Non-Bayesian Parametric Missing-Mass Estimation

Shir Cohen, Tirza Routtenberg, Senior Member, IEEE, and Lang Tong, Fellow, IEEE,

Abstract—We consider the classical problem of missing-mass estimation, which deals with estimating the total probability of unseen elements in a sample. The missing-mass estimation problem has various applications in machine learning, statistics, language processing, ecology, sensor networks, and others. The naive, constrained maximum likelihood (CML) estimator is inappropriate for this problem since it tends to overestimate the probability of the observed elements. Similarly, the conventional constrained Cramér-Rao bound (CCRB), which is a lower bound on the mean-squared-error (MSE) of unbiased estimators of the entire probability mass function (pmf) vector, does not provide a relevant bound on the performance for the problem of missing-mass estimation. In this paper, we introduce a frequentist, non-Bayesian parametric model of the problem of missing-mass estimation. We introduce the concept of missing-mass unbiasedness by using the Lehmann unbiasedness definition. We derive a non-Bayesian CCRB-type lower bound on the missing-mass MSE (mmMSE), named the missing-mass CCRB (mmCCRB), based on the missing-mass unbiasedness. The proposed mmCCRB can be used for system design and for the performance evaluation of existing estimators. Moreover, based on the new mmCCRB, we propose a new method to improve existing estimators by an iterative missing-mass Fisher-scoring method. Finally, we demonstrate via numerical simulations that the biased version of the mmCCRB is a valid and informative lower bound on the mmMSE of state-of-the-art estimators for this problem: the CML, asymptotic profile maximum likelihood (aPML), Good-Turing, and Laplace estimators. We also show that the performance of the Laplace estimator is improved, in terms of mmMSE and missing-mass bias, by using the new missing-mass Fisher-scoring method.

Index Terms—Non-Bayesian estimation, Good-Turing estimator, probability of missing mass, constrained Cramér-Rao bound, Lehmann unbiasedness

I. INTRODUCTION

Given N samples from a population of elements belonging to different types with unknown proportions, how should one estimate the total probability of unseen types? This is a classical problem in statistics, commonly referred to as the missing-mass estimation problem [1, 2]. Missing-mass estimation has gained significant interest in various applications, such as ecological studies [3], sensor networks [4, 5], machine learning, and statistics. In the context of language processing, for example, estimation of new and existing words in text has applications such as language modeling, spelling correction, and word-sense disambiguation [6, 7]. Missingmass estimation is especially important for applied problems where the sampling procedure is expensive, and the need for acquiring more data is determined by the possibility of observing new unobserved elements.

It is well known that the naive, constrained maximum likelihood (CML) estimator of the probability, i.e. the empirical probability, is ineffective if there are insufficient samples [8, 9]. In particular, the CML estimator assigns a zero probability to unseen events, which does not provide relevant information on the missing mass. As a result, the constrained Cramér-Rao bound (CCRB) [10-14], which is associated with the asymptotic performance of the CML estimator, is inappropriate as a bound on the performance of missing-mass estimators outside the asymptotic region. This is because the CCRB is a lower bound on the mean-squared-error (MSE) of the entire probability mass function (pmf) and not on the functional defined by the missing mass. Various estimators of the missing mass have been suggested over the years [2, 8, 15-23]. However, the analysis of these estimators is challenging and there is no comprehensive non-Bayesian estimation theory for estimating the missing mass. In particular, there is a need for appropriate lower bounds on the MSE of the missing mass obtained by any estimator. This theory and these bounds are crucial for system design, error analysis, and performance analysis of existing estimation methods, and for the development of new estimation methods.

A. Summary of results

In this paper, we consider the problem of estimating the missing mass, where it is assumed that we observe samples that are drawn from an unknown distribution. First, we introduce a non-Bayesian parametric formulation of this estimation problem. We use the missing-mass squared-error as a cost function and derive the associated Lehmann unbiasedness. We develop a new non-Bayesian constrained Cramér-Rao bound (CCRB), the missing-mass CCRB (mmCCRB), which is a lower bound on the missing-mass MSE (mmMSE) of any estimator with a specific Lehmann bias. The new bound is obtained by using linear parametric constraints on the probability space and the Lehmann unbiasedness. We investigate the properties of the mmCCRB and some special cases of this bound. Based on the equality condition of the mmCCRB, we propose a new method to improve existing estimators by an iterative missing-mass Fisher-scoring method. The new bound is examined in simulations and compared with the performance of state-of-the-art estimators: the CML, Good-Turing, Laplace, and asymptotic profile maximum likelihood (aPML) estimators. We also show that the performance of the

^{©2022} IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works. This work is partially supported by the ISRAEL SCIENCE FOUNDATION (ISF), grant No. 1173/16. The work of Shir Cohen was supported under a grant from the Ministry of Science and Technology of Israel.

S. Cohen and T. Routtenberg are with the School of Electrical and Computer Engineering Ben-Gurion University of the Negev Beer-Sheva 84105, Israel, e-mail: shiru@post.bgu.ac.il, tirzar@bgu.ac.il. L. Tong is with the School of Electrical and Computer Engineering, Cornell University, Ithaca, NY 14853 USA (e-mail:lt35@cornell.edu).

Laplace estimator is improved by using the new missing-mass Fisher-scoring method.

B. Related works

Various estimators of the missing mass have been suggested in the literature. A fundamental example is the Good-Turing probability estimator [2], which was invented to decipher the Enigma code during World War II. The Good-Turing estimator, its extensions by smoothing techniques [16, 17], and the Laplace estimator [15, 18, 19], have been shown to be useful for the estimation of the probability of unseen elements [8, 20], and have been implemented in many practical applications. More recently, the profile maximum likelihood (PML) approach has been suggested and analyzed [24-26] as an alternative to the CML estimator. The PML estimator is near-optimal in the MSE sense for a uniform distribution and was shown to have impressive statistical properties.

On the theoretical side, various works analyzed the properties of specific estimators. For example, some interpretations of the Good-Turing estimator and its performance in terms of attenuation have been established in [8, 27]. A derivation of the Good-Turing estimator from a Bayesian point of view with a uniform prior was suggested in [2]. Works related to the Good-Turing estimator include analysis of its bias [2, 28], confidence intervals and convergence rate [29], and (un)consistency [30]. The performance of the Good-Turing estimator was analyzed using the theory of large deviations in [4] and the pmf of its worst-case MSE was discussed in [31].

Different performance bounds for the missing-mass estimation problem have been discussed in the literature. For example, lower and upper bounds on the expected missing mass and inequalities on the probability of large deviations of the missing mass are discussed in [32, 33]. Various existing bounds are associated with the performance of a specific estimator and are distribution-free. For example, upper bounds on the MSE of Robbins-type estimators have been proposed in [34]. Bounds on the worst-case MSE of the Good-Turing estimator have been developed in [31, 35]. On the other hand, other studies provide lower and upper bounds on the performance of any estimator for specific distributions. For example, in [31, 35] there are also lower and upper bounds on the minimax MSE of any estimator for specific distributions. It should be noted that the proposed approach in this paper applies to all algorithms/estimators and is based on evaluating the performance for each pmf value. However, there are no Cramér-Rao-type lower bounds on the averaged performance of any estimator of the missing mass. Our recent works on estimation after selection [36-40] suggest that conditional schemes, in which the performance criterion depends on the observed data, require different Cramér-Rao-type bounds for analysis and system design.

C. Organization and notation

The remainder of the paper is organized as follows: Section II presents the non-Bayesian parametric model of missingmass estimation under a multinomial model, including the conventional CCRB and constrained unbiasedness for this model and the appropriate cost function. In Section III, we derive a new CCRB-type lower bound on the mmMSE. In Section IV, we describe the missing-mass Fisher-scoring method. Numerical simulations are presented in Section V. Finally, our conclusions can be found in Section VI.

In the rest of this paper, vectors are denoted by boldface lowercase letters and matrices by boldface uppercase letters. The notations $\mathbb{1}_{\{A\}}$ and I denote the indicator function of an event A and the identity matrix, respectively. The vectors 1 and **0** are column vectors of ones and zeros, respectively, and \mathbf{e}_m is the *m*th column of the identity matrix, all with appropriate dimensions. The matrix $diag(\mathbf{a})$ denotes the diagonal matrix with vector **a** on the diagonal. The *m*th element of the vector **a**, the (m,q)th element of the matrix **A**, and the $(m_1 :$ $m_2 \times q_1 : q_2$) submatrix of **A** are denoted by a_m , $\mathbf{A}_{m,q}$, and $\mathbf{A}_{m_1:m_2,q_1:q_2}$, respectively. The trace of a matrix $\mathbf{A} \in \mathbb{R}^{M \times M}$ is defined as trace(\mathbf{A}) = $\sum_{m=1}^{M} \mathbf{A}_{m,m}$. The gradient of a vector function, \mathbf{c} , of $\boldsymbol{\theta}$, $\nabla_{\boldsymbol{\theta}} \mathbf{c}$, is a matrix in $\mathbb{R}^{K \times M}$, with the (k, m)th element equal to $\frac{\partial c_k}{\partial \theta_m}$, where $\mathbf{c} = [c_1, \dots, c_K]^T$ and $\boldsymbol{\theta} = [\theta_1, \dots, \theta_M]^T$. For a scalar function c, we denote $\nabla_{\theta}^{T} c \stackrel{\triangle}{=} (\nabla_{\theta} c)^{T}$, and $\nabla_{\theta}^{2} c \stackrel{\triangle}{=} \nabla_{\theta} \nabla_{\theta}^{T} c$. The notations $\mathbf{E}_{\theta}[\cdot]$ and $E_{\theta}[\cdot|A]$ represent the expectation and conditional expectation operators, parametrized by a deterministic vector, θ , and given the event A. For a set X, |X| represents its cardinality.

II. NON-BAYESIAN ESTIMATION OF THE MISSING MASS

In this section, we present the problem of estimating the missing mass as a non-Bayesain parameter estimation problem. In Subsection II-A we describe the observation model and the relevant probability functions. In Subsection II-B we develop the χ -unbiasedness and the CCRB for estimating the unknown pmf under this model. Finally, in Subsection II-C we formulate the missing-mass estimation problem, and present the missing-mass squared-error cost function, which is used in this paper.

A. Non-Bayesian model

Assume that there is a set of M symbols, $S = \{s_1, \ldots, s_M\}$, where the alphabet size, $M \ge 1$, is assumed to be finite and known. The elements in S may represent, for example, species in the jungle [8], words in a dictionary [6, 7], or operating sensors [4, 5]. The true probability of observing symbol s_m is denoted by $\theta_m, m = 1, \ldots, M$, where $\theta_m \neq 0$ for all $m = 1, \ldots, M$ and $\sum_{m=1}^M \theta_m = 1$. Thus, $\theta \stackrel{\Delta}{=} [\theta_1, \ldots, \theta_M]^T$ is a pmf vector over the discrete and finite set of symbols, S. As a result, θ is an element of the simplex Ω_{θ} , i.e. $\theta \in \Omega_{\theta}$, where

$$\Omega_{\boldsymbol{\theta}} \stackrel{\triangle}{=} \left\{ \boldsymbol{\theta} \in [0,1]^M \, | f(\boldsymbol{\theta}) = 0 \right\},\tag{1}$$

in which

$$f(\boldsymbol{\theta}) \stackrel{\triangle}{=} \sum_{m=1}^{M} \theta_m - 1.$$
 (2)

In general constrained parameter estimation, the null-space matrix that is orthogonal to the constraints plays an important role [10, 11]. For the considered setting, the gradient of $f(\theta) = m = 1, \dots, M$. By substituting (5) and (10) in (11), we obtain w.r.t. θ is

$$\mathbf{F} \stackrel{\triangle}{=} \nabla_{\boldsymbol{\theta}}^T f(\boldsymbol{\theta}) = \mathbf{1}_M^T. \tag{3}$$

In addition, there exists a null-space matrix $\mathbf{U} \in \mathbb{R}^{M \times (M-1)}$ such that

$$\mathbf{F}\mathbf{U} = \mathbf{1}_M^T \mathbf{U} = \mathbf{0}^T, \ \mathbf{U}^T \mathbf{U} = \mathbf{I}.$$
 (4)

In particular, it can be verified that $\mathbf{U}^T \mathbf{y} = \mathbf{0}_{M-1}$ iff $\mathbf{y} =$ $c\mathbf{1}_M$, where $c \in \mathbb{R}$ is an arbitrary constant.

Under the independent and identically distributed (i.i.d.) multinomial model [2], it is assumed that there are N i.i.d. samples, $\{x_n\}_{n=1}^N$, drawn according to the pmf described by the unknown vector, $\boldsymbol{\theta}$. We consider the problem of estimating the missing mass of the unobserved symbols, which is a function of θ . It can be verified that the pmf of the random observation vector, $\mathbf{x} \stackrel{ riangle}{=} [x_1,\ldots,x_N]^T \in \mathcal{S}^N$, is a binomial distribution:

$$p(\mathbf{x};\boldsymbol{\theta}) = \prod_{m=1}^{M} \theta_m^{C_{N,m}(\mathbf{x})}, \ \mathbf{x} \in \mathcal{S}^N,$$
(5)

where

$$C_{N,m}(\mathbf{x}) \stackrel{\Delta}{=} \sum_{n=1}^{N} \mathbb{1}_{\{x_n = s_m\}}, \ m = 1, \dots, M,$$
(6)

is the number of times that the mth element was observed out of the N samples. Therefore, the vector $[C_{N,0}(\mathbf{x}),\ldots,C_{N,M}(\mathbf{x})]^T$ has a multinomial distribution with parameters N and θ . The pmf in (5) can also be written as

$$p(\mathbf{x};\boldsymbol{\theta}) = \prod_{m=1}^{M} \theta_m^{\sum_{r=0}^{N} r \mathbb{1}_{\{m \in G_{N,r}(\mathbf{x})\}}},$$
(7)

where $G_{N,r}(\mathbf{x})$ is the set of elements that appear exactly r times in the N-length observation vector, x. That is, if $m \in$ $G_{N,r}(\mathbf{x})$, then $C_{N,m}(\mathbf{x}) = r$. In particular, the set

$$G_{N,0}(\mathbf{x}) \stackrel{\triangle}{=} \{m : m = 1, \dots, M, x_n \neq s_m, \ \forall n = 0, \dots, N\}$$
(8)

is the set of elements that do not appear in the observation vector, x, with $C_{N,m}(\mathbf{x}) = 0$. For example, upon observing the vector $\mathbf{x} = \begin{bmatrix} a, c, c \end{bmatrix}^T$ with N = 3 and $S = \{a, b, c\}$, the histogram values are $C_{3,1}(\mathbf{x}) = 1$, $C_{3,2}(\mathbf{x}) = 0$, and $C_{3,3}(\mathbf{x}) = 2$ according to (6), and the missing mass is the pmf of $\{b\}$, since according to (8) $G_{3,0} = \{b\}$.

Let us define the subspace of all observation vectors that do not include s_m as

$$\mathcal{A}_m \stackrel{\triangle}{=} \{ \mathbf{x} \in \mathcal{S}^N : m \in G_{N,0}(\mathbf{x}) \}, \ m = 1, \dots, M.$$
(9)

For a given number of measurements, N, the probability of the *m*th element being unobserved in these measurements is

$$\Pr(\mathbf{x} \in \mathcal{A}_m; \boldsymbol{\theta}) = \mathcal{E}_{\boldsymbol{\theta}} \left[\mathbb{1}_{\{m \in G_{N,0}(\mathbf{x})\}} \right] = (1 - \theta_m)^N, \quad (10)$$

 $\forall m = 1, \dots, M$. By using Bayes rule it can be seen that

$$p(\mathbf{x}|\mathbf{x} \in \mathcal{A}_m; \boldsymbol{\theta}) = \begin{cases} \frac{p(\mathbf{x}; \boldsymbol{\theta})}{\Pr(\mathbf{x} \in \mathcal{A}_m; \boldsymbol{\theta})} & \text{if } \mathbf{x} \in \mathcal{A}_m \\ 0 & \text{otherwise} \end{cases}, \quad (11)$$

$$p(\mathbf{x}|\mathbf{x} \in \mathcal{A}_m; \boldsymbol{\theta}) = \begin{cases} \frac{\prod_{l=1}^M \theta_l^{C_{N,l}(\mathbf{x})}}{(1-\theta_m)^N} & \text{if } \mathbf{x} \in \mathcal{A}_m \\ 0 & \text{otherwise} \end{cases}, \quad (12)$$

 $m = 1, \ldots, M.$

We denote by $\hat{\boldsymbol{\theta}} : S^N \to \Omega_{\boldsymbol{\theta}}$ an arbitrary estimator of the pmf vector, θ , based on the observation vector, x. The CML estimator of $\boldsymbol{\theta}$ under the parametric constraint $f(\boldsymbol{\theta}) = 0$ from (2) is given by

$$\hat{\theta}_m^{\text{CML}} = \frac{C_{N,m}(\mathbf{x})}{N}, \ m = 1, \dots, M.$$
(13)

In particular, the CML estimator assigns zero probability for unseen elements, i.e. for the missing mass. Some alternative estimators are presented in Subsection III-D.

B. CCRB and constrained unbiasedness

In this subsection, we develop the conventional CCRB and the unbiasedness condition for estimating θ under the considered model. The CCRB [10, 11] provides a lower bound on the MSE of any locally χ -unbiased estimator [13, 14, 41], which is a weaker requirement than ordinary mean unbiasedness, and is defined as follows.

Definition 1: An estimator $\hat{\boldsymbol{\theta}} : S^N \to \Omega_{\boldsymbol{\theta}}$ is said to be a locally χ -unbiased estimator in the neighborhood of $\tilde{\theta} \in \Omega_{\theta}$ if it satisfies

$$\mathbf{U}^T \mathbf{E}_{\tilde{\boldsymbol{\theta}}} [\hat{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}] = \mathbf{0}_{M-1} \tag{14}$$

and

$$\left\{ \nabla_{\boldsymbol{\theta}}^{T} \mathbf{E}_{\boldsymbol{\theta}} [\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}] \right\} \Big|_{\boldsymbol{\theta} = \tilde{\boldsymbol{\theta}}} \mathbf{U} = \mathbf{0}_{M \times (M-1)}, \quad (15)$$

where \mathbf{U} is defined in (4).

It should be noted that in this paper the notation $\tilde{\theta}$ represents a specific value (or "local" value) of the unknown parameter vector in the simplex Ω_{θ} , while θ is used as a general parameter in the different functions. For the CML estimator in (13) we obtain that

$$\mathbf{E}_{\boldsymbol{\theta}} \left[\hat{\theta}_m^{\text{CML}} - \theta_m \right] = \mathbf{E}_{\boldsymbol{\theta}} \left[\frac{C_{N,m}(\mathbf{x})}{N} - \theta_m \right] = 0, \qquad (16)$$

for all $m = 1, \ldots, M$ and for any $\theta \in \Omega_{\theta}$, where the last equality follows from the mean of a variable with a multinomial distribution. Thus, (16) implies that the CML estimator satisfies Definition 1 and that it is a locally χ unbiased estimator for any $\theta \in \Omega_{\theta}$; thus, it is a uniformly χ -unbiased estimator. The CML estimator is also a C-unbiased estimator in the Lehmann sense [13, 42], since the constraint $\boldsymbol{\theta} \in \Omega_{\boldsymbol{\theta}}$ is linear.

The CCRB on the MSE of any unbiased estimator in the sense of Definition 1 at $\theta \in \Omega_{\theta}$ is given by [11, 13, 14, 41]

$$\mathbf{E}_{\tilde{\boldsymbol{\theta}}}\left[(\hat{\boldsymbol{\theta}} - \tilde{\boldsymbol{\theta}})(\hat{\boldsymbol{\theta}} - \tilde{\boldsymbol{\theta}})^T\right] \succeq \mathbf{U}(\mathbf{U}^T \mathbf{J}(\tilde{\boldsymbol{\theta}})\mathbf{U})^{-1}\mathbf{U}^T, \quad (17)$$

where the conventional Fisher information matrix (FIM) is

$$\mathbf{J}(\boldsymbol{\theta}) = \mathbf{E}_{\boldsymbol{\theta}} \left[\nabla_{\boldsymbol{\theta}} \log p(\mathbf{x}; \boldsymbol{\theta}) \nabla_{\boldsymbol{\theta}}^T \log p(\mathbf{x}; \boldsymbol{\theta}) \right], \ \boldsymbol{\theta} \in \Omega_{\boldsymbol{\theta}}.$$
 (18)

In Appendix A it is shown that the CCRB on the trace MSE under the considered model is given by

$$\sum_{m=1}^{M} \mathrm{E}_{\tilde{\boldsymbol{\theta}}} \left[(\hat{\theta}_m - \tilde{\theta}_m)^2 \right] \ge B^{\mathrm{CCRB}}(\tilde{\boldsymbol{\theta}}), \tag{19}$$

where

$$B^{\text{CCRB}}(\boldsymbol{\theta}) \stackrel{\triangle}{=} \frac{1}{N} \text{trace}\left(\left(\mathbf{U}^T \left(\text{diag}(\boldsymbol{\theta})\right)^{-1} \mathbf{U}\right)^{-1}\right). \quad (20)$$

However, missing-mass estimators, such as the Good-Turing and add-constant estimators, are χ -biased. Thus, the performance of these estimators should be assessed by the biased CCRB [43], which is a function of the estimator's bias gradient. Moreover, the CCRB in (20) is a lower bound on the MSE of estimators of the entire pmf vector, and does not provide a relevant bound on the performance for the missing-mass estimation problem. This is similar to the mismatch of the naive CML estimator from (13), which tends to overestimate the probability of the observed elements. In the following section, we develop a new CCRB-type bound on the missingmass estimation.

C. Non-Bayesian paradigm and mmMSE risk

In this subsection, we explain the rationale behind the considered approach, which is: 1) purely non-Bayesian estimation of a *deterministic* parameter; 2) an estimation of a parameter of interest in the presence of nuisance parameters that affect the accuracy of estimation; 3) based on the estimation of the entire pmf, θ ; and 4) based on the mmMSE risk.

The missing mass, namely the total probability mass of the outcomes not observed in the samples in x, is defined as

$$p_0(\mathbf{x}, \boldsymbol{\theta}) = \sum_{m=1}^M \theta_m \mathbb{1}_{\{m \in G_{N,0}(\mathbf{x})\}}.$$
 (21)

The missing mass in (21) is a hybrid (mixture of random and deterministic) scalar parameter, which is a function of both the *deterministic* pmf vector, $\boldsymbol{\theta}$, and the *random* observation vector, \mathbf{x} . Thus, various papers in the literature (see, e.g. [31, 44]) treat the estimation problem as the estimation of the *hybrid* parameter, $p_0(\mathbf{x}, \boldsymbol{\theta})$, which allegedly has both random and deterministic parts. However, since the random observation vector, \mathbf{x} , is known, the true unknown part in $p_0(\mathbf{x}, \boldsymbol{\theta})$ is only the deterministic vector, $\boldsymbol{\theta}$. Therefore, in this work we adopt the non-Bayesian approach for the estimation of *deterministic* parameters. Moreover, since all the elements of the pmf vector, $\boldsymbol{\theta}$, are unknown, we treat this estimation problem as the estimation of the parameters of interest in (21) that include the probabilities of unseen events, and refer to the other (seen) parameters in $\boldsymbol{\theta}$ as nuisance parameters [45, 46].

Direct calculation of the MSE of $p_0(\mathbf{x}, \boldsymbol{\theta})$ from (21),

$$\mathbf{E}_{\boldsymbol{\theta}}\left[\left(\sum_{m=1}^{M} \hat{\theta}_m \mathbb{1}_{\{m \in G_{N,0}(\mathbf{x})\}} - \sum_{m=1}^{M} \theta_m \mathbb{1}_{\{m \in G_{N,0}(\mathbf{x})\}}\right)^2\right],\$$

requires the calculation of all cross-correlations of estimation errors of any θ_m and θ_l , $l, m \in G_{N,0}(\mathbf{x})$. This, in turn, requires computing the expectation of a sum of 2^M possible events (that represent the binary options that m and/or l is within/without $G_{N,0}(\mathbf{x})$, for any \mathbf{x} and any $m, l = 1, \ldots, M$). While this approach is feasible for calculating the MSE of specific missing-mass estimators (see, e.g. in [31, 35]), we found it infeasible for the calculation of the associated modified FIM and unbiasedness, which leads to an intractable bound.

In order to capture the relevant errors both meaningfully and in a way that can be easily computed, we use here an alternative cost function, which is based on the missing-mass squared-error cost function:

$$C(\hat{\boldsymbol{\theta}}, \boldsymbol{\theta}) \stackrel{\Delta}{=} \sum_{m=1}^{M} (\hat{\theta}_m - \theta_m)^2 \mathbb{1}_{\{m \in G_{N,0}(\mathbf{x})\}}, \qquad (22)$$

for any estimator $\hat{\boldsymbol{\theta}} = [\hat{\theta}_1, \dots, \hat{\theta}_M]^T$ of the pmf vector, $\boldsymbol{\theta} = [\theta_1, \dots, \theta_M]^T$. The associated mmMSE risk, which is the expected value of (22), is

$$E_{\boldsymbol{\theta}}\left[C(\hat{\boldsymbol{\theta}}, \boldsymbol{\theta})\right] = \sum_{m=1}^{M} E_{\boldsymbol{\theta}}\left[(\hat{\theta}_m - \theta_m)^2 \mathbb{1}_{\{m \in G_{N,0}(\mathbf{x})\}}\right]$$
$$= \sum_{m=1}^{M} E_{\boldsymbol{\theta}}\left[(\hat{\theta}_m - \theta_m)^2 | \mathbf{x} \in \mathcal{A}_m\right] \Pr(\mathbf{x} \in \mathcal{A}_m; \boldsymbol{\theta}), \quad (23)$$

where the last equality is obtained by using the law of total probability and the conditional distribution from (11).

It can be seen that in order to evaluate the mmMSE performance of the estimator of $p_0(\mathbf{x}, \boldsymbol{\theta})$ over all possible observation vectors x, we need to sum over the errors of all the elements of θ (i.e. to compute M terms). This is a significant reduction in computational cost compared with the direct calculation of the MSE that requires computing the expectation over 2^M possible events. In addition, since we assume that M is known and finite, the sum in (22) is finite. This cost takes into account all possible estimation errors by summing over all the errors in a similar manner to existing different bounds on various cost functions (see, e.g. in [31, 35]) and to performance evaluation of specific estimators (see, e.g. in [28, 47]). The use of the indicator functions in the missing-mass squarederror cost function in (22) implies that the error of the *m*th parameter, $\theta_m - \theta_m$, affects the mmMSE only for observations x such that s_m has not been observed. Thus, it can be seen that the mmMSE is the sum of the MSEs of the parameters of interest, i.e. only the estimation errors of elements with the indices that are in $G_{N,0}(\mathbf{x})$ from (8). It should be noted that other parametric statistical analyses of the missing-mass estimation in the literature are based on the estimation of the entire pmf vector, θ . For example, in [28], the full estimator of θ is used to analyze the bias of missing-mass estimators. Similarly, in the development of the minimax bound in [35], a tight bound is derived by replacing the problem of missingmass estimation with that of distribution estimation. Further, some missing-mass estimators are based on estimating θ and then applying different smoothing approaches to obtain the estimator of the missing mass (see, e.g. [8]). Finally, it should be noted that the mmMSE is computed where the expectation is only over the randomness in the estimator, since the pmf is a deterministic vector in the considered model.

III. MISSING-MASS CONSTRAINED CRAMÉR-RAO (MMCCRB) BOUND

In this section, a CCRB-type lower bound is derived. In Subsection III-A we develop the uniform and local unbiasedness in the Lehmann sense under the missing-mass squarederror cost function and under the probability-space parametric constraints. In Subsection III-B, we derive the proposed bound, which is a lower bound on the mmMSE and is a function of the Lehmann bias of the estimators. For the sake of generality, the unbiasedness and the mmCCRB are first derived for a general observation-model distribution, $p(\mathbf{x}; \boldsymbol{\theta})$. Thus, the missingmass unbiasedness in Subsection III-A and the mmCCRB in Subsection III-B can be used for various variations of the missing-mass estimation problem, such as estimating an unknown Markov chain from its sample [33, 48, 49]. Then, in Subsection III-C, we develop the closed-form mmCCRB for the classical i.i.d. model, given in (5), as well as the mmCCRB for missing-mass unbiased estimators. Finally, in Subsection III-D we present some special cases of the mmCCRB.

A. Lehmann unbiasedness

The mean-unbiasedness constraint is commonly used in non-Bayesian parameter estimation [50]. Lehmann [42] proposed a generalization of the unbiasedness concept, which is based on the considered cost function, as follows.

Definition 2: An estimator $\hat{\boldsymbol{\theta}} : S^N \to \mathbb{R}^M$ is an unbiased estimator of $\boldsymbol{\theta}$ in the Lehmann sense [42] w.r.t. a given cost function, $C(\hat{\boldsymbol{\theta}}, \boldsymbol{\theta})$, if

$$\mathbf{E}_{\boldsymbol{\theta}}[C(\hat{\boldsymbol{\theta}}, \boldsymbol{\eta})] \ge \mathbf{E}_{\boldsymbol{\theta}}[C(\hat{\boldsymbol{\theta}}, \boldsymbol{\theta})], \ \forall \boldsymbol{\eta}, \boldsymbol{\theta} \in \Omega_{\boldsymbol{\theta}}, \tag{24}$$

where the simplex Ω_{θ} is the parameter space.

The Lehmann unbiasedness definition implies that an estimator is unbiased if, on average, it is "closer" to the true parameter θ than to any other value in the parameter space (here, denoted by an arbitrary vector, η). The measure of closeness is determined by the considered cost function, $C(\hat{\theta}, \theta)$. Examples for Lehmann unbiasedness with different cost functions and under parametric constraints can be found in [13, 36, 37, 42, 51, 52]. The following lemma states the Lehmann unbiasedness for the estimation problem of missing-mass probability. To this end, we define the elements of the missing-mass bias vector, $\mathbf{b}_{N,0}(\theta) \in \mathbb{R}^M$, as follows:

$$\begin{bmatrix} \mathbf{b}_{N,0}(\boldsymbol{\theta}) \end{bmatrix}_m \stackrel{\triangle}{=} \mathbf{E}_{\boldsymbol{\theta}} \begin{bmatrix} (\hat{\theta}_m - \theta_m) \mathbb{1}_{\{m \in G_{N,0}(\mathbf{x})\}} \end{bmatrix}$$
$$= \mathbf{E}_{\boldsymbol{\theta}} \begin{bmatrix} \hat{\theta}_m - \theta_m | \mathbf{x} \in \mathcal{A}_m \end{bmatrix} \Pr(\mathbf{x} \in \mathcal{A}_m; \boldsymbol{\theta}), \quad (25)$$

 $\forall m = 1, \dots, M$, where the last equality is obtained by using the law of total probability.

Lemma 1: An estimator $\hat{\boldsymbol{\theta}} : S^N \to \Omega_{\boldsymbol{\theta}}$ is said to be a *uniformly* Lehmann-unbiased estimator of $\boldsymbol{\theta} \in \Omega_{\boldsymbol{\theta}}$ w.r.t. the missing-mass squared-error cost function from (22) if

$$\mathbf{U}^T \mathbf{b}_{N,0}(\boldsymbol{\theta}) = \mathbf{0}_{M-1}, \ \forall \boldsymbol{\theta} \in \Omega_{\boldsymbol{\theta}},$$
(26)

where U and $\mathbf{b}_{N,0}(\boldsymbol{\theta})$ are defined in (4) and (25), respectively. *Proof:* The proof appears in Appendix B.

The CCRB is a *local* bound, meaning that it determines the achievable performance at a particular value of θ , denoted here

by $\hat{\theta}$, based on the statistics in its neighborhood. Similar to the local χ -unbiasedness in Definition 1, we can define the *local* missing-mass unbiasedness as follows.

Definition 3: An estimator $\hat{\theta} : S^N \to \Omega_{\theta}$ is said to be a *locally* Lehmann-unbiased estimator [42] in the neighborhood of $\tilde{\theta} \in \Omega_{\theta}$ w.r.t. the missing-mass squared-error cost function from (22) if it satisfies

$$\mathbf{U}^T \mathbf{b}_{N,0}(\hat{\boldsymbol{\theta}}) = \mathbf{0}_{M-1} \tag{27}$$

and

$$\left\{\nabla_{\boldsymbol{\theta}}^{T} \mathbf{b}_{N,0}(\boldsymbol{\theta})\right\}\Big|_{\boldsymbol{\theta}=\tilde{\boldsymbol{\theta}}} \mathbf{U} = \mathbf{0}_{M \times (M-1)}.$$
(28)

It should be noted that the condition in (26) requires a *uniform* unbiasedness, for any $\theta \in \Omega_{\theta}$, while the conditions in (27) and (28) are *local* conditions that are required to be satisfied only at the specific θ , denoted here by $\hat{\theta}$, for which the bound is developed. Both the local and the uniform missing-mass unbiasedness definitions restrict only the values that belong to the set of unseen symbols, i.e. elements that belong to the set $G_{N,0}(\mathbf{x})$ from (8), in S to be unbiased. In addition, by comparing Definition 1 and Definition 3 it can be seen that the differences between the local χ -unbiasedness and the local missing-mass unbiasedness follow from the difference of the cost functions, where both definitions use the null-space matrix, U, which is due to the parametric constraint. However, while there exist estimators that are χ unbiased, such as the CML estimator as shown in (16), there are no estimators that are missing-mass unbiased in the nonasymptotic region, without splitting the data or taking extra draws [44, 53, 54]. It is known that even when unbiased methods do not exist in a particular setting, such as in the case of various nonlinear models, meaningful biased techniques with good performance can still be found [43, 55-57].

The uniform missing-mass unbiasedness from (26) can be interpreted as follows. From the definition of U in (4), it can be verified that $\mathbf{U}^T \mathbf{y} = \mathbf{0}_{M-1}$ iff $\mathbf{y} = c\mathbf{1}_M$, where $c \in \mathbb{R}$ is an arbitrary constant. Thus, the condition in (26) implies that for a uniformly Lehmann unbiased estimator, the missing-mass bias vector satisfies

$$\mathbf{b}_{N,0}(\boldsymbol{\theta}) = \beta_{N,0}(\boldsymbol{\theta})\mathbf{1}_M, \ \forall \boldsymbol{\theta} \in \Omega_{\boldsymbol{\theta}},$$
(29)

where $\beta_{N,0}(\theta) \in \mathbb{R}$ is a constant. The condition in (29) is that the *m*th element of the missing-mass bias is identical for any *m*. This property recalls the notion of natural estimators [27], since it assigns the same bias requirements to all symbols appearing with the same probability.

B. mmCCRB

Lower bounds on the mmMSE are useful for performance analysis and system design. In this subsection, a constrained Cramér-Rao-type lower bound on the mmMSE from (23) is derived. The new bound is based on the missing-mass bias in the Lehmann sense, as defined in Subsection III-A. Thus, it is a bound on the MSE of the missing mass of all estimators having a given bias function, $\mathbf{b}_{N,0}(\tilde{\boldsymbol{\theta}})$, at each point, $\tilde{\boldsymbol{\theta}} \in \Omega_{\boldsymbol{\theta}}$.

Let us define the following missing-mass Fisher information matrix (mmFIM) :

$$\mathbf{J}^{(0)}(\boldsymbol{\theta}) \stackrel{\Delta}{=} \mathbf{E}_{\boldsymbol{\theta}} \left[\boldsymbol{\Delta}(\mathbf{x}, \boldsymbol{\theta}) \boldsymbol{\Delta}^{T}(\mathbf{x}, \boldsymbol{\theta}) \right], \tag{30}$$

 $\forall \boldsymbol{\theta} \in \Omega_{\boldsymbol{\theta}}$, where the *m*th column of the matrix $\boldsymbol{\Delta}(\mathbf{x}, \boldsymbol{\theta}) \in \mathbb{R}^{M \times M}$ is defined as

$$\boldsymbol{\Delta}_{1:M,m}(\mathbf{x},\boldsymbol{\theta}) \stackrel{\triangle}{=} \nabla_{\boldsymbol{\theta}} \log p(\mathbf{x}|\mathbf{x} \in \mathcal{A}_m;\boldsymbol{\theta}) \mathbb{1}_{\{\mathbf{x} \in \mathcal{A}_m\}}, \quad (31)$$

 $m = 1, \dots, M$. In addition, we define the auxiliary matrix $\mathbf{S}(\boldsymbol{\theta}) \in \mathbb{R}^{M \times M}$, in which the *m*th row is defined as

$$\mathbf{S}_{m,1:M}(\boldsymbol{\theta}) \\ \stackrel{\triangle}{=} \left(\nabla_{\boldsymbol{\theta}}^{T} \left\{ \frac{\left[\mathbf{b}_{N,0}(\boldsymbol{\theta}) \right]_{m}}{\Pr(\mathbf{x} \in \mathcal{A}_{m}; \boldsymbol{\theta})} \right\} + \mathbf{e}_{m}^{T} \right) \Pr(\mathbf{x} \in \mathcal{A}_{m}; \boldsymbol{\theta}), \quad (32)$$

 $\forall \boldsymbol{\theta} \in \Omega_{\boldsymbol{\theta}}.$

We define the following regularity condition:

C.1) The likelihood gradient vector, $\Delta_{1:M,m}(\mathbf{x}, \boldsymbol{\theta})$, defined in (31), exists and is finite $\forall \boldsymbol{\theta} \in \Omega_{\boldsymbol{\theta}}$ and $\forall m = 1, \dots, M$. That is, the matrix $\mathbf{U}^T \mathbf{J}^{(0)}(\boldsymbol{\theta})\mathbf{U}$ is a well-defined, non-singular, and non-zero matrix for any $\boldsymbol{\theta} \in \Omega_{\boldsymbol{\theta}}$.

Theorem 1: Let the regularity condition C.1 be satisfied and $\hat{\theta}$ be an estimator of $\theta \in \Omega_{\theta}$ with a local missing-mass bias vector in the neighborhood of $\tilde{\theta} \in \Omega_{\theta}$ given by $\mathbf{b}_{N,0}(\tilde{\theta})$, as defined in (25). Then, the mmMSE from (23) satisfies

$$\mathbf{E}_{\tilde{\boldsymbol{\theta}}}\left[C(\hat{\boldsymbol{\theta}}, \tilde{\boldsymbol{\theta}})\right] \ge B^{\mathrm{mmCCRB}}(\tilde{\boldsymbol{\theta}}), \tag{33}$$

where the mmCCRB evaluated at the local point, θ , is

$$B^{\text{mmCCRB}}(\tilde{\boldsymbol{\theta}}) \stackrel{\triangle}{=} \text{trace} \left(\mathbf{S}^{T}(\tilde{\boldsymbol{\theta}}) \mathbf{U} (\mathbf{U}^{T} \mathbf{J}^{(0)}(\tilde{\boldsymbol{\theta}}) \mathbf{U})^{-1} \mathbf{U}^{T} \mathbf{S}(\tilde{\boldsymbol{\theta}}) \right) \\ + \sum_{m=1}^{M} \frac{\left[\mathbf{b}_{N,0}(\tilde{\boldsymbol{\theta}}) \right]_{m}^{2}}{\Pr(\mathbf{x} \in \mathcal{A}_{m}; \tilde{\boldsymbol{\theta}})}.$$
(34)

Moreover, equality is achieved in (34) if

$$\hat{\theta}_{m} - \tilde{\theta}_{m} = \frac{[\mathbf{b}_{N,0}(\tilde{\boldsymbol{\theta}})]_{m}}{\Pr(\mathbf{x} \in \mathcal{A}_{m}; \tilde{\boldsymbol{\theta}})} + \left[\mathbf{S}^{T}(\tilde{\boldsymbol{\theta}})\mathbf{U}(\mathbf{U}^{T}\mathbf{J}^{(0)}(\tilde{\boldsymbol{\theta}})\mathbf{U})^{-1}\mathbf{U}^{T}\boldsymbol{\Delta}(\mathbf{x}, \tilde{\boldsymbol{\theta}})\right]_{m,m}, \quad (35)$$

for any m = 1, ..., M such that $m \in G_{N,0}(\mathbf{x})$. *Proof:* The proof appears in Appendix C.

Theorem 1 provides a lower bound on the MSE of missingmass estimators that have a specified bias function, $\mathbf{b}_{N,0}(\tilde{\boldsymbol{\theta}})$. This is similar to MSE bounds on biased estimators in the general setting [56, 57]. Biased bounds can be used to explore the fundamental tradeoff between bias and variance, as well as for system design. The specification of the biased mmCCRB requires an *a-priori* choice of the bias gradient. The biased mmCCRB can be used for cases where we consider an estimator with a tractable bias gradient. For the case of the i.i.d. model, we show in Lemma 2 that the biased mmCCRB can be computed without the need for a bias gradient, with simple expectation terms that make the biased mmCCRB more tractable. It should be noted that in the following we use the same notation, $B^{\text{mmCCRB}}(\tilde{\boldsymbol{\theta}})$, for the mmCCRB with different bias specifications.

It can be seen that the equality condition in (35), which is the requirement for the achievability of the mmCCRB, only determines the values of the missing-mass estimation errors. In addition, it can be seen that the estimator defined in (35), $\hat{\theta}_m$, $m = 1, \ldots, M$, may assign a different value for each element in the missing mass. That is, it is not necessarily a natural estimator [27], in the sense that elements that appeared the same number of times will not necessarily get the same estimated probability. Moreover, this estimator is a function of the (local) unknown parameter vector, $\tilde{\theta}$, in the general case. Only if it is independent of $\tilde{\theta}$, then it is an efficient estimator and its mmMSE is equal to the mmCCRB. The equality condition of the mmCCRB in (35) is the basis for a new estimation method developed in Section IV.

The mmFIM in (30) is a function of the entire pmf, θ . It can be verified that $\mathbf{J}^{(0)}(\boldsymbol{\theta})$ is not a diagonal matrix (see also in (39) below). This is because there is a coupling between the different elements in θ , and the estimation of one parameter affects the accuracy in the estimation of the others. Finally, it has been shown in recent works that the profile likelihood, i.e. the empirical distribution up to permutation of the lexicon, can considered to be a sufficient statistic for the problem of missing-mass estimation [24-26, 47]. Thus, in general, a new bound can be developed based on the likelihood of the profile instead of the likelihood in (31). According to the extension of the data processing inequality for Fisher information [58], such a lower bound may result in a tighter bound than the mmCCRB. However, it is a valid lower bound only on profilebased estimators. In addition, the derivation of such a bound is not straightforward since the entire pmf, θ , cannot be estimated based on the profile. Thus, we leave this topic for further investigation.

C. mmCCRB for the i.i.d. model

The mmCCRB in Theorem 1 is a lower bound on the mmMSE from (23), which has been developed for the general observation model, $p(\mathbf{x}; \boldsymbol{\theta}), \boldsymbol{\theta} \in \Omega_{\boldsymbol{\theta}}$. In this subsection, we develop the closed-form expression of the mmCCRB for the classical i.i.d. model, as described by (5). In addition, we develop the mmCCRB for the special case of missing-mass unbiased estimators for the classical model described by (5).

The following corollary describes the closed-form mmFIM for the i.i.d. case.

Corollary 1: Let the conditions of Theorem 1 be satisfied and assume the model described in Subsection II-A with the observation pmf given in (5). Then, the mmCCRB for this model is

$$B^{\text{mmCCRB}}(\tilde{\boldsymbol{\theta}}) = \frac{1}{N} \text{trace} \left(\mathbf{S}^{T}(\tilde{\boldsymbol{\theta}}) \mathbf{U} (\mathbf{U}^{T} \mathbf{D}(\tilde{\boldsymbol{\theta}}) \mathbf{U})^{-1} \mathbf{U}^{T} \mathbf{S}(\tilde{\boldsymbol{\theta}}) \right) + \sum_{m=1}^{M} \frac{\left[\mathbf{b}_{N,0}(\tilde{\boldsymbol{\theta}}) \right]_{m}^{2}}{\Pr(\mathbf{x} \in \mathcal{A}_{m}; \tilde{\boldsymbol{\theta}})},$$
(36)

where $\mathbf{D}(\boldsymbol{\theta})$ is a $M \times M$ diagonal matrix with the following elements on its diagonal:

$$[\mathbf{D}(\boldsymbol{\theta})]_{m,m} = -\frac{(1-\theta_m)^N}{(1-\theta_m)^2} + \frac{1}{\theta_m} \sum_{l=1,l\neq m}^M \frac{(1-\theta_l)^N}{1-\theta_l}, \quad (37)$$

ъ*л*

 $m = 1, \ldots, M$. The associated equality condition is given by

$$\hat{\theta}_{m} - \tilde{\theta}_{m} = \frac{[\mathbf{b}_{N,0}(\boldsymbol{\theta})]_{m}}{\Pr(\mathbf{x} \in \mathcal{A}_{m}; \tilde{\boldsymbol{\theta}})} + \frac{1}{N} \left[\mathbf{S}^{T}(\tilde{\boldsymbol{\theta}}) \mathbf{U} (\mathbf{U}^{T} \mathbf{D}(\tilde{\boldsymbol{\theta}}) \mathbf{U})^{-1} \mathbf{U}^{T} \boldsymbol{\Delta}(\mathbf{x}, \tilde{\boldsymbol{\theta}}) \right]_{m,m}, \quad (38)$$

for any $m = 1, \ldots, M$ such that $m \in G_{N,0}(\mathbf{x})$,

Proof: In Appendix D, it is proved that the mmFIM from (30) for the model described in Subsection II-A with the observation pmf in (5) is:

$$\mathbf{J}^{(0)}(\boldsymbol{\theta}) = \sum_{m=1}^{M} \frac{N(N-1)}{(1-\theta_m)^2} (1-\theta_m)^N \mathbf{1}_M \mathbf{1}_M^T + \sum_{m=1}^{M} \frac{N(1-\theta_m)^N}{(1-\theta_m)^2} \left(\mathbf{e}_m \mathbf{1}_M^T + \mathbf{1}_M \mathbf{e}_m^T\right) + N\mathbf{D}(\boldsymbol{\theta}), (39)$$

where $\mathbf{D}(\boldsymbol{\theta})$ is defined in (37). By using (39) and the null-space property of the matrix \mathbf{U} from (4), $\mathbf{1}_{M}^{T}\mathbf{U} = \mathbf{0}^{T}$, we obtain

$$\mathbf{U}^T \mathbf{J}^{(0)}(\boldsymbol{\theta}) \mathbf{U} = N \mathbf{U}^T \mathbf{D}(\boldsymbol{\theta}) \mathbf{U}.$$
 (40)

By substituting (40) in (34), we obtain the mmCCRB for the classical model in (36). Moreover, by substituting (40) in (35), we obtain that the equality condition of the mmCCRB for the classical model is given by (38).

The auxiliary matrix, $\mathbf{S}(\boldsymbol{\theta})$, in (32) involves the gradient of the missing-mass bias, which makes it intractable for many estimators with implicit bias gradient function. The following lemma presents a tractable form of the auxiliary matrix for the i.i.d. case, which can be evaluated numerically. This tractable form is a function of the bias of the estimator and of its correlation with the empirical histogram of the observations.

Lemma 2: The *m*th row of the auxiliary matrix $S(\theta)$ from (32) under the model described in Subsection II-A with the observation pmf given in (5) can be calculated as

$$\mathbf{S}_{m,1:M}(\boldsymbol{\theta}) = \mathbf{E}_{\boldsymbol{\theta}} \left[(\hat{\theta}_m(\mathbf{x}) - \theta_m) \mathbf{v}^T(\mathbf{x}, \boldsymbol{\theta}) \mathbb{1}_{\{\mathbf{x} \in \mathcal{A}_m\}} \right] \\ + \frac{N}{1 - \theta_m} [\mathbf{b}_{N,0}(\boldsymbol{\theta})]_m (\mathbf{e}_m - \mathbf{1}_M)^T, \qquad (41)$$

where

$$\mathbf{v}(\mathbf{x},\boldsymbol{\theta}) \triangleq \left[\frac{C_{N,1}(\mathbf{x})}{\theta_1}, \dots, \frac{C_{N,M}(\mathbf{x})}{\theta_M}\right]^T.$$
(42)

Proof: The proof appears in Appendix E.

By substituting the auxiliary matrix, $S(\theta)$, from Lemma 2 in (36), we obtained a tractable version of the mmCCRB on the mmMSE of biased estimators. In contrast with the traditional biased CRB, which requires *a priori* specification of the desired bias gradient [56, 57], here, we need to specify the expectations in (41), which enables an numerical calculation if needed.

In the following, we describe the mmCCRB for missingmass unbiased estimators. For the sake of simplicity of derivation, we assume in the following that $\mathbf{b}_{N,0}(\boldsymbol{\theta}) = \mathbf{0}$. According to Lemma 1, this condition is a *sufficient* condition for the Lehmann unbiasedness in (27) and (15). Corollary 2: Let the conditions of Theorem 1 be satisfied and assume the model described in Subsection II-A with $\mathbf{b}_{N,0}(\boldsymbol{\theta}) = \mathbf{0}$. Then, the mmCCRB for this model and missingmass unbiased estimators is:

$$B^{\text{mmCCRB}}(\tilde{\boldsymbol{\theta}}) = \frac{1}{N} \sum_{m=1}^{M} (1 - \tilde{\theta}_m)^{2N} \left[\mathbf{U} (\mathbf{U}^T \mathbf{D}(\tilde{\boldsymbol{\theta}}) \mathbf{U})^{-1} \mathbf{U}^T \right]_{m,m}, \quad (43)$$

where $\mathbf{D}(\boldsymbol{\theta})$ is defined in (37).

Proof: By substituting $\mathbf{b}_{N,0}(\boldsymbol{\theta}) = \mathbf{0}$ in (32), we obtain that in this case $\mathbf{S}(\boldsymbol{\theta})$ is a diagonal matrix with the diagonal elements

$$[\mathbf{S}(\boldsymbol{\theta})]_{m,m} = \Pr(\mathbf{x} \in \mathcal{A}_m; \boldsymbol{\theta}) = (1 - \theta_m)^N, \qquad (44)$$

 $m = 1, \ldots, M$, where the last equality is obtained by substituting (10). By substituting $\mathbf{b}_{N,0}(\boldsymbol{\theta}) = \mathbf{0}$ and (44) in (36), we obtain that the mmCCRB on the mmMSE of a missing-mass unbiased estimator is

$$B^{\text{mmCCRB}}(\tilde{\boldsymbol{\theta}}) = \frac{1}{N} \text{trace} \left(\mathbf{U}^T \mathbf{P}(\tilde{\boldsymbol{\theta}}) \mathbf{U} (\mathbf{U}^T \mathbf{D}(\tilde{\boldsymbol{\theta}}) \mathbf{U})^{-1} \right), \quad (45)$$

where $\mathbf{P}(\boldsymbol{\theta})$ is a diagonal $M \times M$ matrix with the following elements on its diagonal:

$$[\mathbf{P}(\tilde{\boldsymbol{\theta}})]_{m,m} \stackrel{\triangle}{=} (1 - \tilde{\theta}_m)^{2N}, \ m = 1, \dots, M.$$
(46)

By substituting (46) in (45) and using the trace operator properties, the mmCCRB for this case is given by (43).

The main advantage of the mmCCRB for missing-mass unbiased estimators in Corollary 2, is that it is only a function of the symbol generation system via the true pmf, θ , and the number of observations, N. While the minimax MSE is lowerbounded by $\frac{c}{N}$ for a constant c [31, 35], the lower bound on the mmMSE provided by the mmCCRB in (43) has a more complicated structure as a function of N. It should be noted also that the minimax MSE approach is derived for a specific algorithm (e.g. Good-Turing estimator) or for the worst-case pmf, while the proposed mmCCRB applies to all algorithms and should be evaluated for each value of θ . Finally, it can be seen that for each different pmf, θ , the bound in (43) requires only the computation of the diagonal matrix $\mathbf{D}(\theta)$, which is defined in (37).

D. Special cases

In this subsection, we develop some important special cases of the mmCCRB for the i.i.d. model from Subsection III-C.

1) mmCCRB on the mmMSE of the CML estimator: The CML estimator from (13) assigns a zero probability to unseen events:

$$\hat{\theta}_m^{\text{CML}} = 0, \ \forall m \in G_{N,0}(\mathbf{x}).$$
 (47)

By substituting (47) in (25), one obtains that the missing-mass bias of the CML estimator satisfies

$$[\mathbf{b}_{N,0}^{\mathrm{CML}}(\boldsymbol{\theta})]_{m} = \mathbf{E}_{\boldsymbol{\theta}}[\hat{\theta}_{m}^{\mathrm{CML}} - \theta_{m} | \mathbf{x} \in \mathcal{A}_{m}] \operatorname{Pr}(\mathbf{x} \in \mathcal{A}_{m}; \boldsymbol{\theta})$$
$$= -\theta_{m} \operatorname{Pr}(\mathbf{x} \in \mathcal{A}_{m}; \boldsymbol{\theta}), \qquad (48)$$

for any m = 1, ..., M. By substituting (48) in (32), we obtain that for the CML estimator, the auxiliary matrix satisfies

 $\mathbf{S}(\boldsymbol{\theta}) = \mathbf{0}_{M \times M}$. By substituting this result and (48) in (36), we obtain that the mmCCRB on the mmMSE of the CML estimator (or any other estimator with the same bias function as in (48)) is

$$B^{\text{mmCCRB}}(\tilde{\boldsymbol{\theta}}) = \sum_{m=1}^{M} \tilde{\theta}_m^2 \Pr(\mathbf{x} \in \mathcal{A}_m; \tilde{\boldsymbol{\theta}}).$$
(49)

On the other hand, by substituting (47) in (23), it can be seen that the mmMSE of the CML estimator evaluated at the local point, $\hat{\theta}$, is

$$\mathbf{E}_{\tilde{\boldsymbol{\theta}}}\left[C(\hat{\boldsymbol{\theta}}^{\mathrm{CML}}, \tilde{\boldsymbol{\theta}})\right] = \sum_{m=1}^{M} \tilde{\theta}_{m}^{2} \operatorname{Pr}(\mathbf{x} \in \mathcal{A}_{m}; \tilde{\boldsymbol{\theta}}).$$
(50)

Thus, in this case, the biased mmCCRB with the bias function of the CML estimator coincides with the mmMSE of the CML estimator. Therefore, we can conclude that there is no other estimator with the same missing-mass bias as that of the CML estimator, given in (48), that achieves a lower mmMSE than the CML estimator. It should be noted that since we used the mmCCRB from Theorem 1, which was developed for the general observation model, $p(\mathbf{x}; \boldsymbol{\theta}), \boldsymbol{\theta} \in \Omega_{\boldsymbol{\theta}}$, this result also holds for non-i.i.d. sampling with a general structure of $p(\mathbf{x}; \boldsymbol{\theta})$ [33, 48, 49].

2) mmCCRB on the mmMSE of a general missing-mass estimator: Various estimators of the missing mass have been suggested in the literature. In this paper we consider three estimators that are described in Section V: The Good-Turing, Laplace, and aPML estimators. In the following, we describe how to obtain the biased mmCCRB from Corollary 1 for a general missing-mass estimator.

Consider an estimator $\hat{p}_0(\mathbf{x}, \boldsymbol{\theta})$ for the missing mass from (21). Under the assumption of a natural estimator [27], the associated estimator of a specific element in $G_{N,0}(\mathbf{x})$ is

$$\hat{\theta}_m = \begin{cases} \frac{\hat{p}_0(\mathbf{x}, \hat{\boldsymbol{\theta}})}{|G_{N,0}(\mathbf{x})|} & \text{if } G_{N,0}(\mathbf{x}) \neq \emptyset \\ 0 & \text{if } G_{N,0}(\mathbf{x}) = \emptyset \end{cases},$$
(51)

for any $m \in G_{N,0}(\mathbf{x})$. By substituting (51) in (25), one obtains that the missing-mass bias of an arbitrary estimator $\hat{p}_0(\mathbf{x}, \boldsymbol{\theta})$ satisfies

-

$$\begin{bmatrix} \mathbf{b}_{N,0}(\boldsymbol{\theta}) \end{bmatrix}_{m} = \mathbf{E}_{\boldsymbol{\theta}} \left[\frac{\hat{p}_{0}(\mathbf{x}, \hat{\boldsymbol{\theta}})}{|G_{N,0}(\mathbf{x})|} - \theta_{m} | \mathbf{x} \in \mathcal{A}_{m} \right] \Pr(\mathbf{x} \in \mathcal{A}_{m}; \boldsymbol{\theta}). \quad (52)$$

While a closed-form expression of the missing-mass bias of an arbitrary estimator in (52) is intractable, the proposed bound can be used by numerically calculating the auxiliary matrix for this case. That is, by substituting (51) in (41), one obtains

$$\mathbf{S}_{m,1:M}(\boldsymbol{\theta}) = \mathbf{E}_{\boldsymbol{\theta}} \left[\left(\frac{\hat{p}_0(\mathbf{x}, \hat{\boldsymbol{\theta}})}{|G_{N,0}(\mathbf{x})|} - \theta_m \right) \mathbf{v}^T(\mathbf{x}, \boldsymbol{\theta}) | \mathbf{x} \in \mathcal{A}_m \right] \Pr(\mathbf{x} \in \mathcal{A}_m; \boldsymbol{\theta}) + \frac{N}{1 - \theta_m} [\mathbf{b}_{N,0}(\boldsymbol{\theta})]_m (\mathbf{e}_m - \mathbf{1}_M)^T.$$

Then, by substituting (52) and (53) in (36), we obtain the associated mmCCRB, which can be evaluated numerically. In particular, a Monte Carlo approach can be applied to approximate the expectation in (53), in a similar manner to the empirical FIM approximation described in [59].

3) Uniform distribution: For the special case where θ = $\frac{1}{M} \mathbf{1}_M$, the diagonal elements of the matrix $\mathbf{D}(\boldsymbol{\theta})$ from (37) are given by

$$[\mathbf{D}(\boldsymbol{\theta})]_{m,m} = -\left(\frac{M-1}{M}\right)^{N-2} + M^2 \left(\frac{M-1}{M}\right)^N$$
$$= M(M-2) \left(\frac{M-1}{M}\right)^{N-2}, \qquad (54)$$

 $m = 1, \ldots, M$. Similarly, for this case (46) is reduced to

$$\mathbf{P}(\boldsymbol{\theta}) = \left(\frac{M-1}{M}\right)^{2N} \mathbf{I}_M.$$
 (55)

By substituting (54) and (55) in (45) and using $\mathbf{U}^T \mathbf{U} = \mathbf{I}$ from (4), we obtain that for this case the mmCCRB missing-mass unbiased estimator is given by

$$B^{\text{mmCCRB}}(\tilde{\boldsymbol{\theta}}) = \frac{1}{N} \left(\frac{M-1}{M}\right)^{N+3} \frac{1}{M-2},$$
 (56)

for M > 2, where we used the cyclic property of the trace. We can see that the mmCCRB for the uniform pmf from (56) decreases as the number of samples, N, increases, since we have more information. In general, the rate of decrease is a function of M. For large values of M, the mmCCRB has approximately the order of $\frac{1}{N}$, similar to the minimax results [31, 35]. The lower bounds on the minimax MSE are independent of M, which is assumed unknown in [31, 35]. As a function of M, the mmCCRB in (56) increases as Mincreases if $1 < M \leq \frac{N+4+\sqrt{N^2+4}}{2}$, and decreases as M increases otherwise. On the other hand, for this case of uniform pmf the CCRB in (20) on the trace MSE is reduced to

$$B^{\text{CCRB}}(\boldsymbol{\theta}) = \frac{1}{N} - \frac{1}{MN}.$$
(57)

Thus, the CCRB is almost independent of the number of elements, M, for large value of MN, in contrast to the missing-mass CCRB in (56), and, thus, is less informative for the problem of missing-mass estimation.

IV. MISSING-MASS FISHER-SCORING-TYPE ESTIMATION

Obtaining the minimum mmMSE estimator among all unbiased estimators is usually intractable. Moreover, in most nonlinear parameter estimation problems, such an estimator does not exist. Therefore, in this section we describe a new iterative algorithm, the missing-mass Fisher-scoring algorithm, that further improves the performance of existing estimators by using the proposed bound. Similar to the Fisher-scoring method [60] and the constrained Fisher-scoring method [12, 61], the equality condition in (35) can be used to obtain an iterative estimation procedure. In this case, the estimator at the kth iteration, $\hat{\theta}^{(k)}$, is obtained by substituting the estimator from the previous iteration, $\tilde{\boldsymbol{\theta}} = \hat{\boldsymbol{\theta}}^{(k-1)}$, in (35) to obtain

$$\hat{\theta}_{m}^{(k)} - \hat{\theta}_{m}^{(k-1)} = \psi^{(k)} \left\{ \frac{[\mathbf{b}_{N,0}(\boldsymbol{\theta})]_{m}}{\Pr(\mathbf{x} \in \mathcal{A}_{m}; \tilde{\boldsymbol{\theta}})} + \left[\mathbf{S}^{T}(\tilde{\boldsymbol{\theta}}) \mathbf{U}(\mathbf{U}^{T} \mathbf{J}^{(0)}(\tilde{\boldsymbol{\theta}}) \mathbf{U})^{-1} \mathbf{U}^{T} \boldsymbol{\Delta}(\mathbf{x}, \tilde{\boldsymbol{\theta}}) \right]_{m,m} \right\} \Big|_{\tilde{\boldsymbol{\theta}} = \hat{\boldsymbol{\theta}}^{(k-1)}}, (58)$$

for any $m = 1, \ldots, M$ such that $m \in G_{N,0}(\mathbf{x})$, where $\psi^{(k)}$ is the step size at the kth iteration. Similarly, for the i.i.d. model, the equality condition in (38) results in the following *k*th iteration:

$$\hat{\theta}_{m}^{(k)} - \hat{\theta}_{m}^{(k-1)} = \psi^{(k)} \begin{cases} \frac{[\mathbf{b}_{N,0}(\tilde{\boldsymbol{\theta}})]_{m}}{\Pr(\mathbf{x} \in \mathcal{A}_{m}; \tilde{\boldsymbol{\theta}})} & \mathbf{3} \\ \frac{1}{N} \left[\mathbf{S}^{T}(\tilde{\boldsymbol{\theta}}) \mathbf{U} (\mathbf{U}^{T} \mathbf{D}(\tilde{\boldsymbol{\theta}}) \mathbf{U})^{-1} \mathbf{U}^{T} \boldsymbol{\Delta}(\mathbf{x}, \tilde{\boldsymbol{\theta}}) \right]_{m,m} \end{cases} \begin{vmatrix} \mathbf{z} \\ \mathbf{z} \\$$

for any m = 1, ..., M such that $m \in G_{N,0}(\mathbf{x})$, where $\psi^{(k)}$ is the step size at the kth iteration. An appropriate step-size rule for $\psi^{(k)}$, k = 1, 2, ..., should be chosen to guarantee usability (such that the resulting iterate reduces the mmMSE) and to stabilize the convergence. In comparison with the classical method of Fisher-scoring or with the constrained Fisher-scoring [12], the proposed iteration in (59) is essentially a replacement of the Cramér-Rao bound (in classical Fisherscoring) or the CCRB (in the constrained Fisher-scoring) with the new mmCCRB from (36) in the Fisher-scoring iteration.

The initial estimator, $\hat{\theta}^{(0)}$, can be chosen to be any existing estimator, such as the CML, Good-Turing, Laplace, or aPML estimator, all described in Subsections V. In order to obtain reasonable estimation, the initial estimator: 1) should satisfy the constraint $\boldsymbol{\theta} \in \Omega_{\boldsymbol{\theta}}$ in (1), i.e. $\mathbf{1}_{M}^{T} \hat{\boldsymbol{\theta}}^{(0)} = 1$; and 2) should be a natural estimator [27], i.e. $\hat{\theta}^{(0)}$ assigns the same probabilities to symbols appearing the same number of times. Similarly, after each iteration, we project the solution to the constraint set and to be a natural estimator, by the following steps:

. (1)

$$\hat{\theta}_{m}^{(k)} = \frac{\hat{\theta}_{m}^{(k)}}{\sum_{l=1}^{M} \hat{\theta}_{l}^{(k)}}, \ m = 1, \dots, M$$
(60)

and

$$\hat{\theta}_{m}^{(k)} = \frac{1}{|G_{N,C_{N,m}}(\mathbf{x})|} \sum_{l \in G_{N,C_{N,m}}(\mathbf{x})}^{M} \hat{\theta}_{l}^{(k)},$$
(61)

 $m = 1, \ldots, M$. For a desired tolerance v, the algorithm exits when the condition $||\hat{\boldsymbol{\theta}}^{(k)} - \hat{\boldsymbol{\theta}}^{(k-1)}|| < v$ is met. Since the lexicon size is assumed to be known, in a case where we observe all symbols, $G_{N,0}(\mathbf{x}) = \emptyset$, we set $\hat{p}_0(\mathbf{x}, \boldsymbol{\theta}) = 0$. Finally, the algorithm is summarized in Algorithm 1.

V. SIMULATIONS

In this section, we evaluate the proposed bound and the missing-mass Fisher-scoring method. In Subsection V-A, we describe the estimators and bounds that are evaluated in the simulations. In Subsections V-B and V-C we evaluate the performance for uniform and Zipf distributions, respectively.

Algorithm 1: missing-mass Fisher-scoring algorithm for improving missing-mass estimators

Input:

- M Number of symbols
- \mathbf{x} Observation vector $\hat{\boldsymbol{\theta}}^{(0)}$ Initial estimator
- $\psi^{(k)}, k = 1, \ldots$ step sizes
- v Tolerance
- $K_{\rm max}$ Maximum iteration number

Output: $\hat{p}_0(\mathbf{x}, \boldsymbol{\theta})$ - Estimator of the missing mass 1 Initialize k = 0f C (\mathbf{v}) -() 4 h an

Return:
$$\hat{p}_0(\mathbf{x}, \boldsymbol{\theta}) = 0$$

Update $k \leftarrow k + 1$ Update $\hat{\boldsymbol{\theta}}^{(k)} \leftarrow \hat{\boldsymbol{\theta}}^{(k-1)}$ by (59) with step size $\psi^{(k)}$ Correct $\hat{\boldsymbol{\theta}}^{(k)}$ by the projection in (60) 7 Correct $\hat{\boldsymbol{\theta}}^{(k)}$ by the projection in (61) 8 if $||\hat{\boldsymbol{\theta}}^{(k)} - \hat{\boldsymbol{\theta}}^{(k-1)}|| < v$ and/or $k > K_{\max}$ then 9 **Return:** 10 $\hat{p}_0(\mathbf{x}, \boldsymbol{\theta}) = \frac{1}{|G_{N,0}(\mathbf{x})|} \sum_{m=1, m \in G_{N,0}(\mathbf{x})}^M \hat{\theta}_m^{(k)}$ (62)else 11 Repeat to step 5. 12

A. Estimators and bounds

In the following simulations, we evaluate the performance of four estimators of the missing mass:

I. CML estimator from (13).

II. Good-Turing estimator - The Good-Turing estimator [2] of the missing mass from (21) is defined as the fraction of symbols occurring exactly once in the observed samples divided by the length of the observation vector. It is well known that smoothing of the Good-Turing estimator may improve the estimation performance [2, 8]. Here we use a smooth modified version of the Good-Turing estimator, described in in [8]. This modified Good-Turing estimator is given by

$$\hat{p}_0^{GT}(\mathbf{x}, \hat{\boldsymbol{\theta}}) = \frac{\varphi(|G_{N,1}(\mathbf{x})|)}{\zeta}, \tag{63}$$

where $\varphi(t) = \max\{t, 1\}, \forall t \in \mathbb{R}, |G_{N,1}(\mathbf{x})|$ is the number of elements that appear exactly once in the N-length observation vector, x, and

$$\zeta \stackrel{\Delta}{=} \varphi(|G_{N,1}(\mathbf{x})|) + \sum_{r \in \{r: |G_{N,r}(\mathbf{x})| > 0\}} |G_{N,r}(\mathbf{x})|(r+1) \frac{\varphi(|G_{N,r+1}(\mathbf{x})|)}{|G_{N,r}(\mathbf{x})|}$$
(64)

is a normalization factor.

III. Laplace estimator - The add-constant estimator of the

missing mass from (21) is defined as [16]

$$\hat{p}_0^{add-c}(\mathbf{x}, \hat{\boldsymbol{\theta}}) = \frac{c}{N + c(M - |G_{N,0}(\mathbf{x})| + 1)},$$
(65)

for a positive constant c. The add-constant estimator has been applied and studied extensively and has been shown to have some optimality properties [8]. For c = 1, we obtain the special case of the Laplace estimator [15, 16], used in the simulations.

IV. aPML estimator- The PML estimator [24-26] has impressive statistical properties, but is computationally challenging. Consequently, an efficiently computable approximation for the PML distribution was proposed in [47]. In the simulations, we use the code from [62] for computing the aPML distributions, where the support size is used as a given input. Then, we take the smallest value in this distribution as the aPML estimator of the missing mass.

The performance of these estimators is evaluated using 500, 000 Monte-Carlo simulations that are used to evaluate the mmMSE, $E_{\theta}[C(\hat{\theta}, \theta)]$, and the absolute value of the missingmass total bias, $|\sum_{m=1}^{M} [\mathbf{b}_{N,0}(\theta)]_m|$, as defined in (23) and (25), respectively.

We compare the performance of these estimators with the following bounds:

- The CCRB from (20), which is a lower bound on the MSE of the entire pmf vector and is presented here in order to compare its behavior with that of the proposed mmMSE bounds.
- The mmCCRB with the bias of the CML estimator, as given in (49).
- Three versions of the biased mmCCRB from Corollary 1 with the empirical bias and the empirical auxiliary matrix, S(θ), as described in Subsection III-D2 for the Good-Turing, Laplace, and aPML estimators.
- The mmCCRB on missing-mass unbiased estimators from Corollary 2.

It can be verified that in all the simulations the regularity condition C.1 is satisfied.

B. Example 1: Uniform distribution

In the first experiment we examine the case of a uniform pmf with equally-likely elements, i.e. where $\theta = \frac{1}{M} \mathbf{1}_M$, as described in Subsection III-D3. In Figs. 1a and 1b we present the missing-mass bias and the mmMSE, respectively, of the different estimators versus the number of elements, M, for N = 30. Similarly, In Figs. 2a and 2b we present the missing-mass bias and the mmMSE, respectively, of the different estimators versus the number of samples, N, for M = 15. The CCRB and the mmCCRB of unbiased estimators are also presented in Figs. 1b and 2b. The performance of the aPML estimator for this case is not presented since the default of this estimator, where the observations are insufficient, is to choose the uniform estimator. Thus, the aPML estimator is unstable and not suitable for comparison in the uniform distribution case. It can be seen that for this case, the Good-Turing estimator outperforms the two other estimators in both missing-mass bias and mmMSE terms, where the gap is larger where M is larger and where N is smaller. The differences between the performance of the CML and the Laplace estimators are insignificant. In addition, it can be seen in Fig. 2b that the mmMSE of the Good-Turing estimator coincides with the proposed mmCCRB on the mmMSE of *unbiased* estimators from Corollary 2 for small values of N. While the unbiased mmCCRB is not a tight bound for this case, its curve demonstrates the influence of the system parameters, M (in Fig. 1b) and N (in Fig. 2b), on the practical mmMSE. In contrast, the CCRB is not a useful tool for performance analysis and system design. For example, it does not reflect the influence of M, as also shown analytically in Subsection III-D3.

C. Example 2: Zipf distribution

In the second experiment, we consider a Zipf's law distribution, $\theta_m = \frac{m^{-s}}{\sum_{k=1}^{M} k^{-s}}$, $m = 1, \ldots, M$, where s is the skewness parameter. The Zipf's law distribution is a heavy-tailed distribution that is widely used in physical and social sciences, linguistics, economics, and other fields [63]. In Figs. 3a and 3b we present the bias and mmMSE, respectively, of the different estimators versus the number of elements, M, for N = 100 and s = 1. Similarly, in Figs. 4a and 4b we present the bias and mmMSE, respectively, versus the number method.sof samples, N, for M = 15 and s = 1. In addition, in Figs. 3b and 4b we also present the CCRB, unbiased mmCCRB, and the different biased mmCCRBs associated with the considered estimators.

It can be seen in Figs. 3-4 that the CML estimator has the largest missing-mass bias and the largest mmMSE. Additionally, the aPML estimator has the smallest missing-mass bias and the smallest mmMSE for all N and M. In Figs. 3b and 4b we can see that the mmCCRB on the mmMSE of missing-mass unbiased estimators is lower than the actual mmMSE of theestimators, since these estimators are *biased* in the Lehmann sense. However, the curve of the mmCCRB demonstrates the influence of the system parameters, M (in Fig. 3b) and N (in Fig. 4b), on the practical mmMSE. In contrast, the CCRB does not reflect the influence of M also in this case.

It can be seen that the biased mmCCRB with the CML bias coincides with the mmMSE of the CML estimator, as shown analytically in (49)-(50). Similarly, the biased mmC-CRB associated with the Laplace estimator coincides with the mmMSE of the Laplace estimator. Thus, the CML and Laplace estimators achieve the lowest mmMSE for their associated bias function. However, for the Good-Turing and aPML estimators there is a gap between the associated mmCCRB and the mmMSE. As the sample size, N, increases, the mmMSE of these estimators achieves the associated mmCCRBs. Thus, in this case, there can be estimators with the same bias as the Good-Turing or aPML estimator but with a lower mmMSE. For large values of M and small values of N, the mmMSE of the aPML and the Good-Turing estimators coincides. These regions can also be deduced from a comparison between the biased mmCCRBs associated with these estimators.

In Figs. 5a and 5b, we compare the missing-mass bias and the mmMSE of the Laplace estimator and the estimators that





Fig. 1: Example 1 (uniform pmf): The performance of the CML, Good-Turing, and Laplace estimators versus the number of elements, M, in terms of missing-mass bias (a) and the mmMSE (b). In (b) we also present the CCRB and the biased and unbiased versions of the mmCCRB.

Fig. 2: Example 1 (uniform pmf): The performance of the CML, Good-Turing, and Laplace estimators versus the number of samples, N, in terms of missing-mass bias (a) and the mmMSE (b). In (b) we also present the CCRB and the biased and unbiased versions of the mmCCRB.

are obtained after 1-5 iterations of the proposed missing-mass Fisher-scoring method in Algorithm 1 with initialization by the Laplace estimator. We set $\psi^{(k)} = \frac{1}{N}$, $\forall k \ge 1$ in Algorithm 1. It can be seen that the proposed missing-mass Fisher-scoring method reduces the missing-mass bias and the mmMSE of the Laplace estimator. In addition, the proposed method is consistent in the sense that by using more iterations, we obtain better estimators. These results demonstrate that the proposed mmCCRB can be used to further improve the performance of existing estimators of the missing mass.

VI. CONCLUSION

In this paper, we consider the problem of estimation of the missing mass. Similar to the naive CML estimator, which overestimates the probability of the observed elements, the CCRB on the MSE of χ -unbiased estimators does not provide a relevant bound for the missing-mass estimation problem. Hence, we adopt a new non-Bayesian approach, which is based on using the mmMSE risk function that only penalizes the estimation errors of elements that belong to the missing mass. The missing-mass unbiasedness, which is based on Lehmann's concept of unbiasedness and the mmMSE risk function, is proposed. We develop a new Cramér-Rao-type bound for this problem, the mmCCRB, which is a lower bound on the mmMSE of any locally missing-mass unbiased estimators. In addition, the biased mmCCRB on the mmMSE of for missing-mass biased estimators is developed. By using the mmCCRB on the mmMSE of the CML estimator, we show analytically that the CML estimator has the smallest mmMSE among all estimators that have the same missing-mass bias as the CML estimator. Based on the equality condition of the new mmCCRB, we derive a new method to improve existing estimators by an iterative missing-mass Fisher-scoring method.

In the simulations, we show that the unbiased mmCCRB is not a tight bound, but it can predict the behavior of the estimators w.r.t. the different system parameters. The biased versions of the mmCCRB are tight for the CML and Laplace estimators. Thus, the CML and the Laplace estimators achieve the lowest mmMSE possible for their bias. In contrast, for the Good-Turing and aPML estimators, there are regions in which one can theoretically find a better estimator (in the mmMSE sense) with the same bias. In addition, the different biased mmCCRBs can also be useful for setting the order relation between the different estimators and for exploring the bias-variance tradeoff. It is also shown that the proposed missing-mass Fisher-scoring method reduces the missing-mass bias and the mmMSE of the Laplace estimator. Future work should include missing-mass estimation with an unknown and infinite alphabet size. In addition, further investigation is needed regarding development of lower bounds based on the profile likelihood that may be tighter than the mmCCRB.





Fig. 3: Example 2 (Zipf distribution): The performance of the CML, Good-Turing, Laplace, and aPML estimators versus the number of elements, M, for N = 100 and s = 1 in terms of missing-mass bias (a) and the mmMSE (b). In (b) we also present the CCRB and the biased and unbiased versions of the mmCCRB.

APPENDIX A Derivation of (20)

In this appendix, we develop the CCRB on the MSE for pmf estimation. First, we note that taking the logarithm of (5) yields the following log-likelihood function:

$$\log p(\mathbf{x}; \boldsymbol{\theta}) = \sum_{m=1}^{M} C_{N,m}(\mathbf{x}) \log \theta_{m}, \ \mathbf{x} \in \mathcal{S}^{N}.$$
(66)

By substituting the derivative of (66) w.r.t. θ in (18), we obtain that the (l, m) element of the FIM is given by

$$\begin{bmatrix} \mathbf{J}(\boldsymbol{\theta}) \end{bmatrix}_{l,m} = \frac{1}{\theta_l \theta_m} \mathbf{E}_{\boldsymbol{\theta}} \begin{bmatrix} C_{N,l}(\mathbf{x}) C_{N,m}(\mathbf{x}) \end{bmatrix} \\ = \begin{cases} N(N-1) & m \neq l \\ N^2 + N \frac{1-\theta_m}{\theta_m} & m = l \end{cases},$$
(67)

 $\forall m, l = 1, \ldots, M$, where $C_{N,m}(\mathbf{x})$, $m = 1, \ldots, M$, are defined in (6) and the last equality holds by using known results on the moments of the multinomial distributed variables [64]. By using the elements in (67), the FIM for the probability estimation model can be written in a matrix form as

$$\mathbf{J}(\boldsymbol{\theta}) = N(N-1)\mathbf{1}_M\mathbf{1}_M^T + N\left(\operatorname{diag}(\boldsymbol{\theta})\right)^{-1}.$$
 (68)

It should be noted that $\mathbf{J}(\boldsymbol{\theta})$ from (68) is a well-defined, non-singular matrix, since we assume that $\theta_m \neq 0, \forall m =$

Fig. 4: Example 2 (Zipf distribution): The performance of the CML, Good-Turing, Laplace, and aPML estimators versus the number of samples, N, for M = 15 and s = 1 in terms of missing-mass bias (a) and the mmMSE (b). In (b) we also present the CCRB and the biased and unbiased versions of the mmCCRB.

1,..., *M*. By substituting (68) in (17), one obtains the following closed-form CCRB on the MSE under the constraint $\theta \in \Omega_{\theta}$:

$$\mathbf{E}_{\tilde{\boldsymbol{\theta}}} \left[(\hat{\boldsymbol{\theta}} - \tilde{\boldsymbol{\theta}}) (\hat{\boldsymbol{\theta}} - \tilde{\boldsymbol{\theta}})^T \right] \\
\succeq \frac{1}{N(N-1)} \mathbf{U} (\mathbf{U}^T \mathbf{1}_M \mathbf{1}_M^T \mathbf{U})^{-1} \mathbf{U}^T \\
+ \frac{1}{N} \mathbf{U} \left(\mathbf{U}^T \left(\operatorname{diag}(\tilde{\boldsymbol{\theta}}) \right)^{-1} \mathbf{U} \right)^{-1} \mathbf{U}^T \\
= \frac{1}{N} \mathbf{U} \left(\mathbf{U}^T \left(\operatorname{diag}(\tilde{\boldsymbol{\theta}}) \right)^{-1} \mathbf{U} \right)^{-1} \mathbf{U}^T, \quad (69)$$

where the last equality is obtained by substituting (4), which implies $\mathbf{1}_{M}^{T}\mathbf{U} = \mathbf{0}_{M-1}$. By applying the trace operator on the CCRB from (69) and using $\mathbf{U}^{T}\mathbf{U} = \mathbf{I}$ from (4) we obtain the bound on the trace MSE in (19)-(20).

APPENDIX B Proof of Lemma 1

In this appendix, we develop the missing-mass Lehmann unbiasedness. By substituting the missing-mass squared-error cost function from (22) and Ω_{θ} from (1) in (24), one obtains that the Lehmann-unbiasedness condition for the missing-mass



Fig. 5: Example 2 (Zipf distribution): The performance of the Laplace estimator and its improvement by the missing-mass Fisher-scoring method (from Algorithm 1) after 1-5 iterations versus the number of samples, N, in terms of missing-mass bias (a) and the mmMSE (b).

estimation problem is given by

$$\sum_{m=1}^{M} \mathbf{E}_{\boldsymbol{\theta}} \left[(\hat{\theta}_m - \eta_m)^2 \mathbb{1}_{\{m \in G_{N,0}(\mathbf{x})\}} \right]$$

$$\geq \sum_{m=1}^{M} \mathbf{E}_{\boldsymbol{\theta}} \left[(\hat{\theta}_m - \theta_m)^2 \mathbb{1}_{\{m \in G_{N,0}(\mathbf{x})\}} \right], \quad (70)$$

 $\forall \theta, \eta \in \Omega_{\theta}$. By using the definition of the constrained set in (1) and since U from (4) is the null-space matrix of this constrained set, then, for a given $\theta \in \Omega_{\theta}$, any $\eta \in \Omega_{\theta}$ can be written as (see, e.g. Section 4.2.4 in [65])

$$\eta = \theta + \mathbf{U}\mathbf{w},\tag{71}$$

where $\mathbf{w} \in \mathbb{R}^{M-1}$ is an arbitrary vector. By substituting (71) in (70), we obtain

$$\sum_{m=1}^{M} \mathbf{E}_{\boldsymbol{\theta}} \left[(\hat{\theta}_{m} - \theta_{m} - \mathbf{e}_{m}^{T} \mathbf{U} \mathbf{w}))^{2} \mathbb{1}_{\{m \in G_{N,0}(\mathbf{x})\}} \right]$$
$$\geq \sum_{m=1}^{M} \mathbf{E}_{\boldsymbol{\theta}} \left[(\hat{\theta}_{m} - \theta_{m})^{2} \mathbb{1}_{\{m \in G_{N,0}(\mathbf{x})\}} \right], \quad (72)$$

 $\forall \boldsymbol{\theta} \in \Omega_{\boldsymbol{\theta}}, \mathbf{w} \in \mathbb{R}^{M-1}$. By using (10), the unbiasedness condition from (72) can be rewritten as:

$$\sum_{m=1}^{M} (\mathbf{e}_{m}^{T} \mathbf{U} \mathbf{w})^{2} \operatorname{Pr}(\mathbf{x} \in \mathcal{A}_{m}; \boldsymbol{\theta})$$

$$\geq 2 \sum_{m=1}^{M} \operatorname{E}_{\boldsymbol{\theta}} \left[(\hat{\theta}_{m} - \mathbf{e}_{m}^{T} \boldsymbol{\theta}) \mathbb{1}_{\{m \in G_{N,0}(\mathbf{x})\}} \right] \mathbf{e}_{m}^{T} \mathbf{U} \mathbf{w}, \quad (73)$$

 $\forall \boldsymbol{\theta} \in \Omega_{\boldsymbol{\theta}}, \mathbf{w} \in \mathbb{R}^{M-1}$. Since the condition in (73) should be satisfied for any $\mathbf{w} \in \mathbb{R}^{M-1}$, it should be satisfied in particular for both $\mathbf{w} = \epsilon \mathbf{e}_k$ and $\mathbf{w} = -\epsilon \mathbf{e}_k$, where $\epsilon > 0$. By summing the separate substitution of $\mathbf{w} = \pm \epsilon \mathbf{e}_k$ (that is, the result of substituting $\mathbf{w} = \epsilon \mathbf{e}_k$ into (73) and the result of substituting

 $\mathbf{w} = \epsilon \mathbf{e}_k$ into the same equation), we obtain the following *necessary* condition for (73) to hold:

$$\sum_{m=1}^{M} \mathbf{E}_{\boldsymbol{\theta}} \left[(\hat{\theta}_m - \theta_m) \mathbb{1}_{\{m \in G_{N,0}(\mathbf{x})\}} \right] \mathbf{e}_m^T \mathbf{U} = \mathbf{0}_{M-1}^T, \quad (74)$$

 $\forall \theta \in \Omega_{\theta}$. Since the l.h.s. of (73) is a quadratic term, it can be verified that (74) is also a *sufficient* condition for unbiasedness in this case. Therefore, by applying the transpose operator on (74), one obtains that the missing-mass unbiasedness in (26) is the Lehmann unbiasedness under the missing-mass squared error cost function.

Appendix C

PROOF OF THEOREM 1

In this appendix, we develop the new mmCCRB from Theorem 1. The proof is divided into: 1) the development of Lemma 3 in Subsection C-A; 2) the main development of the bound based on the covariance inequality in Subsection C-B; and 3) derivation of the equality condition in Subsection C-C. To this end, we define $\Gamma(\mathbf{x}, \theta)$ as a diagonal $M \times M$ matrix with the following elements on its diagonal:

$$[\mathbf{\Gamma}(\mathbf{x},\boldsymbol{\theta})]_{m,m} \stackrel{\Delta}{=} \epsilon_m(\boldsymbol{\theta}) \mathbb{1}_{\{m \in G_{N,0}(\mathbf{x})\}}, \ m = 1..., M, \quad (75)$$

where

$$\epsilon_m(\boldsymbol{\theta}) \stackrel{\triangle}{=} \hat{\theta}_m - \mathbf{E}_{\boldsymbol{\theta}} \left[\hat{\theta}_m | \mathbf{x} \in \mathcal{A}_m \right], \ m = 1, \dots, M.$$
(76)

A. Lemma 3

In this subsection, we prove the following Lemma: *Lemma 3:*

$$\mathbf{E}_{\boldsymbol{\theta}}\left[\boldsymbol{\Gamma}(\mathbf{x},\boldsymbol{\theta})\boldsymbol{\Delta}^{T}(\mathbf{x},\boldsymbol{\theta})\right] = \mathbf{S}(\boldsymbol{\theta}), \tag{77}$$

where $\Delta(\mathbf{x}, \theta)$, $\mathbf{S}(\theta)$, and $\Gamma(\mathbf{x}, \theta)$ are defined in (31), (32), and (75), respectively.

Proof: By substituting (31), (75), and (76) in (77) one obtains that the *m*th row of $E_{\theta} [\Gamma(\mathbf{x}, \theta) \Delta^T(\mathbf{x}, \theta)]$ on the r.h.s. of (77) satisfies

$$E_{\boldsymbol{\theta}} \left[\epsilon_{m}(\boldsymbol{\theta}) \boldsymbol{\Delta}_{1:M,m}^{T}(\mathbf{x},\boldsymbol{\theta}) \right] \\= E_{\boldsymbol{\theta}} \left[\epsilon_{m}(\boldsymbol{\theta}) \nabla_{\boldsymbol{\theta}} \log p(\mathbf{x} | \mathbf{x} \in \mathcal{A}_{m}; \boldsymbol{\theta}) \mathbb{1}_{\{\mathbf{x} \in \mathcal{A}_{m}\}} \right] \\= \Pr(\mathbf{x} \in \mathcal{A}_{m}; \boldsymbol{\theta}) \sum_{\boldsymbol{\alpha} \in \mathcal{A}_{m}} \epsilon_{m}(\boldsymbol{\theta}) \nabla_{\boldsymbol{\theta}}^{T} \Pr(\mathbf{x} = \boldsymbol{\alpha} | \mathbf{x} \in \mathcal{A}_{m}; \boldsymbol{\theta}), (78)$$

 $m = 1, \ldots, M$, where the last equality stems from

$$\nabla_{\boldsymbol{\theta}} \log p(\mathbf{x} | \mathbf{x} \in \mathcal{A}_m; \boldsymbol{\theta}) = \frac{\nabla_{\boldsymbol{\theta}} p(\mathbf{x} | \mathbf{x} \in \mathcal{A}_m; \boldsymbol{\theta})}{p(\mathbf{x} | \mathbf{x} \in \mathcal{A}_m; \boldsymbol{\theta})}, \ \forall \mathbf{x} \in \mathcal{A}_m,$$

in which A_m is defined in (9). Then, by applying the product rule on the r.h.s. of (78), we obtain

$$\sum_{\boldsymbol{\alpha}\in\mathcal{A}_{m}} \epsilon_{m}(\boldsymbol{\theta}) \nabla_{\boldsymbol{\theta}}^{T} \operatorname{Pr}(\mathbf{x}=\boldsymbol{\alpha}|\mathbf{x}\in\mathcal{A}_{m};\boldsymbol{\theta})$$

$$= \nabla_{\boldsymbol{\theta}}^{T} \left\{ \sum_{\boldsymbol{\alpha}\in\mathcal{A}_{m}} \epsilon_{m}(\boldsymbol{\theta}) \operatorname{Pr}(\mathbf{x}=\boldsymbol{\alpha}|\mathbf{x}\in\mathcal{A}_{m};\boldsymbol{\theta}) \right\}$$

$$- \sum_{\boldsymbol{\alpha}\in\mathcal{A}_{m}} \nabla_{\boldsymbol{\theta}}^{T} \{\epsilon_{m}(\boldsymbol{\theta})\} \operatorname{Pr}(\mathbf{x}=\boldsymbol{\alpha}|\mathbf{x}\in\mathcal{A}_{m}(\mathbf{x});\boldsymbol{\theta})$$

$$= \nabla_{\boldsymbol{\theta}}^{T} \{ \operatorname{E}_{\boldsymbol{\theta}} [\epsilon_{m}(\boldsymbol{\theta})|\mathbf{x}\in\mathcal{A}_{m}] \}$$

$$- \sum_{\boldsymbol{\alpha}\in\mathcal{A}_{m}} \nabla_{\boldsymbol{\theta}}^{T} \{\epsilon_{m}(\boldsymbol{\theta})\} \operatorname{Pr}(\mathbf{x}=\boldsymbol{\alpha}|\mathbf{x}\in\mathcal{A}_{m}(\mathbf{x});\boldsymbol{\theta}). \quad (79)$$

Computing the conditional expectation of (76), given the event that $m \in G_{N,0}(\mathbf{x})$, results in

$$E_{\boldsymbol{\theta}} \left[\epsilon_m(\boldsymbol{\theta}) | \mathbf{x} \in \mathcal{A}_m \right] = 0, \ m = 1, \dots, M.$$
(80)

In addition, computing the gradient of (76) results in

$$\nabla_{\boldsymbol{\theta}}^{T} \epsilon_{m}(\boldsymbol{\theta}) = -\nabla_{\boldsymbol{\theta}}^{T} \mathbf{E}_{\boldsymbol{\theta}} \left[\hat{\theta}_{m} | \mathbf{x} \in \mathcal{A}_{m} \right], \qquad (81)$$

m = 1, ..., M. By using the conditional expectation definition, and then substituting (80) and (81) in (79), one obtains

$$\sum_{\boldsymbol{\alpha}\in\mathcal{A}_{m}} \epsilon_{m}(\boldsymbol{\theta}) \nabla_{\boldsymbol{\theta}}^{T} \operatorname{Pr}(\mathbf{x}=\boldsymbol{\alpha}|\mathbf{x}\in\mathcal{A}_{m}(\mathbf{x});\boldsymbol{\theta})$$
$$= \nabla_{\boldsymbol{\theta}}^{T} \left\{ \operatorname{E}_{\boldsymbol{\theta}} \left[\hat{\theta}_{m} | \mathbf{x}\in\mathcal{A}_{m} \right] \right\} \sum_{\boldsymbol{\alpha}\in\mathcal{A}_{m}} \operatorname{Pr}(\mathbf{x}=\boldsymbol{\alpha}|\mathbf{x}\in\mathcal{A}_{m};\boldsymbol{\theta})$$
$$= \nabla_{\boldsymbol{\theta}}^{T} \left\{ \operatorname{E}_{\boldsymbol{\theta}} \left[\hat{\theta}_{m} | \mathbf{x}\in\mathcal{A}_{m} \right] \right\}, \tag{82}$$

where the last equality stems from the fact that for a conditional pmf we have $\sum_{\boldsymbol{\alpha}\in\mathcal{A}_m} \Pr(\mathbf{x}=\boldsymbol{\alpha}|m\in G_{N,0}(\mathbf{x});\boldsymbol{\theta})=1$. By substituting the definition of the missing-mass bias vector, $\mathbf{b}_{N,0}(\boldsymbol{\theta})$, from (25) in (82) and using the fact that $\nabla_{\boldsymbol{\theta}}^T \boldsymbol{\theta}_m = \mathbf{e}_m^T$, one obtains

$$\sum_{\boldsymbol{\alpha}\in\mathcal{A}_{m}} \epsilon_{m}(\boldsymbol{\theta}) \nabla_{\boldsymbol{\theta}}^{T} \operatorname{Pr}(\mathbf{x}=\boldsymbol{\alpha}|\mathbf{x}\in\mathcal{A}_{m}(\mathbf{x});\boldsymbol{\theta})$$
$$= \nabla_{\boldsymbol{\theta}}^{T} \left\{ \frac{\left[\mathbf{b}_{N,0}(\boldsymbol{\theta})\right]_{m}}{\operatorname{Pr}(\mathbf{x}\in\mathcal{A}_{m};\boldsymbol{\theta})} \right\} + \mathbf{e}_{m}^{T}, \quad (83)$$

m = 1, ..., M. Substitution of (32) in (83) and then substituting the result in (78) results in (77).

B. Covariance inequality

The following part of the proof is along the path of the proof from [11] for the CCRB on the MSE in a conventional estimation problem. Let $\mathbf{W} \in \mathbb{R}^{M \times M}$ be an arbitrary matrix and $\tilde{\boldsymbol{\theta}} \in \Omega_{\boldsymbol{\theta}}$ is a local parameter vector. Then,

$$\mathbf{0} \leq \mathbf{E}_{\tilde{\boldsymbol{\theta}}} \left[\left(\boldsymbol{\Gamma}(\mathbf{x}, \tilde{\boldsymbol{\theta}}) - \mathbf{W}^{T} \mathbf{U} \mathbf{U}^{T} \boldsymbol{\Delta}(\mathbf{x}, \tilde{\boldsymbol{\theta}}) \right) \\ \times \left(\boldsymbol{\Gamma}(\mathbf{x}, \tilde{\boldsymbol{\theta}}) - \mathbf{W}^{T} \mathbf{U} \mathbf{U}^{T} \boldsymbol{\Delta}(\mathbf{x}, \tilde{\boldsymbol{\theta}}) \right)^{T} \right] \\ = \mathbf{E}_{\tilde{\boldsymbol{\theta}}} \left[\boldsymbol{\Gamma}(\mathbf{x}, \tilde{\boldsymbol{\theta}}) \boldsymbol{\Gamma}^{T}(\mathbf{x}, \tilde{\boldsymbol{\theta}}) \right] - \mathbf{S}^{T}(\tilde{\boldsymbol{\theta}}) \mathbf{U} \mathbf{U}^{T} \mathbf{W} \\ - \mathbf{W}^{T} \mathbf{U} \mathbf{U}^{T} \mathbf{S}(\tilde{\boldsymbol{\theta}}) + \mathbf{W}^{T} \mathbf{U} \mathbf{U}^{T} \mathbf{J}^{(0)}(\tilde{\boldsymbol{\theta}}) \mathbf{U} \mathbf{U}^{T} \mathbf{W}, \quad (84)$$

where we substitute (77) from Lemma 3 and the mmFIM definition from (30). By rearranging (84), we obtain

$$\mathbf{0} \leq \mathbf{E}_{\tilde{\boldsymbol{\theta}}} \left[\boldsymbol{\Gamma}(\mathbf{x}, \tilde{\boldsymbol{\theta}}) \boldsymbol{\Gamma}^{T}(\mathbf{x}, \tilde{\boldsymbol{\theta}}) \right] - \mathbf{S}^{T}(\tilde{\boldsymbol{\theta}}) \mathbf{U} \mathbf{U}^{T} \mathbf{W} - \mathbf{W}^{T} \mathbf{U} \mathbf{U}^{T} \mathbf{S}(\tilde{\boldsymbol{\theta}}) + \mathbf{W}^{T} \mathbf{U} \mathbf{U}^{T} \mathbf{J}^{(0)}(\tilde{\boldsymbol{\theta}}) \mathbf{U} \mathbf{U}^{T} \mathbf{W}.$$
(85)

By applying the trace operator on (85), it can be verified that the matrix inequality in (85) provides a family of bounds on the mmMSE from (23), which depends on the specific choice of the matrix **W**. Theorem 1 is obtained by choosing the optimal member from this family, as described in the following.

Since $\mathbf{U}^T \mathbf{J}^{(0)}(\tilde{\boldsymbol{\theta}}) \mathbf{U}$ is a non-singular matrix under regularity Condition C.1, then it is shown in [11] that the greatest lower bound, i.e. the supremum of the r.h.s. of (85) over W, is obtained by a matrix W which satisfies

$$\mathbf{W}^{T}\mathbf{U} = \mathbf{S}^{T}(\tilde{\boldsymbol{\theta}})\mathbf{U}\left(\mathbf{U}^{T}\mathbf{J}^{(0)}(\tilde{\boldsymbol{\theta}})\mathbf{U}\right)^{-1}.$$
 (86)

By substituting (86) into (85), one obtains

$$\mathbf{E}_{\tilde{\boldsymbol{\theta}}} \left[\boldsymbol{\Gamma}(\mathbf{x}, \tilde{\boldsymbol{\theta}}) \boldsymbol{\Gamma}^{T}(\mathbf{x}, \tilde{\boldsymbol{\theta}}) \right] \\ \succeq \mathbf{S}^{T}(\tilde{\boldsymbol{\theta}}) \mathbf{U} (\mathbf{U}^{T} \mathbf{J}^{(0)}(\tilde{\boldsymbol{\theta}}) \mathbf{U})^{-1} \mathbf{U}^{T} \mathbf{S}(\tilde{\boldsymbol{\theta}}).$$
(87)

By applying the trace operator on (87) and substituting (75), we obtain

$$\sum_{m=1}^{M} \mathbf{E}_{\tilde{\boldsymbol{\theta}}} \left[\epsilon_{m}(\tilde{\boldsymbol{\theta}})^{2} \mathbb{1}_{\{m \in G_{N,0}(\mathbf{x})\}} \right]$$

$$\geq \operatorname{trace}(\mathbf{S}^{T}(\tilde{\boldsymbol{\theta}}) \mathbf{U}(\mathbf{U}^{T} \mathbf{J}^{(0)}(\tilde{\boldsymbol{\theta}}) \mathbf{U})^{-1} \mathbf{U}^{T} \mathbf{S}(\tilde{\boldsymbol{\theta}})). \quad (88)$$

The following equation relates the mmMSE from (23) and the l.h.s. of (88). From (76), it can be seen that

where the second equality stems from (80) and the last equality is obtained by substituting (25). By multiplying (89) by $Pr(\mathbf{x} \in \mathcal{A}_m; \tilde{\boldsymbol{\theta}})$, then summing the result over $m = 1, \ldots, M$, and substituting the result in (88), one obtains the mmBCRB on the mmMSE of biased estimator evaluated at the local point, $\tilde{\boldsymbol{\theta}}$, in (33)-(34).

C. Derivation of the equality condition in (35)

According to Cauchy-Schwartz inequality properties (or covariance inequality properties), equality in (84) for W that satisfies (86) holds if

$$\begin{split} & \mathbf{E}_{\tilde{\boldsymbol{\theta}}} \Big[(\boldsymbol{\Gamma}(\mathbf{x}, \tilde{\boldsymbol{\theta}}) - \mathbf{S}^{T}(\tilde{\boldsymbol{\theta}}) \mathbf{U} (\mathbf{U}^{T} \mathbf{J}^{(0)}(\tilde{\boldsymbol{\theta}}) \mathbf{U})^{-1} \mathbf{U}^{T} \boldsymbol{\Delta}(\mathbf{x}, \tilde{\boldsymbol{\theta}})) \\ & \times (\boldsymbol{\Gamma}(\mathbf{x}, \tilde{\boldsymbol{\theta}}) - \mathbf{S}^{T}(\tilde{\boldsymbol{\theta}}) \mathbf{U} (\mathbf{U}^{T} \mathbf{J}^{(0)}(\tilde{\boldsymbol{\theta}}) \mathbf{U})^{-1} \mathbf{U}^{T} \boldsymbol{\Delta}(\mathbf{x}, \tilde{\boldsymbol{\theta}}))^{T} \Big] \\ &= \mathbf{0}. \end{split}$$
(90)

The condition in (90) holds if

$$\Gamma(\mathbf{x}, \tilde{\boldsymbol{\theta}}) = \mathbf{S}^T(\tilde{\boldsymbol{\theta}}) \mathbf{U} (\mathbf{U}^T \mathbf{J}^{(0)}(\tilde{\boldsymbol{\theta}}) \mathbf{U})^{-1} \mathbf{U}^T \boldsymbol{\Delta}(\mathbf{x}, \tilde{\boldsymbol{\theta}})$$
(91)

which, by using (75), implies that

$$\epsilon_m(\tilde{\boldsymbol{\theta}}) = \left[\mathbf{S}(\tilde{\boldsymbol{\theta}})^T \mathbf{U} (\mathbf{U}^T \mathbf{J}^{(0)}(\tilde{\boldsymbol{\theta}}) \mathbf{U})^{-1} \mathbf{U}^T \boldsymbol{\Delta}(\mathbf{x}, \tilde{\boldsymbol{\theta}}) \right]_{m,m},$$
(92)

for any m = 1, ..., M such that $m \in G_{N,0}\mathbf{x}$. By substituting (25) and (76) in (92) we get (35).

APPENDIX D Derivation of the FIM in (39)

In this appendix, we develop the closed-form mmFIM, defined in (30), for the observation model from (5). By substituting (12) in (31), one obtains

$$\begin{aligned} \mathbf{\Delta}_{1:M,m}(\mathbf{x}, \boldsymbol{\theta}) &= \\ \nabla_{\boldsymbol{\theta}} \left(\sum_{l=1}^{M} C_{N,l}(\mathbf{x}) \log \theta_l - N \log(1 - \theta_m) \right) \mathbb{1}_{\{m \in G_{N,0}(\mathbf{x})\}} \\ &= \left(\mathbf{v}^T(\mathbf{x}, \boldsymbol{\theta}) - \frac{N}{1 - \theta_m} \mathbf{e}_m^T \right) \mathbb{1}_{\{m \in G_{N,0}(\mathbf{x})\}}, \end{aligned}$$
(93)

where $\mathbf{v}(\mathbf{x}, \boldsymbol{\theta})$ is defined in (42). By substituting (93) in the mmFIM definition from (30), we obtain that the (k, l)th element of the mmFIM for our model is given by

$$\begin{bmatrix} \mathbf{J}^{(0)}(\boldsymbol{\theta}) \end{bmatrix}_{k,l} = \sum_{m=1}^{M} \mathbf{E}_{\boldsymbol{\theta}} \left[\left[\mathbf{\Delta}(\mathbf{x}, \boldsymbol{\theta}) \right]_{k,m} \left[\mathbf{\Delta}(\mathbf{x}, \boldsymbol{\theta}) \right]_{l,m} \right] \\ = \sum_{m=1}^{M} \mathbf{E}_{\boldsymbol{\theta}} \left[\left(\frac{C_{N,k}(\mathbf{x})}{\theta_{k}} + \frac{N}{1 - \theta_{m}} \delta_{k,m} \right) \\ \times \left(\frac{C_{N,l}(\mathbf{x})}{\theta_{l}} + \frac{N}{1 - \theta_{m}} \delta_{l,m} \right) \mathbb{1}_{\{m \in G_{N,0}(\mathbf{x})\}} \right] \\ = \sum_{m=1}^{M} \left\{ \frac{1}{\theta_{k} \theta_{l}} \mathbf{E}_{\boldsymbol{\theta}} \left[C_{N,k}(\mathbf{x}) C_{N,l}(\mathbf{x}) | m \in G_{N,0}(\mathbf{x}) \right] \\ + \frac{N \delta_{k,m}}{(1 - \theta_{m}) \theta_{l}} \mathbf{E}_{\boldsymbol{\theta}} \left[C_{N,k}(\mathbf{x}) | \mathbf{x} \in \mathcal{A}_{m} \right] \\ + \frac{N \delta_{l,m}}{(1 - \theta_{m}) \theta_{k}} \mathbf{E}_{\boldsymbol{\theta}} \left[C_{N,k}(\mathbf{x}) | \mathbf{x} \in \mathcal{A}_{m} \right] \\ + \frac{N^{2} \delta_{k,m} \delta_{l,m}}{(1 - \theta_{m})^{2}} \right\} \Pr(\mathbf{x} \in \mathcal{A}_{m}; \boldsymbol{\theta}), \tag{94}$$

where the last equality is obtained by using the law of total probability. In the following, we compute the conditional expectation terms from (94). First, it can be seen that

$$\sum_{\boldsymbol{\alpha}\in\mathcal{A}_m} \Pr(\mathbf{x}=\boldsymbol{\alpha};\boldsymbol{\theta}) = \Pr(\mathbf{x}\in\mathcal{A}_m;\boldsymbol{\theta}) = \left(\sum_{n=1,n\neq m}^N \theta_n\right)^N, (95)$$

 $m = 1, \ldots, M$, where A_m is defined in (9) and where the probability of the *m*th element to be unobserved on the r.h.s. of (95) is an alternative way of writing the r.h.s of (10). Then, by using (12), it can be verified that

$$E_{\boldsymbol{\theta}} [C_{N,k}(\mathbf{x}) | \mathbf{x} \in \mathcal{A}_{m}] = \sum_{\boldsymbol{\alpha} \in \mathcal{A}_{m}} C_{N,k}(\boldsymbol{\alpha}) \frac{\prod_{n=1}^{M} \theta_{n}^{C_{N,n}(\boldsymbol{\alpha})}}{(1-\theta_{m})^{N}} = \frac{\theta_{k}}{(1-\theta_{m})^{N}} \sum_{\boldsymbol{\alpha} \in \mathcal{A}_{m}} \frac{\partial}{\partial \theta_{k}} \prod_{n=1}^{M} \theta_{n}^{C_{N,n}(\boldsymbol{\alpha})} = \frac{\theta_{k}}{(1-\theta_{m})^{N}} \frac{\partial}{\partial \theta_{k}} \sum_{\boldsymbol{\alpha} \in \mathcal{A}_{m}} \Pr(\mathbf{x} = \boldsymbol{\alpha}; \boldsymbol{\theta}), \quad (96)$$

m, k = 1, ..., M. It should be noted that the last equality in (96) is only valid if we use the probability term, $\Pr(\mathbf{x} = \boldsymbol{\alpha}; \boldsymbol{\theta})$,

as it is written in (95). By substituting (95) in (96), one obtains

$$E_{\boldsymbol{\theta}} [C_{N,k}(\mathbf{x}) | \mathbf{x} \in \mathcal{A}_{m}] = \frac{\theta_{k}}{(1 - \theta_{m})^{N}} \frac{\partial}{\partial \theta_{k}} \Big(\sum_{n=1, n \neq m}^{N} \theta_{n} \Big)^{N} = \begin{cases} \frac{N \theta_{k}}{1 - \theta_{m}} & m \neq k \\ 0 & m = k \end{cases}.$$
(97)

Similar to the derivation of (96), by using (12), we obtain

$$E_{\boldsymbol{\theta}} [C_{N,k}(\mathbf{x})C_{N,l}(\mathbf{x})|\mathbf{x} \in \mathcal{A}_{m}]$$

$$= \sum_{\boldsymbol{\alpha} \in \mathcal{A}_{m}} C_{N,k}(\boldsymbol{\alpha})C_{N,l}(\boldsymbol{\alpha})\frac{\prod_{n=1}^{M}\theta_{n}^{C_{N,n}(\boldsymbol{\alpha})}}{(1-\theta_{m})^{N}}$$

$$= \sum_{\boldsymbol{\alpha} \in \mathcal{A}_{m}} \frac{\theta_{k}\theta_{l}}{(1-\theta_{m})^{N}}\frac{\partial^{2}}{\partial\theta_{k}\partial\theta_{l}}\prod_{n=1}^{M}\theta_{n}^{C_{N,n}(\boldsymbol{\alpha})}$$

$$+ \frac{\delta_{k,l}}{(1-\theta_{m})^{N}}\sum_{\boldsymbol{\alpha} \in \mathcal{A}_{m}}C_{N,k}(\boldsymbol{\alpha})\prod_{n=1}^{M}\theta_{n}^{C_{N,n}(\boldsymbol{\alpha})}$$

$$= \frac{\theta_{k}\theta_{l}}{(1-\theta_{m})^{N}}\frac{\partial^{2}}{\partial\theta_{k}\partial\theta_{l}}\sum_{\boldsymbol{\alpha} \in \mathcal{A}_{m}}\Pr(\mathbf{x}=\boldsymbol{\alpha};\boldsymbol{\theta})$$

$$+ \frac{N\delta_{k,l}}{1-\theta_{m}}E_{\boldsymbol{\theta}}[C_{N,k}(\mathbf{x})|\mathbf{x} \in \mathcal{A}_{m}], \qquad (98)$$

where we replace the order of the derivative and the sum, and we use (5). Again, it should be noted that the last equality in (98) is only valid if we use the probability term, $Pr(\mathbf{x} = \alpha; \theta)$, in (95). By substituting (95) and (97) in (98), one obtains

$$E_{\theta} [C_{N,k}(\mathbf{x})C_{N,l}(\mathbf{x})|\mathbf{x} \in \mathcal{A}_{m}] = \frac{\theta_{k}\theta_{l}}{(1-\theta_{m})^{N}} \frac{\partial^{2}}{\partial\theta_{k}\partial\theta_{l}} \Big(\sum_{n=1,n\neq m}^{N} \theta_{n}\Big)^{N} + \frac{N\theta_{k}\delta_{k,l}(1-\delta_{m,k})}{1-\theta_{m}}$$
$$= \begin{cases} \frac{N(N-1)\theta_{k}\theta_{l}}{(1-\theta_{m})^{2}} & \text{if } m\neq k, \ l\neq m, \ k\neq l\\ \frac{N(N-1)\theta_{k}^{2}}{(1-\theta_{m})^{2}} + \frac{N\theta_{k}}{1-\theta_{m}} & \text{if } k=l\neq m\\ 0 & \text{otherwise} \end{cases}$$
(99)

Thus, by substituting (10), (97), and (99) in (94), one obtains that the elements of the mmFIM are given by

$$\begin{aligned} [\mathbf{J}^{(0)}(\boldsymbol{\theta})]_{k,k} \\ &\sum_{m=1,m\neq k}^{M} \left(\frac{N(N-1)}{(1-\theta_m)^2} + \frac{N}{\theta_k(1-\theta_m)} \right) (1-\theta_m)^N \\ &+ \frac{N^2}{(1-\theta_k)^2} (1-\theta_k)^N \\ &= \sum_{m=1}^{M} \left(\frac{N(N-1)}{(1-\theta_m)^2} + \frac{N}{\theta_k(1-\theta_m)} \right) (1-\theta_m)^N \\ &+ \frac{N}{(1-\theta_k)^2} (1-\theta_k)^N \\ &- \frac{N}{\theta_k(1-\theta_k)} (1-\theta_k)^N, \end{aligned}$$
(100)

for $k = 1, \ldots, M$, and

$$[\mathbf{J}^{(0)}(\boldsymbol{\theta})]_{k,l} = \sum_{m=1,m\neq k,m\neq l}^{M} \frac{N(N-1)}{(1-\theta_m)^2} (1-\theta_m)^N + \frac{N^2}{(1-\theta_k)^2} (1-\theta_k)^N + \frac{N^2}{(1-\theta_l)^2} (1-\theta_l)^N = \sum_{m=1}^{M} \frac{N(N-1)}{(1-\theta_m)^2} (1-\theta_m)^N + \frac{N}{(1-\theta_k)^2} (1-\theta_k)^N + \frac{N}{(1-\theta_l)^2} (1-\theta_l)^N, (101)$$

for k, l = 1, ..., M, $k \neq l$. By rearranging the elements in (100) and (101), and substituting (10), we obtain the closed-form matrix $\mathbf{J}^{(0)}(\boldsymbol{\theta})$ in its matrix representation in (39).

$$D_{m,m} = -\frac{(1-\theta_m)^N}{(1-\theta_m)^2} + \sum_{l=1, l \neq m}^M \frac{(1-\theta_l)^N}{\theta_m (1-\theta_l)}.$$
 (102)

APPENDIX E Proof of Lemma 2

Taking the logarithm of (12) yields the following conditional log-likelihood function:

$$\log p(\mathbf{x}|\mathbf{x} \in \mathcal{A}_m; \boldsymbol{\theta}) = \sum_{l=1}^{M} C_{N,l}(\mathbf{x}) \log \theta_l - N \log(1 - \theta_m) (103)$$

for any $\mathbf{x} \in \mathcal{A}_m$. By substituting (76) and the derivative of (103) w.r.t. $\boldsymbol{\theta}$ in (78), we obtain

$$\begin{aligned} \mathbf{E}_{\boldsymbol{\theta}} \left[\epsilon_{m}(\boldsymbol{\theta}) \boldsymbol{\Delta}_{1:M,m}^{T}(\mathbf{x},\boldsymbol{\theta}) \right] &= \mathbf{E}_{\boldsymbol{\theta}} \left[\left(\hat{\theta}_{m} - \mathbf{E}_{\boldsymbol{\theta}} [\hat{\theta}_{m} | \mathbf{x} \in \mathcal{A}_{m}] \right) \right. \\ & \left. \times (\mathbf{v}^{T}(\mathbf{x},\boldsymbol{\theta}) + N \mathbf{e}_{m}^{T} \frac{1}{1 - \theta_{m}}) \mathbb{1}_{\{\mathbf{x} \in \mathcal{A}_{m}\}} \right] \\ &= \mathbf{E}_{\boldsymbol{\theta}} \left[\left(\hat{\theta}_{m} - \theta_{m} \right) (\mathbf{v}^{T}(\mathbf{x},\boldsymbol{\theta}) + N \mathbf{e}_{m}^{T} \frac{1}{1 - \theta_{m}}) \mathbb{1}_{\{\mathbf{x} \in \mathcal{A}_{m}\}} \right] \\ &- \left[\mathbf{b}_{N,0}(\boldsymbol{\theta}) \right]_{m} \left(\mathbf{E}_{\boldsymbol{\theta}} \left[\mathbf{v}^{T}(\mathbf{x},\boldsymbol{\theta}) | \mathbf{x} \in \mathcal{A}_{m} \right] + N \mathbf{e}_{m}^{T} \frac{1}{1 - \theta_{m}} \right) (104) \end{aligned}$$

where the last equality is obtained by substituting (25). By substituting (25) and (97) in (104), we obtain

$$E_{\boldsymbol{\theta}} \left[\epsilon_{m}(\boldsymbol{\theta}) \boldsymbol{\Delta}_{1:M,m}^{T}(\mathbf{x}, \boldsymbol{\theta}) \right] \\= E_{\boldsymbol{\theta}} \left[(\hat{\theta}_{m} - \theta_{m}) \mathbf{v}^{T}(\mathbf{x}, \boldsymbol{\theta}) \mathbb{1}_{\{\mathbf{x} \in \mathcal{A}_{m}\}} \right] \\+ \frac{N}{1 - \theta_{m}} [\mathbf{b}_{N,0}(\boldsymbol{\theta})]_{m} \mathbf{e}_{m}^{T} - N \frac{1}{1 - \theta_{m}} [\mathbf{b}_{N,0}(\boldsymbol{\theta})]_{m} \mathbf{1}_{M}^{T}, (105)$$

 $m = 1, \ldots, M$. By substituting (105) in (77), we obtain (41).

REFERENCES

- H. Robbins, "Estimating the total probability of the unobserved outcomes of an experiment," *The Annals of Mathematical Statistics*, vol. 39, no. 1, pp. 256–257, 1968.
- [2] I. J. Good, "The population frequencies of species and the estimation of population parameters," *Biometrika*, vol. 40, pp. 237–264, 1953.
- [3] B. Efron and R. Thisted, "Estimating the number of unseen species: How many words did Shakespeare know?" *Biometrika*, vol. 63, no. 3, pp. 435–447, 12 1976.
- [4] C. Budianu and L. Tong, "Good-Turing estimation of the number of operating sensors: a large deviations analysis," in *Proc. ICASSP*, vol. 2, May 2004, pp. 1029–1032.

- [5] C. Budianu, S. Ben-David, and L. Tong, "Estimation of the number of operating sensors in large-scale sensor networks with mobile access," *IEEE Trans. Signal Processing*, vol. 54, no. 5, pp. 1703–1715, May 2006.
- [6] S. F. Chen and J. Goodman, "An empirical study of smoothing techniques for language modeling," in *Annual Meeting on Association for Computational Linguistics*, 1996, pp. 310–318.
- [7] S. Katz, "Estimation of probabilities from sparse data for the language model component of a speech recognizer," *IEEE Trans. acoustics, speech, and signal processing*, vol. 35, no. 3, pp. 400–401, 1987.
- [8] A. Orlitsky, N. P. Santhanam, and J. Zhang, "Always Good Turing: Asymptotically optimal probability estimation," *Science*, vol. 302, no. 5644, pp. 427–431, 2003.
- [9] G. Valiant and P. Valiant, "Estimating the unseen: An n/log(n)-sample estimator for entropy and support size, shown optimal via new CLTs," in the 43rd ACM Symposium on Theory of Computing, 2011, pp. 685–694.
- [10] J. Gorman and A. Hero, "Lower bounds for parametric estimation with constraints," *IEEE Trans. Information Theory*, vol. 36, no. 6, pp. 1285– 1301, Nov. 1990.
- [11] P. Stoica and B. C. Ng, "On the Cramér-Rao bound under parametric constraints," *IEEE Signal Processing Letters*, vol. 5, no. 7, pp. 177–179, July 1998.
- [12] T. J. Moore, B. M. Sadler, and R. J. Kozick, "Maximum-likelihood estimation, the Cramér-Rao bound, and the method of scoring with parameter constraints," *IEEE Trans. Signal Processing*, vol. 56, no. 3, pp. 895–908, 2008.
- [13] E. Nitzan, T. Routtenberg, and J. Tabrikian, "Cramér-Rao bound for constrained parameter estimation using Lehmann-unbiasedness," *IEEE Trans. Signal Processing*, vol. 67, no. 3, pp. 753–768, Feb 2019.
- [14] E. Nitzan, T. Routtenberg, and J. Tabrikian, "Cramér-Rao bound under norm constraint," *IEEE Signal Processing Letters*, vol. 26, no. 9, pp. 1393–1397, 2019.
- [15] P. S. Laplace, Pierre-Simon Laplace Philosophical Essay on Probabilities: Translated from the fifth French edition of 1825 With Notes by the Translator. Springer Science & Business Media, 2012, vol. 13.
- [16] A. Nadas, "On Turing's formula for word probabilities," *IEEE Trans.* Acoustics, Speech, and Signal Processing, vol. 33, no. 6, pp. 1414–1416, Dec. 1985.
- [17] W. A. Gale and G. Sampson, "Good-Turing frequency estimation without tears," *Journal of quantitative linguistics*, vol. 2, no. 3, pp. 217– 237, 1995.
- [18] D. Braess and T. Sauer, "Bernstein polynomials and learning theory," *Journal of Approximation Theory*, vol. 128, no. 2, pp. 187 – 206, 2004.
 [19] R. Krichevsky and V. Trofimov, "The performance of universal encod-
- [19] R. Krichevsky and V. Trofimov, "The performance of universal encoding," *IEEE Trans. Information Theory*, vol. 27, no. 2, pp. 199–207, March 1981.
- [20] W. Gale and K. Church, "What's wrong with adding one," Corpus-Based Research into Language: In honour of Jan Aarts, pp. 189–200, 1994.
- 4) [21] K. P. Burnham and W. S. Overton, "Robust estimation of population size when capture probabilities vary among animals," *Ecology*, vol. 60, no. 5, pp. 927–936, 1979.
 - [22] I. Good and G. Toulmin, "The number of new species, and the increase in population coverage, when a sample is increased," *Biometrika*, vol. 43, no. 1-2, pp. 45–63, 1956.
 - [23] C. X. Mao and B. G. Lindsay, "Estimating the number of classes," the Annals of Statistics, pp. 917–930, 2007.
 - [24] J. Acharya, H. Das, A. Orlitsky, and A. T. Suresh, "A unified maximum likelihood approach for estimating symmetric properties of discrete distributions," in *International Conference on Machine Learning (PMLR)*, 2017, pp. 11–21.
 - [25] A. Orlitsky, N. P. Santhanam, K. Viswanathan, and J. Zhang, "On modeling profiles instead of values," in *Proceedings of the 20th conference* on Uncertainty in artificial intelligence, 2004, pp. 426–435.
 - [26] Y. Hao and A. Orlitsky, "The broad optimality of profile maximum likelihood," Advances in Neural Information Processing Systems, vol. 32, pp. 10991–11003, 2019.
 - [27] A. Orlitsky and A. T. Suresh, "Competitive distribution estimation: Why is Good-Turing good," in *Advances in Neural Information Processing Systems* 28, C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, Eds. Curran Associates, Inc., 2015, pp. 2143–2151.
 - [28] B. H. Juang and S. H. Lo, "On the bias of the Turing-Good estimate of probabilities," *IEEE Trans. Signal Processing*, vol. 42, no. 2, pp. 496–498, Feb. 1994.
 - [29] D. A. McAllester and R. E. Schapire, "On the convergence rate of Good-Turing estimators," in *Proceedings of the Thirteenth Annual Conference* on Computational Learning Theory. Morgan Kaufmann Publishers Inc., 2000, pp. 1–6.

- [30] A. B. Wagner, P. Viswanath, and S. R. Kulkarni, "Probability estimation in the rare-events regime," *IEEE Trans. Information Theory*, vol. 57, no. 6, pp. 3207–3229, June 2011.
- [31] J. Acharya, Y. Bao, Y. Kang, and Z. Sun, "Improved bounds for minimax risk of estimating missing mass," in *International Symposium* on Information Theory (ISIT), 2018, pp. 326–330.
- [32] D. Berend and A. Kontorovich, "The missing mass problem," *Statistics & Probability Letters*, vol. 82, no. 6, pp. 1102 1110, 2012.
- [33] —, "On the concentration of the missing mass," *Electronic Communications in Probability*, vol. 18, 2013.
- [34] Y. G. Yatracos, "On the rare species of a population," Journal of Statistical Planning and Inference, vol. 48, no. 3, pp. 321 – 329, 1995.
- [35] N. Rajaraman, A. Thangaraj, and A. T. Suresh, "Minimax risk for missing mass estimation," in *IEEE International Symposium on Information Theory (ISIT)*, 2017, pp. 3025–3029.
- [36] T. Routtenberg and L. Tong, "Estimation after parameter selection: Performance analysis and estimation methods," *IEEE Trans. Signal Processing*, vol. 64, no. 20, pp. 5268–5281, Oct. 2016.
- [37] E. Meir and T. Routtenberg, "Cramér-Rao bound for estimation after model selection and its application to sparse vector estimation," *IEEE Trans. Signal Process.*, vol. 69, pp. 2284–2301, Mar. 2021.
- [38] N. Harel and T. Routtenberg, "Low-complexity methods for estimation after parameter selection," *IEEE Trans. Signal Processing*, vol. 68, pp. 1152–1167, 2020.
- [39] —, "Bayesian post-model-selection estimation," *IEEE Signal Process. Lett.*, vol. 28, pp. 175–179, Jan. 2021.
- [40] T. Weiss, T. Routtenberg, and H. Messer, "Total performance evaluation of intensity estimation after detection," *Signal Processing*, vol. 183, p. 108042, 2021.
- [41] Z. Ben-Haim and Y. C. Eldar, "The Cramér-Rao bound for estimating a sparse parameter vector," *IEEE Trans. Signal Processing*, vol. 58, no. 6, pp. 3384–3389, June 2010.
- [42] E. L. Lehmann and J. P. Romano, *Testing Statistical Hypotheses*, 3rd ed. New York: Springer Texts in Statistics, 2005.
- [43] Z. Ben-Haim and Y. C. Eldar, "On the constrained Cramér-Rao bound with a singular Fisher information matrix," *IEEE Signal Processing Letters*, vol. 16, no. 6, pp. 453–456, 2009.
- [44] A. Cohen and H. B. Sackrowitz, "Admissibility of estimators of the probability of unobserved outcomes," *Annals of the Institute of Statistical Mathematics*, vol. 42, no. 4, pp. 623–636, 1990.
- [45] F. Gini, "Estimation strategies in the presence of nuisance parameters," *Signal processing*, vol. 55, no. 2, pp. 241–245, 1996.
- [46] S. Bar and J. Tabrikian, "Bayesian estimation in the presence of deterministic nuisance parameters; Part I: Performance bounds," *IEEE Trans. Signal Processing*, vol. 63, no. 24, pp. 6632–6646, Dec. 2015.
- [47] D. S. Pavlichin, J. Jiao, and T. Weissman, "Approximate profile maximum likelihood." *Journal of Machine Learning Research*, vol. 20, no. 122, pp. 1–55, 2019.
- [48] Y. Hao, A. Orlitsky, and V. Pichapati, "On learning Markov chains," in Advances in Neural Information Processing Systems, 2018, pp. 648–657.
- [49] M. Skorski, "Missing mass concentration for Markov chains," arXiv preprint arXiv:2001.03603, 2020.
- [50] S. M. Kay, Fundamentals of Statistical Signal Processing: Estimation Theory. Upper Saddle River, NJ, USA: Prentice Hall, 1993.
- [51] T. Routtenberg and J. Tabrikian, "Non-Bayesian periodic Cramér-Rao bound," *IEEE Trans. Signal Processing*, vol. 61, no. 4, pp. 1019–1032, Feb. 2013.
- [52] —, "Cyclic Barankin-type bounds for non-Bayesian periodic parameter estimation," *IEEE Trans. Signal Processing*, vol. 62, no. 13, pp. 3321–3336, July 2014.
- [53] S.-H. Lo, "From the species problem to a general coverage problem via a new interpretation," *The Annals of Statistics*, pp. 1094–1109, 1992.
- [54] S. Engen, Stochastic abundance models, with emphasis on biological communities and species diversity. London: Chapman and Hall; New York: Wiley, 1978.
- [55] E. Nitzan, T. Routtenberg, and J. Tabrikian, "Optimal biased estimation using Lehmann-unbiasedness," in *Proc. ICASSP*, 2017, pp. 4496–4500.
- [56] Y. C. Eldar, "Minimum variance in biased estimation: bounds and asymptotically optimal estimators," *IEEE Trans. Signal Processing*, vol. 52, no. 7, pp. 1915–1930, July 2004.
- [57] S. Kay and Y. C. Eldar, "Rethinking biased estimation," *IEEE Signal Processing Magazine*, vol. 25, no. 3, pp. 133–136, May 2008.
- [58] R. Zamir, "A proof of the Fisher information inequality via a data processing argument," *IEEE Trans. Inf. Theory*, vol. 44, no. 3, pp. 1246– 1250, 1998.

- [59] V. Berisha and A. O. Hero, "Empirical non-parametric estimation of the Fisher information," *IEEE Signal Processing Letters*, vol. 22, no. 7, pp. 988–992, 2014.
- [60] C. R. Rao, Linear statistical inference and its applications. Wiley New York, 1973, vol. 2.
- [61] G. T. Whipps, E. Ertin, and R. L. Moses, "Constrained Fisher scoring for a mixture of factor analyzers," US Army Research Laboratory Adelphi United States, Tech. Rep., 2016.
- [62] D. S. Pavlichin, J. Jiao, and T. Weissman. (2017) Approximate profile maximum likelihood. [Online]. Available: https://doi.org/10.5281/zenodo.1043617
- [63] S. T. Piantadosi, "Zipf's word frequency law in natural language: a critical review and future directions," *Psychon Bull Rev.*, vol. 21, no. 5, pp. 1112–1130, Oct. 2014.
- [64] H. Kesten and N. Morse, "A property of the multinomial distribution," *The Annals of Mathematical Statistics*, vol. 30, no. 1, pp. 120–127, 1959.
- [65] S. Boyd and L. Vandenberghe, *Convex Optimization*. New York, NY, USA: Cambridge University Press, 2004.