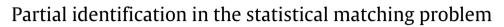
Contents lists available at ScienceDirect

**Computational Statistics and Data Analysis** 

journal homepage: www.elsevier.com/locate/csda





<sup>a</sup> Department of Mathematics, University of Queensland, Australia

<sup>b</sup> Public Health Foundation of India, IIPH Hyderabad, India

<sup>c</sup> CR Rao Advanced Institute of Mathematics, Statistics and Computer Science, Hyderabad, India

## ARTICLE INFO

Article history: Received 20 November 2015 Received in revised form 15 June 2016 Accepted 15 June 2016 Available online 22 June 2016

Keywords: Data integration Missing data Positive-definite matrix completion Statistical matching

## ABSTRACT

The statistical matching problem involves the integration of multiple datasets where some variables are not observed jointly. This missing data pattern leaves most statistical models unidentifiable. Statistical inference is still possible when operating under the framework of partially identified models, where the goal is to bound the parameters rather than to estimate them precisely. In many matching problems, developing feasible bounds on the parameters is equivalent to finding the set of positive-definite completions of a partially specified covariance matrix. Existing methods for characterising the set of possible completions do not extend to high-dimensional problems. A Gibbs sampler to draw from the set of possible completions is proposed. The variation in the observed samples gives an estimate of the feasible region of the parameters. The Gibbs sampler extends easily to high-dimensional statistical matching problems.

© 2016 The Authors. Published by Elsevier B.V. This is an open access article under the CC BY-NC-ND license (http://creativecommons.org/licenses/by-nc-nd/4.0/).

### 1. Introduction

The statistical matching problem involves the integration of multiple datasets where we have a set of variables common to all datasets, and other variables which only appear in some datasets. In the simplest terms, we have two samples A and Bof  $n_A$  and  $n_B$  independent observations, respectively, from the same population. In sample A we have measurements on sets of variables X and Y, and in sample B we have observations on variables X and Z. Our objective is to recover the joint density function f(x, y, z) from the lower dimensional datasets. The statistical matching problem is a special class of a missing data problem, where the defining characteristic is that we have no joint observations of Y and Z.

We often assume that the joint density function belongs to some parametric family { $f(\mathbf{x}, \mathbf{y}, \mathbf{z}; \theta) : \theta \in \Omega$ }, where  $\Omega$  denotes some parameter space. The objective is to perform statistical inference on the parameter  $\theta$ . Because of the missing data structure in the statistical matching scenario some of the parameters may be unidentifiable. Statistical inference is still possible if the model is viewed as a partially identified model. The concept of partially identified models stems from the belief that identification is not a simple binary issue. In a partially identified model, the range of values that the parameter  $\theta$  can take while leaving the observed data likelihood function unchanged is some non-trivial set. Informally, given an infinite dataset, under an identifiable model we can recover the true value of the parameters. In a partially identified model, given an infinite dataset, we are limited to being able to restrict the parameters to some feasible set. In a partially identified model, some elements of  $\theta$  may be point-wise identifiable while others are only partially identifiable.

Standard estimation approaches for missing data problems can exhibit pathological behaviour when the model is only partially identified. Use of the EM algorithm (Dempster et al., 1977) is complicated by the fact that the observed data

\* Correspondence to: MRC Biostatistics Unit, Cambridge, UK.

E-mail address: daniel.ahfock@uqconnect.edu.au (D. Ahfock).

http://dx.doi.org/10.1016/j.csda.2016.06.005





CrossMark

<sup>0167-9473/© 2016</sup> The Authors. Published by Elsevier B.V. This is an open access article under the CC BY-NC-ND license (http://creativecommons.org/licenses/by-nc-nd/4.0/).

#### Table 1

Missing data structure in the canonical statistical matching problem. Observed dimensions for each observation have been shaded.

		X	Y	Z
$s_{i}$	$_1^A$	$oldsymbol{s}_{1X}^A$	$oldsymbol{s}_{1Y}^A$	-
$s_{\tilde{z}}$	$A_2$	$oldsymbol{s}_{2X}^A$	$oldsymbol{s}_{2Y}^A$	-
:		••••		:
$oldsymbol{s}_n^{\mathbb{A}}$	$A_{\nu_A}$	$oldsymbol{s}^A_{n_AX}$	$oldsymbol{s}^A_{n_AY}$	-
s	$_1^B$	$oldsymbol{s}^B_{1X}$	-	$oldsymbol{s}_{1Z}^B$
$s_{2}^{2}$	$B_2$	$oldsymbol{s}^B_{2X}$	-	$oldsymbol{s}^B_{2Z}$
:		•••		
$oldsymbol{s}_n^E$	3 8	$s^B_{n_BX}$	-	$oldsymbol{s}^B_{n_BZ}$

likelihood will not have a unique maximiser, the likelihood will have a ridge over the allowable range of the partially identified parameters. It can be shown that the EM parameter estimates will converge to values that are located on likelihood ridge, and that the limiting estimates depend on the choice of initial values (Schafer, 1997, p. 53). This phenomenon is illustrated for the statistical matching problem in Section 4.2.2 of Rässler (2002). The high sensitivity to the initial values complicates the interpretation of the single point estimate returned by the EM algorithm. Bayesian approaches easily extend to partially identified models, however the posterior distribution can be highly sensitive to the prior, even in large samples (Gustafson, 2015). Credible intervals can also have very poor frequentist coverage (Moon and Schorfheide, 2012). We will pursue a frequentist strategy for estimating the partially identified parameters.

In the statistical matching problem, the partially identified parameters are often elements of a covariance matrix. It is typical to have all elements of the covariance matrix identifiable, other than the values that require joint observations on **Y** and **Z**. In this setting, estimating the identified set corresponds to determining the set of positive-definite completions of a partially specified covariance matrix. Existing methods for doing so are not applicable when both **Y** and **Z** are multivariate (D'Orazio, 2015). We take a new sampling based approach to characterising the identified set which is easily applicable to high-dimensional problems. We propose a Gibbs sampler to draw values uniformly from the identified set of covariance parameters. The range of the sampled values gives a direct measure of the uncertainty attached to the partially identified parameters. The Gibbs sampler extends the range of datasets that can be analysed using the statistical matching methodology.

## 2. The statistical matching problem

A standard mathematical description of the statistical matching problem is as follows (Rässler, 2002). Let X, Y, Z be multivariate random variables with joint density function  $f(x, y, z; \theta)$ . Assume we have a sample of  $n_A$  i.i.d. observations distributed according to  $f(x, y, z; \theta)$ , which we will call file A, and another independent sample of size  $n_B$  from  $f(x, y, z; \theta)$ , which we will call file B. Let  $s_i^A$  be a row vector representing the *i*th observation in file A for  $i = 1, ..., n_A$ . Similarly, let  $s_j^B$  be a row vector representing the *j*th observation in file B for  $j = 1, ..., n_B$ . The *i*th observation in file A can be written as  $s_i^A = (s_{i\chi}^A, s_{i\gamma}^A, s_{iZ}^A)$ , where  $s_{i\chi}^A$  is a row vector representing the value of X and  $s_{i\chi}^A, s_{iZ}^A$  are row vectors representing the values of Y and Z, respectively. We can also form an identical partition  $s_j^B = (s_{j\chi}^B, s_{j\chi}^B, s_{jZ}^B)$  for observation j in file B. Let the observations in file A have the Z values missing and the observations in file B have the Y values missing. Table 1 represents the data matrix in the statistical matching problem. We can consider inference in the statistical matching problem to be inference under a partially identified model. We call a model partially identified if the observed data likelihood is flat for a range of the parameters (Tamer, 2010). The identified set for a parameter is the range of values it can take without altering the observed data likelihood function. We use the notation  $\Theta(\theta_j)$  to denote the identified set for parameter  $\theta_j$ . When analysing a partially identified model we are interested in forming a set of plausible values for the non point-identified parameters. For example, assume we have observations  $(X, Y, Z)^T$  from a trivariate normal distribution,

$$N_3\left(\boldsymbol{\mu} = \begin{bmatrix} \mu_X \\ \mu_Y \\ \mu_Z \end{bmatrix}, \boldsymbol{\Sigma} = \begin{bmatrix} \sigma_{XX} & \sigma_{XY} & \sigma_{XZ} \\ \sigma_{YX} & \sigma_{YY} & \sigma_{YZ} \\ \sigma_{ZX} & \sigma_{ZY} & \sigma_{ZZ} \end{bmatrix}\right)$$

and the standard statistical matching problem applies. The likelihood function formed from the observed data will not depend on  $\sigma_{YZ}$ , and so  $\sigma_{YZ}$  can be considered to be a partially identified parameter. All the parameters are point-wise identifiable other than  $\sigma_{YZ}$ . Even though we do not have any data to estimate  $\sigma_{YZ}$  from, as we do not observe Y and Z jointly, our modelling assumptions induce non-trivial bounds on the parameter. Given the other parameters, the possible

Download English Version:

# https://daneshyari.com/en/article/6868943

Download Persian Version:

https://daneshyari.com/article/6868943

Daneshyari.com