infty]$ of Legendre type. Then the *Bregman proximal operator* is defined as

$$661 \quad (5.2) \quad \text{prox}_g^h(\bar{x}) := \text{argmin} \{g(x) + D_h(x, \bar{x}) : x \in \mathbb{R}^n\} \quad (\bar{x} \in \text{int}(\text{dom } h)).$$

662 The BPG method is applicable under the following assumption.

664 **Assumption 5.1.** Consider problem (5.1) and assume that there exists a function of Le-
 665 gendre type $h : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ such that:

- 666 1. h is a smooth adaptable kernel with respect to f .
- 667 2. h induces a computationally efficient Bregman proximal operator with respect to g .

668 The BPG method reads:

669 **(BPG Method)** Pick $t \in (0, 1/L]$ and $x^0 \in \text{int}(\text{dom } h)$. For $k = 0, 1, 2, \dots$ compute

$$x^{k+1} = \text{prox}_{tg}^h(\nabla h^*(\nabla h(x^k) - t\nabla f(x^k))).$$

670 For $h = (1/2)\|\cdot\|_2^2$ and f convex, $Lh - f$ is convex if and only if ∇f is L -Lipschitz. In this
 671 case, the Bregman proximal operator reduces to the classical proximal operator and the BPG
 672 method is the well-known proximal gradient algorithm [11].

673 The BPG method for solving (5.1) exhibits a sublinear convergence rate [6]. Under suitable
 674 assumptions, the convergence improves to linear [5]. Accelerated variants, which improve
 675 practical performance and have superior theoretical guarantees under additional assumptions,
 676 are also available [3, 12]. For simplicity's sake, we confine ourselves with the basic BPG
 677 scheme, but the operators to be presented can be readily applied to the enhanced algorithms.

678 In order to customize the method to a particular instance of problem (5.1), a smooth
 679 adaptable kernel and corresponding Bregman proximal operator must be specified. To illus-
 680 trate this idea for MEM estimation, we focus on the linear models discussed in the previous
 681 section. In particular, we consider the model (4.4) where $\varphi = \psi_R^*$. We assume that **Assump-**
 682 **tion 4.6** holds and that the prior reference measure $R \in \mathcal{P}(\Omega)$ satisfies **Assumptions 3.1** and **3.9**.
 683 Furthermore, we assume that $\text{dom } \psi_R \subseteq X$ which allows us to disregard the constraint $x \in X$.
 684 The latter assumption holds in many practical situations and we assume it here for simplicity.
 685 Otherwise, one can simply apply the BPG method with $g = \psi_R^* + \delta_X$ (under the appropriate
 686 adjustments to the proximal operator). In **Table 3** below, we summarize the smooth adaptable
 687 kernels suitable for the models described in the previous section, see **Table 2**. In all cases,
 688 the smooth adaptable function admits a separable structure of the form $h(x) = \sum_{j=1}^d h_j(x_j)$
 689 where $h_j : \mathbb{R} \rightarrow (-\infty, +\infty]$ ($j = 1, \dots, d$) is a (univariate) function of Legendre type. As we
 690 will see in what follows, this property is very desirable as it give rise to a computationally

691 efficient implementation of the Bregman proximal operator. For completeness, we include the
 692 explicit formulas for the operators involved in the BPG method.

| Reference family | Kernel (h_j) | Constant (L) | $[\nabla h(x)]_j$ | $[\nabla h^*(z)]_j$ |
|-----------------------|------------------|--|-------------------|---------------------|
| Normal | $(1/2)x_j^2$ | $\ A\ _2 := \sqrt{\lambda_{\max}(A^T A)}$ | x_j | z_j |
| Poisson | $x_j \log(x_j)$ | $\ A\ _1 := \max_{j=1,2,\dots,d} \sum_{i=1}^m A_{i,j} $ | $\log(x_j) + 1$ | $\exp(z_j - 1)$ |
| Gamma ($\beta = 1$) | $-\log(x_j)$ | $\ \hat{y}\ _1 := \sum_{i=1}^m \hat{y}_i $ | $-1/x_j$ | $-1/z_j$ |

Table 3: Smooth adaptable kernels and related operators that correspond to the objective function ($f = \psi_{P_\theta}^* \circ A$) of the linear models listed in Table 2.

693 The kernel and related constant that correspond to the normal reference family is a well-known
 694 consequence due to the Lipschitz gradient continuity, a special case of the smooth adaptability
 695 property considered here.⁷ The kernel and related constant that correspond to the Poisson
 696 reference family is due to [6, Lemma 8]. The kernel and related constant that correspond to
 697 the Gamma distribution follows from [6, Lemma 7].

698 We now discuss the special form of the Bregman proximal operator in the setting of the
 699 linear model (4.4) with $\varphi = \psi_R^*$. According to (5.2), for any $t > 0$, the Bregman proximal
 700 operator is defined by the smooth adaptable kernel h and the regularizer $g = \psi_R^*$ as follows:

$$701 \quad (5.3) \quad \text{prox}_{t\psi_R^*}^h(\bar{x}) = \operatorname{argmin} \left\{ t\psi_R^*(u) + D_h(u, \bar{x}) : u \in \mathbb{R}^d \right\}.$$

703 The following theorem records that, in our setting, the above operator is well defined.

704 **Theorem 5.2 (Well-definedness of the Bregman proximal operator).** *Let $h : \mathbb{R}^d \rightarrow (-\infty, +\infty]$*
 705 *be of Legendre type and let $R \in \mathcal{P}(\Omega)$ be a reference distribution satisfying the conditions in*
 706 *Assumptions 3.1 and 3.9. Assume further that $\operatorname{int}(\operatorname{dom} h) \cap \operatorname{dom} \psi_R^* \neq \emptyset$. Then, for any $t > 0$*
 707 *and $\bar{x} \in \operatorname{int}(\operatorname{dom} h)$, the Bregman proximal operator defined in (5.3) produces a unique point*
 708 *in $\operatorname{int}(\operatorname{dom} h) \cap \operatorname{dom} \psi_R^*$.*

709 *Proof.* Since $\bar{x} \in \operatorname{int}(\operatorname{dom} h)$, the function $D_h(\cdot, \bar{x})$ is proper. In addition, since h is of
 710 Legendre type, so is $D_h(\cdot, \bar{x})$. Finally, $D_h(\cdot, \bar{x})$ is bounded below (by zero) by convexity of
 711 h . The result follows from Lemma 4.3 with $\phi = D_h$ and $\varphi = t\psi_R^*$ due to the aforementioned
 712 properties of D_h and the coercivity of $t\psi_R^*$ (Theorem 3.10 and $t > 0$). ■

713 We now show that this operator is also computationally tractable. For many reference distri-
 714 butions, this fact stems from the following separability property.

⁷More precisely, the equivalence holds for convex functions such as the ones considered here. For the nonconvex case see an extension of the smooth adaptability condition presented in [16].

715 *Remark 5.3 (Separability of the Bregman proximal operator).* In all cases under con-
 716 sideration, the smooth adaptable kernel $h : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ admits a separable struc-
 717 ture $h(x) = \sum_{j=1}^d h_j(x_j)$. Therefore, by (2.4), the induced Bregman distance satisfies:
 718 $D_h(x, y) = \sum_{i=1}^d D_{h_i}(x_i, y_i)$. If, in addition, the Cramér rate function admits a separable
 719 structure $\psi_R^* = \sum_{i=1}^d \psi_{R_i}^*$ (cf. Remark 3.12), then the optimization problem defining the
 720 Bregman proximal operator is separable and can be evaluated for each component of \bar{x} . \diamond

721 Given a particular instance of problem (5.1), with fidelity term $f = \psi_{P_\lambda}^* \circ A$ and regularizer
 722 $g = \psi_R^*$, one can derive a formula for the corresponding Bregman proximal operator. These
 723 formulas are summarized in Tables 4, 5, and 6 for each of the combinations of linear models
 724 (by using a compatible kernel generating distance from Table 3) and regularizers from Table 1.
 725 Some formulas are given in a closed form, others must be evaluated numerically through a
 726 solution of a nonlinear system.⁸ Due to Remark 5.3, for most of the regularizer reference
 727 distributions (excluding only the multivariate normal, multinomial and negative multinomial)
 728 the resulting subproblem is separable. Thus, for the sake of simplicity and without loss of gen-
 729 erality, we assume that $d = 1$, i.e., the resulting formulas correspond to one entry of the vector
 730 produced by the operator. The general case follows by applying the operator components-
 731 wise on all the elements of a vector $\bar{x} \in \mathbb{R}^d$. An implementation of the operators along with
 732 selected algorithms, applications, and detailed derivations of the operators can be found under:

733

<https://github.com/yakov-vaisbourd/MEMshared>.

734

735 The following table lists the formulas of Bregman proximal operators for the normal linear
 736 family. In this case, the operator reduces to the classical proximal operator [41].

| Reference Distribution (R) | Proximal Operator ($x^+ = \text{prox}_{t\psi_R^*}(\bar{x})$) |
|---|---|
| Multivariate Normal ($\mu \in \mathbb{R}^d, \Sigma \in \mathbb{S}^d : \Sigma \succ 0$) | $x^+ = (tI + \Sigma)^{-1}(\Sigma\bar{x} + t\mu)$ |
| Multivariate Normal-inverse Gaussian ($\mu, \beta \in \mathbb{R}^d, \alpha, \delta \in \mathbb{R},$ $\Sigma \in \mathbb{R}^{d \times d} : \delta > 0, \Sigma \succ 0,$ $\alpha^2 \geq \beta^T \Sigma \beta, \gamma := \sqrt{\alpha^2 - \beta^T \Sigma \beta}$) | $x^+ = (I + \rho \Sigma^{-1})^{-1} (t\beta + \bar{x} + \rho \Sigma^{-1} \mu)$, where $\rho \in \mathbb{R}_+ :$ $(\rho\delta)^2 + \ (\rho^{-1}I + \Sigma^{-1})^{-1} (t\beta + \bar{x} - \mu)\ _{\Sigma^{-1}}^2 = (\alpha t)^2$ |
| Gamma ($\alpha, \beta \in \mathbb{R}_{++}$) | $x^+ = (\bar{x} - t\beta + \sqrt{(\bar{x} - t\beta)^2 + 4t\alpha}) / 2$ |

continued ...

⁸The solution of the nonlinear system can be efficiently approximated by various methods. In our implementation, building upon the fact that the systems involve monotonic functions (since they stem from the optimality conditions of a convex problem), we used a variant of safeguarded Newton-Raphson method.

... continued

| Reference Distribution (R) | Proximal Operator ($x^+ = \text{prox}_{t\psi_R^*}(\bar{x})$) |
|--|--|
| Laplace ($\mu \in \mathbb{R}$, $b \in \mathbb{R}_{++}$) | $x^+ = \begin{cases} \mu, & \bar{x} = \mu, \\ \mu + b\rho, & \bar{x} \neq \mu, \end{cases}$ <p>where $\rho \in \mathbb{R}$: $\alpha_1\rho^3 + \alpha_2\rho^2 + \alpha_3\rho + \alpha_4 = 0$, with $\alpha_1 = (b/t)^2b^2$, $\alpha_2 = 2(b/t)^2b(\mu - \bar{x})$, $\alpha_3 = (b/t)^2(\mu - \bar{x})^2 - 2(b/t)b - 1$, $\alpha_4 = -2(b/t)(\mu - \bar{x})$</p> |
| Poisson ⁹ ($\lambda \in \mathbb{R}_{++}$) | $x^+ = tW\left(\frac{\lambda e^{\bar{x}/t}}{t}\right)$ |
| Multinomial ($n \in \mathbb{N}, p \in \Delta(d)$: $\sum_{i=1}^d p_i < 1$) | $x^+ \in \mathbb{R}_+^d \cap I(p): (x_i^+ - \bar{x}_i)/t + \log\left(\frac{x_i^+(1 - \sum_{j=1}^d p_j)}{p_i(n - \sum_{j=1}^d x_j^+)}\right) = 0$ |
| Negative Multinomial ($p \in [0, 1]^d$, $x_0 \in \mathbb{R}_{++}$, $p_0 := 1 - \sum_{i=1}^d p_i > 0$) | $x^+ \in \mathbb{R}_+^d \cap I(p): (x_i^+ - \bar{x}_i)/t + \log\left(\frac{x_i^+}{p_i(x_0 + \sum_{j=1}^d x_j^+)}\right) = 0,$ |
| Discrete Uniform ($a, b \in \mathbb{R} : a < b$) | $x^+ = \bar{x} - t\theta^+ \text{ where } \theta^+ = 0 \text{ if } \bar{x} = (a + b)/2,$ <p>otherwise: $\theta^+ \in \mathbb{R} \setminus \{0\}$:</p> $t(\theta^+ - \bar{x}/t) + \frac{(b+1)e^{(b+1)\theta^+} - ae^{a\theta^+}}{e^{(b+1)\theta^+} - e^{a\theta^+}} = \frac{e^{\theta^+}}{e^{\theta^+} - 1}$ |
| Continuous Uniform ($a, b \in \mathbb{R} : a \leq b$) | $x^+ = \bar{x} - t\theta^+ \text{ where } \theta^+ = 0 \text{ if } \bar{x} = (a + b)/2,$ <p>otherwise: $\theta^+ \in \mathbb{R} \setminus \{0\}$:</p> $t(\theta^+ - \bar{x}/t) + \frac{be^{b\theta^+} - ae^{a\theta^+}}{e^{b\theta^+} - e^{a\theta^+}} = \frac{1}{\theta^+}$ |
| Logistic ($\mu \in \mathbb{R}$, $s \in \mathbb{R}_{++}$): | $x^+ = \bar{x} - t\theta^+ \text{ where } \theta^+ = 0 \text{ if } \bar{x} = \mu,$ <p>otherwise: $\theta^+ \in \mathbb{R} \setminus \{0\}$:</p> $t\theta^+ + \frac{1}{\theta^+} + \frac{\pi s}{\tan(-\pi s\theta^+)} = \bar{x} - \mu$ |

Table 4: Bregman Proximal Operators - Normal Linear Model ($h = \frac{1}{2}\|\cdot\|^2$).

737 Recall that the Cramér rate function induced by a uniform (discrete/continuous) or logistic
738 reference distribution does not admit a closed form. To compute their proximal operator
739 we appeal to the corresponding dual of the subproblem in (5.3). This is done via Moreau
740 decomposition (see, e.g., [11, Theorem 6.45]) which applies when the Bregman proximal op-
741 erator (5.3) reduces to the classical proximal operator (i.e., when $h = (1/2)\|\cdot\|^2$). For the

⁹We denote by $W : \mathbb{R} \rightarrow \mathbb{R}$ the Lambert W function (see, for example, [23]).

742 general case, we will employ a result summarized in [Lemma 5.4](#) and [Corollary 5.5](#) below. The
 743 proofs of both results can be found in [Appendix A](#). Some notation is needed: for a function
 744 $g : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ proper, closed and convex and of $h : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ of Legendre type
 745 we set

$$746 \quad (5.4) \quad \text{iconv}_g^h(\bar{x}) := \operatorname{argmin} \left\{ g(x) + h(\bar{x} - x) : x \in \mathbb{R}^d \right\}.$$

748 This is the (possibly empty) solution of the optimization problem defining the *infimal convo-*
 749 *lution* $(g \square h)(\bar{x}) := \inf \{ g(x) + h(\bar{x} - x) : x \in \mathbb{R}^d \}$.

750 **Lemma 5.4.** *Let $g : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ be proper, closed and convex and let $h : \mathbb{R}^d \rightarrow$
 751 $(-\infty, +\infty]$ be of Legendre type. Let $\bar{x} \in \operatorname{int}(\operatorname{dom} h)$ and assume that there exists a unique
 752 point $x^+ := \operatorname{prox}_g^h(\bar{x})$ satisfying $x^+ \in \operatorname{int}(\operatorname{dom} h) \cap \operatorname{dom} g$. Then, $y^+ := \operatorname{iconv}_{g^*}^{h^*}(\nabla h(\bar{x}))$ exists
 753 and it holds that $\nabla h(x^+) + y^+ = \nabla h(\bar{x})$.*

754 The following corollary adapts the above lemma to the setting considered in our study. Fur-
 755 thermore, we complement this result with a simple observation which is particularly useful
 756 for Bregman proximal operator computations.

757 **Corollary 5.5.** *Let $h : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ be of Legendre type and let $R \in \mathcal{P}(\Omega)$ satisfy
 758 [Assumptions 3.1](#) and [3.9](#). Assume further that $\operatorname{int}(\operatorname{dom} h) \cap \operatorname{dom} \psi_R^* \neq \emptyset$. For $t > 0$ and $\bar{x} \in$
 759 $\operatorname{int}(\operatorname{dom} h)$, let $x^+ := \operatorname{prox}_{t\psi_R^*}^h(\bar{x})$ and $\theta^+ := \operatorname{iconv}_{t\psi_R^*(\cdot/t)}^{h^*}(\bar{x})$. Then, $\nabla h(x^+) + \theta^+ = \nabla h(\bar{x})$.
 760 In particular, $\theta^+ = 0$ (and $x^+ = \bar{x}$) if and only if $\bar{x} = \mathbb{E}_R$.*

761 The formulas of Bregman proximal operators for the Poisson and Gamma ($\beta = 1$) linear fam-
 762 ilies are included in [Appendix A](#). We close our study with particular models and algorithms.
 763

764 **Barcode Image Deblurring.** Restoration of a blurred and noisy image represented by a
 765 vector $\hat{y} \in \mathbb{R}^d$ can be cast as the following optimization problem:

$$766 \quad (5.5) \quad \min \left\{ \frac{1}{2} \|Ax - \hat{y}\|_2^2 + \tau \varphi_R^*(x) : x \in \mathbb{R}^d \right\}.$$

768 $A \in \mathbb{R}^{d \times d}$ is the blurring operator and $\tau > 0$ is a regularization parameter. The noise is
 769 assumed to be Gaussian which explains the least-squares fidelity term which can be justified
 770 from the viewpoint of both the ML and, as we know from our study, the MEM framework.
 771 If the original image is a 2D barcode, a natural choice for the reference measure $R \in \mathcal{P}(\Omega)$
 772 inducing φ_R^* is a separable Bernoulli distribution with $p = 1/2$ due to the binary nature of
 773 each pixel and no preference at each pixel to take either value.¹⁰ Additional information
 774 (symbology) can be easily incorporated by an appropriate adjustment of the parameter for
 775 each known pixel (see [\[47\]](#)). Using the appropriate proximal operator from [Table 4](#), the BPG
 776 method for solving the model takes the form

$$777 \quad x_i^{k+1} \in \mathbb{R} : \quad x_i^{k+1} + t\tau \log \left(\frac{x_i^{k+1}}{1 - x_i^{k+1}} \right) = x_i^k - t[A^T(Ax^k - \hat{y})]_i, \quad (i = 1, 2, \dots, d).$$

¹⁰As mentioned in [Remark 3.13](#), Bernoulli is a special case of the multinomial distribution. This, one dimensional, distribution is used to form a d -dimensional i.i.d as described in [Remark 3.12](#).

779 As mentioned above, our focus on the Bregman proximal gradient method is only for illustra-
 780 tion purposes. Favorable accelerated algorithms that employ the proximal operators derived
 781 in this work are readily available and should be used in practice. The acceleration scheme
 782 applicable here is known as the Fast Iterative Shrinkage Thresholding Algorithm (FISTA) [12].
 783

784 **Natural Image Deblurring.** For natural image deblurring there is no obvious structure such
 785 as the binary one for barcodes. However, it is customary to assume that the image is piecewise
 786 smooth. A popular model that promotes piecewise constant restoration is the Rudin, Osher
 787 and Fatemi (ROF) model [51] based on the total variation (TV) regularizer $\sum_{i=1}^d g(L_i x)$.
 788 Here, $L_i \in \mathbb{R}^{2 \times d}$ extracts the difference between the pixel i and two adjacent pixels while g
 789 stands for either the l_1 (isotropic TV) or l_2 (anisotropic TV) norm. Variants which admit the
 790 same structure with other choices of g are also considered in the literature: in [21, Subsection
 791 6.2.3], a model with the Huber norm for g was shown to promote restoration prone to artificial
 792 flat areas. Alternatively, one may consider the pseudo-Huber norm that corresponds to an
 793 MEM regularizer induced by the multivariate normal inverse-Gaussian reference distribution
 794 with parameters $\mu = \beta = 0$, $\alpha = 1$ and $\Sigma = I$. The resulting model is similar to (5.5)
 795 where the regularization term is substituted by $\sum_{i=1}^d \psi_R^*(L_i x)$. This model can be tackled by
 796 a primal-dual decomposition method that employs the appropriate proximal operator from
 797 Table 4. For example, using the separability of the proximal operator [11, Theorem 6.6] and
 798 the extended Moreau decomposition [11, Theorem 6.45], the update formula of the Chambolle-
 799 Pock algorithm [21, Algorithm 1] reads

$$y_i^{k+1} = \frac{\rho_i}{1+\rho_i} (y^k + sL_i z^k) \quad (i = 1, 2, \dots, d),$$

$$\text{with } \rho_i \in \mathbb{R}_+ : \rho_i^2 (s\delta)^2 + \left(\frac{\rho_i}{1+\rho_i} \right)^2 \|y_i^k + sL_i z^k\|_2^2 = 1,$$

$$x^{k+1} = (I + \tau A^T A)^{-1} (x^k - \tau (L^T y^{k+1} - A^T \hat{y})),$$

$$z^{k+1} = 2x^{k+1} - x^k,$$

802 where $L^T = [L_1^T, \dots, L_d^T] \in \mathbb{R}^{d \times 2d}$, $y^k \in \mathbb{R}^{2d} : (y^k)^T = [(y_1^k)^T, \dots, (y_d^k)^T]$ with $y_i^k \in \mathbb{R}^2$ for all
 803 $i = 1, 2, \dots, d$ and s, τ are some positive step-sizes satisfying $s\tau \|L\|_2^2 < 1$.

804 We point out that an efficient implementation of the above algorithm that takes into ac-
 805 count the sparse and structured nature of the matrices L and A , respectively, will result in a
 806 per-iteration complexity of the order $O(d \log d)$. The same statement is true with regard to
 807 the BPG method in the previous and following examples.

808

809 **Poisson Linear Inverse Problem.** Poisson linear inverse problems play a prominent role
 810 in various physical and medical imaging applications. The linear model proposed in [6, Sub-
 811 section 5.3] is simply the MEM linear model with Poisson reference distribution. The authors
 812 of [6] suggest l_1 -regularization to deploy their BPG method. Alternatively, one may consider
 813 the MEM function induced by the Laplace distribution with parameters $\mu = 0$ and $b = 1$.

814 This setting leads to the following update formula of the BPG method. For $i = 1, 2, \dots, d$:

$$815 \quad \bar{x}_i^{k+1} = \exp \left(\log(x_i^k) - t \sum_{j=1}^m a_{ji} \log(\langle a_j, x^k \rangle / \hat{y}_j) \right),$$

$$816 \quad x_i^{k+1} \in \mathbb{R} : t^2 x_i^{k+1} + 2t \log \left(\frac{x_i^{k+1}}{\bar{x}_i^{k+1}} \right) = x_i^{k+1} \left[\log \left(\frac{x_i^{k+1}}{\bar{x}_i^{k+1}} \right) \right]^2.$$

817

REFERENCES

- 818 [1] C. AMBLARD, E. LAPALME, AND J.-M. LINA, *Biomagnetic source detection by maximum entropy and*
819 *graphical models*, IEEE T. Biomed. Eng., 51 (2004), pp. 427–442.
- 820 [2] A. AUSLENDER AND M. TEBoulLE, *Asymptotic Cones and Functions in Optimization and Variational*
821 *Inequalities*, Springer Science & Business Media, 2006.
- 822 [3] A. AUSLENDER AND M. TEBoulLE, *Interior gradient and proximal methods for convex and conic opti-*
823 *mization*, SIAM J. Optim., 16 (2006), pp. 697–725.
- 824 [4] O. BARNDORFF-NIELSEN, *Information and Exponential Families: in Statistical Theory*, John Wiley &
825 Sons, 2014.
- 826 [5] H. H. BAUSCHKE, J. BOLTE, J. CHEN, M. TEBoulLE, AND X. WANG, *On linear convergence of non-*
827 *Euclidean gradient methods without strong convexity and Lipschitz gradient continuity*, J. Optimiz.
828 Theory App., 182 (2019), pp. 1068–1087.
- 829 [6] H. H. BAUSCHKE, J. BOLTE, AND M. TEBoulLE, *A descent lemma beyond Lipschitz gradient continuity:*
830 *first-order methods revisited and applications*, Math. Oper. Res., 42 (2017), pp. 330–348.
- 831 [7] H. H. BAUSCHKE AND J. M. BORWEIN, *Joint and separate convexity of the Bregman distance*, in *Studies*
832 *in Computational Mathematics*, vol. 8, Elsevier, 2001, pp. 23–36.
- 833 [8] H. H. BAUSCHKE, J. M. BORWEIN, ET AL., *Legendre functions and the method of random Bregman*
834 *projections*, J. Convex Anal., 4 (1997), pp. 27–67.
- 835 [9] H. H. BAUSCHKE, P. L. COMBETTES, ET AL., *Convex Analysis and Monotone Operator Theory in Hilbert*
836 *Spaces*, vol. 408, Springer, 2011.
- 837 [10] A. BECK, *Introduction to Nonlinear Optimization: Theory, Algorithms, and Applications with MATLAB*,
838 SIAM, 2014.
- 839 [11] A. BECK, *First-order Methods in Optimization*, SIAM, 2017.
- 840 [12] A. BECK AND M. TEBoulLE, *A fast iterative shrinkage-thresholding algorithm for linear inverse problems*,
841 SIAM J. Imaging Sci., 2 (2009), pp. 183–202.
- 842 [13] A. BEN-TAL AND A. CHARNES, *A dual optimization framework for some problems of information theory*
843 *and statistics.*, tech. report, Texas Univ. at Austin Center for Cybernetic Studies, 1979.
- 844 [14] A. BEN-TAL, M. TEBoulLE, AND A. CHARNES, *The role of duality in optimization problems involv-*
845 *ing entropy functions with applications to information theory*, J. Optimiz. Theory App., 58 (1988),
846 pp. 209–223.
- 847 [15] Å. BJÖRCK, *Numerical Methods for Least Squares Problems*, SIAM, 1996.
- 848 [16] J. BOLTE, S. SABACH, M. TEBoulLE, AND Y. VAISBOURD, *First order methods beyond convexity and*
849 *Lipschitz gradient continuity with applications to quadratic inverse problems*, SIAM J. Optim., 28
850 (2018), pp. 2131–2151.
- 851 [17] L. M. BREGMAN, *The relaxation method of finding the common point of convex sets and its application*
852 *to the solution of problems in convex programming*, U.S.S.R. Comp. Math. & Math. Phys., 7 (1967),
853 pp. 200–217.
- 854 [18] L. D. BROWN, *Fundamentals of statistical exponential families: with applications in statistical decision*
855 *theory*, Institute of Mathematical Statistics, 1986.
- 856 [19] Z. CAI, A. MACHADO, R. A. CHOWDHURY, A. SPILKIN, T. VINCENT, Ü. AYDIN, G. PELLEGRINO, J.-M.
857 LINA, AND C. GROVA, *Diffuse optical reconstructions of functional near infrared spectroscopy data*
858 *using maximum entropy on the mean*, Sci. Rep., 12 (2022), pp. 1–18.

- 859 [20] C. CARATHÉODORY, *Über den Variabilitätsbereich der Fourier'schen konstanten von positiven harmonischen*
860 *funktionen*, Rend. Circ. Mat. Palermo, 32 (1911), pp. 193–217.
- 861 [21] A. CHAMBOLLE AND T. POCK, *A first-order primal-dual algorithm for convex problems with applications*
862 *to imaging*, J. Math. Imaging Vis., 40 (2011), pp. 120–145.
- 863 [22] R. A. CHOWDHURY, J. M. LINA, E. KOBAYASHI, AND C. GROVA, *MEG source localization of spatially*
864 *extended generators of epileptic activity: comparing entropic and hierarchical Bayesian approaches*,
865 PLoS one, 8 (2013), p. e55969.
- 866 [23] R. M. CORLESS, G. H. GONNET, D. E. HARE, D. J. JEFFREY, AND D. E. KNUTH, *On the LambertW*
867 *function*, Adv. Comput. Math., 5 (1996), pp. 329–359.
- 868 [24] T. M. COVER, *Elements of Information Theory*, John Wiley & Sons, 1999.
- 869 [25] H. CRAMÉR, *Sur un nouveau théorème-limite de la théorie des probabilités*, Actual. Sci. Ind., 736 (1938),
870 pp. 5–23.
- 871 [26] D. DACUNHA-CASTELLE AND F. GAMBOA, *Maximum d'entropie et problème des moments*, in *Annales de*
872 *l'IHP Probabilités et Statistiques*, vol. 26, 1990, pp. 567–596.
- 873 [27] M. D. DONSKER AND S. S. VARADHAN, *Asymptotic evaluation of certain Markov process expectations*
874 *for large time—III*, Commun. Pur. Appl. Math., 29 (1976), pp. 389–461.
- 875 [28] R. S. ELLIS, *Entropy, Large Deviations, and Statistical Mechanics*, vol. 1431, Taylor & Francis, 2006.
- 876 [29] A. FERMIN, J.-M. LOUBES, AND C. LUDENA, *Bayesian methods for a particular inverse problem seismic*
877 *tomography*, Int. J. Tomogr. Stat., 4 (2006), pp. 1–19.
- 878 [30] F. GAMBOA, *Méthode du maximum d'entropie sur la moyenne et applications*, PhD thesis, Paris 11, 1989.
- 879 [31] F. GAMBOA AND E. GASSIAT, *Bayesian methods and maximum entropy for ill-posed inverse problems*, 25
880 (1997), pp. 328–350.
- 881 [32] F. GAMBOA, C. GUÉNEAU, T. KLEIN, AND E. LAWRENCE, *Maximum entropy on the mean approach*
882 *to solve generalized inverse problems with an application in computational thermodynamics*, RAIRO
883 *Oper. Res.*, 55 (2021), pp. 355–393.
- 884 [33] C. GROVA, J. DAUNIZEAU, J.-M. LINA, C. G. BÉNAR, H. BENALI, AND J. GOTMAN, *Evaluation of EEG*
885 *localization methods using realistic simulations of interictal spikes*, Neuroimage, 29 (2006), pp. 734–
886 753.
- 887 [34] H. GZYL, *Maximum entropy in the mean: A useful tool for constrained linear problems*, in *AIP Conference*
888 *Proceedings*, vol. 659, American Institute of Physics, 2003, pp. 361–385.
- 889 [35] M. HEERS, R. A. CHOWDHURY, T. HEDRICH, F. DUBEAU, J. A. HALL, J.-M. LINA, C. GROVA, AND
890 E. KOBAYASHI, *Localization accuracy of distributed inverse solutions for electric and magnetic source*
891 *imaging of interictal epileptic discharges in patients with focal epilepsy*, Brain Topogr., 29 (2016),
892 pp. 162–181.
- 893 [36] E. T. JAYNES, *Information theory and statistical mechanics*, Phys. Rev., 106 (1957), p. 620.
- 894 [37] S. KULLBACK, *Information Theory and Statistics*, Courier Corporation, 1997.
- 895 [38] S. KULLBACK AND R. A. LEIBLER, *On information and sufficiency*, Ann. Math. Stat., 22 (1951), pp. 79–
896 86.
- 897 [39] G. LE BESNERAIS, J.-F. BERCHER, AND G. DEMOMENT, *A new look at entropy for solving linear inverse*
898 *problems*, IEEE Trans. Inform. Theory, 45 (1999), pp. 1565–1578.
- 899 [40] P. MARÉCHAL AND A. LANNES, *Unification of some deterministic and probabilistic methods for the solu-*
900 *tion of linear inverse problems via the principle of maximum entropy on the mean*, Inverse Problems,
901 13 (1997), p. 135.
- 902 [41] J.-J. MOREAU, *Proximité et dualité dans un espace Hilbertien*, Bulletin de la Société mathématique de
903 France, 93 (1965), pp. 273–299.
- 904 [42] J. NAVAZA, *On the maximum-entropy estimate of the electron density function*, Acta Crystallogr. A, 41
905 (1985), pp. 232–244.
- 906 [43] J. NAVAZA, *The use of non-local constraints in maximum-entropy electron density reconstruction*, Acta
907 Crystallogr. A, 42 (1986), pp. 212–223.
- 908 [44] J. A. NELDER AND R. W. WEDDERBURN, *Generalized linear models*, J. R. Stat. Soc. Ser. A-G., 135
909 (1972), pp. 370–384.
- 910 [45] E. RIETSCH ET AL., *The maximum entropy approach to inverse problems-spectral analysis of short data*
911 *records and density structure of the Earth*, J. Geophys., 42 (1977), pp. 489–506.
- 912 [46] G. RIOUX, R. CHOKSI, T. HOHEISEL, P. MARÉCHAL, AND C. SCARVELIS, *The maximum entropy on the*

- 913 *mean method for image deblurring*, Inverse Problems, 37 (2021), p. 015011.
- 914 [47] G. RIOUX, C. SCARVELIS, R. CHOKSI, T. HOHEISEL, AND P. MARECHAL, *Blind deblurring of barcodes*
915 *via Kullback-Leibler divergence*, IEEE Trans. Pattern Anal. Mach. Intell., 43 (2019), pp. 77–88.
- 916 [48] R. T. ROCKAFELLAR, *Convex Analysis*, vol. 18, Princeton University Press, 1970.
- 917 [49] R. T. ROCKAFELLAR AND R. J.-B. WETS, *Variational Analysis*, vol. 317, Springer Science & Business
918 Media, 2009.
- 919 [50] V. K. ROHATGI AND A. M. E. SALEH, *An Introduction to Probability and Statistics*, John Wiley & Sons,
920 2015.
- 921 [51] L. I. RUDIN, S. OSHER, AND E. FATEMI, *Nonlinear total variation based noise removal algorithms*, Phys.
922 D, 60 (1992), pp. 259–268.
- 923 [52] B. URBAN, *Retrieval of atmospheric thermodynamical parameters using satellite measurements with a*
924 *maximum entropy method*, Inverse problems, 12 (1996), p. 779.
- 925 [53] Y. VARDI, L. A. SHEPP, AND L. KAUFMAN, *A statistical model for positron emission tomography*, J. Am.
926 Stat. Assoc., 80 (1985), pp. 8–20.
- 927 [54] M. J. WAINWRIGHT AND M. I. JORDAN, *Graphical Models, Exponential Families, and Variational Infer-*
928 *ence*, Now Publishers Inc, 2008.

929 Appendix A. Deferred Proofs and Tables.

930 A.1. Deferred Proofs.

931 *Proof (for Lemma 3.11).* For $y \in \text{dom } \psi_P^*$, we have

$$\begin{aligned} \psi_{P_{\hat{\theta}}}^*(y) &\stackrel{(3)}{=} \sup \left\{ \langle y, \theta \rangle - \log \left(M_{P_{\hat{\theta}}}[\theta] \right) : \theta \in \mathbb{R}^d \right\} \\ &\stackrel{(3.9)}{=} \sup \left\{ \langle y, \theta \rangle - [\psi_P(\hat{\theta} + \theta) - \psi_P(\hat{\theta})] : \theta \in \mathbb{R}^d \right\} \\ &= \psi_P^*(y) + \psi_P(\hat{\theta}) - \langle y, \hat{\theta} \rangle. \end{aligned}$$

934 The result follows from the definition of the Bregman distance, (2.2) and $\hat{\theta} \in \text{int}(\text{dom } \psi_P)$. ■

935 *Proof (for Lemma 4.3).* Existence and uniqueness of the solution follows from [9, Corol-
936 lary 11.15]. It remains to show that $y^* \in \text{int}(\text{dom } \phi) \cap \text{dom } \varphi$. Evidently, $y^* \in \text{dom } \phi \cap \text{dom } \varphi$
937 thus it is sufficient to show that $y^* \in \text{int}(\text{dom } \phi)$. Using [9, Theorem 16.2] and [9, Corollary
938 16.38] we have $0 \in \partial\phi(y^*) + \partial\varphi(y^*)$, in particular $\partial\phi(y^*) \neq \emptyset$. Since ϕ is of Legendre type we
939 conclude that $y^* \in \text{int}(\text{dom } \phi)$ [48, Theorem 26.1]. ■

940 *Proof (for Theorem 4.5).* Since \mathcal{F}_P is assumed to be minimal and steep, it is easy to
941 verify (recall (3.9)) that P_{θ} satisfies Assumption 3.1 for any $\theta \in \text{int } \Theta_P$. As we assume
942 $S \cap \text{dom } \psi_P \neq \emptyset$ and $S^* \cap \text{dom } \psi_P^* \neq \emptyset$, the MEM and ML estimator exist due to Theorem 4.4
943 and [18, Theorem 5.7], respectively. We now prove (b). Since \mathcal{F}_P is an exponential family, we
944 have $\log f_{P_{\hat{y}}}(\hat{y}) = \langle \hat{y}, \theta \rangle - \psi_P(\theta)$ and the ML estimator is a solution to

$$\begin{aligned} \max\{\log f_{P_{\hat{y}}}(\hat{y}) : \theta \in S\} &= \max\{\langle \hat{y}, \theta \rangle - \psi_P(\theta) : \theta \in S\} \\ &= -\min\{D_{\psi_P}(\theta, \nabla\psi_P^*(\hat{y})) : \theta \in S\} - \psi_P(\nabla\psi_P^*(\hat{y})) + \langle \hat{y}, \nabla\psi_P^*(\hat{y}) \rangle. \end{aligned}$$

947 Omitting terms independent of the minimization and using that $\hat{\theta} = \nabla\psi_P^*(\hat{y})$, the formulation
948 for the ML estimator follows. To obtain the formulation for the MEM estimator, observe that,

949 due to Lemma 3.11, we have

$$950 \quad \min\{\psi_{P_\theta}^*(y) : y \in S^*\} = \min\{D_{\psi_P^*}(y, \nabla\psi_P(\hat{\theta})) : y \in S^*\}.$$

952 Thus, the result follows by recalling that $\hat{y} = \nabla\psi_P(\hat{\theta})$.

953 We now turn to prove (a). Since $S^* \cap \text{int}(\text{dom}\psi_P^*) \neq \emptyset$ we obtain by Theorem 4.4 that
 954 $y_{MEM} \in S^* \cap \text{int}(\text{dom}\psi_P^*)$. This fact combined with the assumption $\nabla\psi_P^*(S^* \cap \text{int}(\text{dom}\psi_P^*)) =$
 955 $S \cap \text{int}(\text{dom}\psi_P)$ implies that $\nabla\psi_P^*(y_{MEM}) \in S \cap \text{int}(\text{dom}\psi_P)$. Thus, (a) follows from (b) due
 956 to the Bregman distance dual representation property (2.3) and Remark 2.6. ■

957 *Proof (for Lemma 5.4).* By the optimality condition of the optimization problem in the
 958 definition of the Bregman proximal operator (5.2) we obtain that

$$959 \quad \nabla h(\bar{x}) - \nabla h(x^+) \in \partial g(x^+).$$

961 Since g is assumed to be proper, closed and convex, (2.2) yields

$$962 \quad (\text{A.1}) \quad x^+ \in \partial g^*(\nabla h(\bar{x}) - \nabla h(x^+)).$$

964 Setting $\tilde{y} := \nabla h(\bar{x}) - \nabla h(x^+)$ and observing that $x^+ = \nabla h^*(\nabla h(\bar{x}) - \tilde{y})$ we can rewrite (A.1)
 965 as

$$966 \quad \nabla h^*(\nabla h(\bar{x}) - \tilde{y}) \in \partial g^*(\tilde{y}).$$

968 It is now easy to verify that the above is nothing else but the optimality condition for \bar{y} , thus,
 969 $\tilde{y} = y^+$ and we can conclude that $\nabla h(x^+) + y^+ = \nabla h(\bar{x})$, establishing the desired result. ■

970 *Proof (for Corollary 5.5).* By Theorem 3.10 we have that ψ_R^* is proper, closed and convex
 971 and thus $\psi_R^{**} = \psi_R$ due to [11, Theorem 4.8]. By Theorem 5.2 we know that x^+ is well
 972 defined. The proof of the first part then follows directly from Lemma 5.4 (with $g = t\psi_R^*$ and
 973 $y^+ = \theta^+$) and [11, Theorem 4.14(a)]. To see that $\theta^+ = 0$ if and only if $\bar{x} = \mathbb{E}_R$, observe
 974 that the objective function in the subproblem defining the Bregman proximal operator (5.3)
 975 is greater equal than zero, and equality holds if and only if $\bar{x} = \mathbb{E}_R$ with $x^+ = \bar{x}$. Thus, the
 976 statement holds true in view of the first part of the current corollary. ■

977 **A.2. Bregman Proximal Operators for Poisson and Gamma ($\beta = 1$) Linear Families.**

978 The following table lists the formulas of Bregman proximal operators for the Poisson and
 979 Gamma ($\beta = 1$) linear families, respectively. Observe that by Theorem 5.2 the Bregman
 980 proximal operator is well defined if $\text{int}(\text{dom}h) \cap \text{dom}\psi_R^* \neq \emptyset$. Since $\text{int}(\text{dom}h) = \mathbb{R}_{++}^d$ this
 981 implies that for the multinomial and negative multinomial distributions we must assume that
 982 $p_i > 0$ for all $i = 1, 2, \dots, d$. Furthermore, for the sake of simplicity we include the normal and
 983 normal inverse-Gaussian distributions. The multivariate variants can be found in the software
 984 documentation along with further explanations.

| Reference Distribution (R) | Bregman Proximal Operator ($x^+ = \text{prox}_{t\psi_R^*}^h(\bar{x})$) |
|--|--|
| Normal ($\mu, \sigma \in \mathbb{R} : \sigma > 0$) | $x^+ = \frac{\sigma}{t} W\left(\frac{t}{\sigma} \bar{x} e^{\frac{t\mu}{\sigma}}\right)$ |
| Normal-inverse Gaussian ($\mu, \alpha, \beta, \delta \in \mathbb{R} : \delta > 0,$ $\alpha \geq \beta , \gamma := \sqrt{\alpha^2 - \beta^2}$) | $x^+ \in \mathbb{R}_{++} :$ $(t\alpha/\sigma)(x^+ - \mu) = (t\beta - \log(x^+/\bar{x})) \sqrt{\delta^2 + (x^+ - \mu)^2/\sigma}$ |
| Gamma ($\alpha, \beta \in \mathbb{R}_{++}$) | $x^+ = \frac{\alpha t}{W\left(\frac{\alpha t \exp(t\beta)}{\bar{x}}\right)}$ |
| Laplace ($\mu \in \mathbb{R}, b \in \mathbb{R}_{++}$) | $x^+ = \begin{cases} \mu, & \bar{x} = \mu, \\ \mu + b\rho, & \bar{x} \neq \mu, \end{cases}$ where $\rho \in \mathbb{R} : \rho + \frac{2b}{t} \log\left(\frac{\mu + b\rho}{\bar{x}}\right) = \frac{b^2\rho}{t^2} \log^2\left(\frac{\mu + b\rho}{\bar{x}}\right)$ |
| Poisson ($\lambda \in \mathbb{R}_{++}$) | $x^+ = \bar{x}^{1-\tau} \lambda^\tau \quad (\tau := \frac{t}{t+1})$ |
| Multinomial ($n \in \mathbb{N}, p \in \text{int } \Delta_{(d)}$) | $x_i^+ = \gamma_i (n - \rho)^\tau \quad \left(\tau := \frac{t}{t+1}, \gamma_i := \left[\frac{p_i \bar{x}_i^{1/t}}{1 - \sum_{j=1}^d p_j}\right]^\tau\right)$ where $\rho \in \mathbb{R} : \rho = (n - \rho)^{\frac{t}{t+1}} \left(\sum_{i=1}^d \gamma_i\right)$ |
| Negative Multinomial ($p \in (0, 1)^d,$ $x_0 \in \mathbb{R}_{++}, p_0 := 1 - \sum_{i=1}^d p_i > 0$) | $x^+ \in \mathbb{R}_+^d \cap I(p) : \log\left(\frac{x_i^+}{\bar{x}_i}\right) + t \log\left(\frac{x_i^+}{p_i(x_0 + \sum_{j=1}^d x_j^+)}\right) = 0,$ |
| Discrete Uniform ($a, b \in \mathbb{R} : a < b$) | $x^+ = \bar{x} e^{-t\theta^+} \text{ where } \theta^+ = 0 \text{ if } \bar{x} = (a + b)/2,$ otherwise: $\theta^+ \in \mathbb{R} \setminus \{0\}$: $\frac{(b+1)\exp((b+1)\theta^+) - a\exp(a\theta^+)}{\exp((b+1)\theta^+) - \exp(a\theta^+)} = \frac{\exp(\theta^+)}{\exp(\theta^+) - 1} + \exp(\bar{x} - t\theta^+ - 1)$ |
| Continuous Uniform ($a, b \in \mathbb{R} : a \leq b$) | $x^+ = \bar{x} e^{-t\theta^+} \text{ where } \theta^+ = 0 \text{ if } \bar{x} = (a + b)/2,$ otherwise: $\theta^+ \in \mathbb{R} \setminus \{0\}$: $\frac{b\exp(b\theta^+) - a\exp(a\theta^+)}{\exp(b\theta^+) - \exp(a\theta^+)} = \frac{1}{\theta^+} + \exp(\bar{x} - t\theta^+ - 1)$ |
| Logistic ($\mu \in \mathbb{R}, s \in \mathbb{R}_{++}$): | $x^+ = \bar{x} e^{-t\theta^+} \text{ where } \theta^+ = 0 \text{ if } \bar{x} = \mu,$ otherwise: $\theta^+ \in \mathbb{R} \setminus \{0\}$: $\frac{1}{\theta^+} + \frac{\pi s}{\tan(-\pi s \theta^+)} + \mu = \exp(\bar{x} - t\theta^+ - 1)$ |

Table 5: Bregman Proximal Operators - Poisson Linear Model ($h_j(x) = x_j \log x_j$)

| Reference Distribution (R) | Bregman Proximal Operator ($x^+ = \text{prox}_{t\psi_R^*}^h(\bar{x})$) |
|---|---|
| Normal ($\mu, \sigma \in \mathbb{R} : \sigma > 0$) | $x^+ = \left((t/\sigma)\mu - 1/\bar{x} + \sqrt{((t/\sigma)\mu - 1/\bar{x})^2 + 4(t/\sigma)} \right) / (2t/\sigma)$ |
| Normal-inverse Gaussian ($\mu, \alpha, \beta, \delta \in \mathbb{R} : \delta > 0, \alpha \geq \beta , \gamma := \sqrt{\alpha^2 - \beta^2}$) | $x^+ \in \mathbb{R}_{++} :$ $t\alpha(x^+ - \mu)x^+ = ((t\beta - 1/\bar{x})x^+ + 1) \sqrt{\delta^2 + (x^+ - \mu)^2}$ |
| Multivariate Normal-inverse Gaussian ($\mu, \beta \in \mathbb{R}^d, \alpha, \delta \in \mathbb{R}, \Sigma = \sigma I, \sigma > 0 : \delta > 0, \Sigma \succ 0, \alpha^2 \geq \beta^T \Sigma \beta, \gamma := \sqrt{\alpha^2 - \beta^T \Sigma \beta}$) | $x_i^+ = (w_i + \rho\mu_i + \sqrt{(w_i + \rho\mu_i)^2 + 4\rho}) / (2\rho),$ with $w_i = t\beta_i - 1/\bar{x}_i$ and $\rho \in \mathbb{R}_+ :$ $(\rho\delta)^2 + \frac{1}{4\sigma} \sum_{i=1}^d \left(w_i + \sqrt{(w_i + \mu_i\rho)^2 + 4\rho} \right)^2 = (\alpha t/\sigma)^2$ |
| Gamma ($\alpha, \beta \in \mathbb{R}_{++}$) | $x^+ = \bar{x}(t\alpha + 1) / (\bar{x}t\beta + 1)$ |
| Laplace ($\mu \in \mathbb{R}, b \in \mathbb{R}_{++}$) | $x^+ = \begin{cases} \mu, & \bar{x} = \mu, \\ \mu + b\rho, & \bar{x} \neq \mu, \end{cases}$ where $\rho \in \mathbb{R} : \alpha_1\rho^3 + \alpha_2\rho^2 + \alpha_3\rho + \alpha_4 = 0,$ with $\alpha_1 = b^2((b/\bar{x})^2 - t^2), \alpha_2 = 2b(\mu((b/\bar{x})^2 - t^2) - b^2(t+1)/\bar{x}),$ $\alpha_3 = b^2((1 - \mu/\bar{x})^2 + 2t(1 - 2\mu/\bar{x})) - t^2\mu^2, \alpha_4 = 2tb\mu(1 - \mu/\bar{x})$ |
| Poisson ($\lambda \in \mathbb{R}_{++}$) | $x^+ \in \mathbb{R}_+ : t \log \left(\frac{x^+}{\lambda} \right) = \frac{1}{x^+} - \frac{1}{\bar{x}}$ |
| Multinomial ($n \in \mathbb{N}, p \in \text{ri } \Delta_{(d)}$) | $x^+ \in \text{ri } n\Delta_{(d)} : t \log \left(\frac{x_i^+(1 - \sum_{j=1}^d p_j)}{p_i(n - \sum_{j=1}^d x_j^+)} \right) = \frac{1}{x_i^+} - \frac{1}{\bar{x}_i}$ |
| Negative Multinomial ($p \in (0, 1)^d, x_0 \in \mathbb{R}_{++}, p_0 := 1 - \sum_{i=1}^d p_i > 0$) | $x^+ \in \mathbb{R}_{++}^d : t \log \left(\frac{x_i^+}{p_i(x_0 + \sum_{i=1}^d x_j^+)} \right) = \frac{1}{x_i^+} - \frac{1}{\bar{x}_i},$ |
| Discrete Uniform ($a, b \in \mathbb{R} : a < b$) | $x^+ = \bar{x} / (\bar{x}t\theta^+ + 1)$ where $\theta^+ = 0$ if $\bar{x} = (a+b)/2,$ otherwise: $\theta^+ \in \mathbb{R} \setminus \{0\}:$ $\frac{(b+1)\exp((b+1)\theta) - a\exp(a\theta)}{(\exp((b+1)\theta) - \exp(a\theta))} = \frac{\exp(\theta)}{\exp(\theta) - 1} + \frac{\bar{x}}{t\bar{x}\theta^+ + 1}$ |
| Continuous Uniform ($a, b \in \mathbb{R} : a \leq b$) | $x^+ = \bar{x} / (\bar{x}t\theta^+ + 1)$ where $\theta^+ = 0$ if $\bar{x} = (a+b)/2,$ otherwise: $\theta^+ \in \mathbb{R} \setminus \{0\}:$ $\frac{b\exp(b\theta^+) - a\exp(a\theta^+)}{\exp(b\theta^+) - \exp(a\theta^+)} = \frac{1}{\theta^+} + \frac{\bar{x}}{t\bar{x}\theta^+ + 1}$ |
| Logistic ($\mu \in \mathbb{R}, s \in \mathbb{R}_{++}$): | $x^+ = \bar{x} / (\bar{x}t\theta^+ + 1)$ where $\theta^+ = 0$ if $\bar{x} = \mu,$ otherwise: $\theta^+ \in \mathbb{R} \setminus \{0\}:$ $\frac{1}{\theta^+} + \frac{\pi s}{\tan(-\pi s\theta^+)} + \mu = \frac{\bar{x}}{\bar{x}t\theta^+ + 1}$ |

Table 6: Bregman Proximal Operators - Gamma ($\beta = 1$) Linear Model ($h_j(x) = -\log(x_j)$)