

Nielsen, B. & Nielsen, J. P. (2014). Identification and forecasting in mortality models. The Scientific World Journal, 2014, 347043 - ?. doi: 10.1155/2014/347043



**CITY UNIVERSITY  
LONDON**

[City Research Online](#)

**Original citation:** Nielsen, B. & Nielsen, J. P. (2014). Identification and forecasting in mortality models. The Scientific World Journal, 2014, 347043 - ?. doi: 10.1155/2014/347043

**Permanent City Research Online URL:** <http://openaccess.city.ac.uk/4638/>

### **Copyright & reuse**

City University London has developed City Research Online so that its users may access the research outputs of City University London's staff. Copyright © and Moral Rights for this paper are retained by the individual author(s) and/ or other copyright holders. All material in City Research Online is checked for eligibility for copyright before being made available in the live archive. URLs from City Research Online may be freely distributed and linked to from other web pages.

### **Versions of research**

The version in City Research Online may differ from the final published version. Users are advised to check the Permanent City Research Online URL above for the status of the paper.

### **Enquiries**

If you have any enquiries about any aspect of City Research Online, or if you wish to make contact with the author(s) of this paper, please email the team at [publications@city.ac.uk](mailto:publications@city.ac.uk).

## Research Article

# Identification and Forecasting in Mortality Models

**Bent Nielsen<sup>1,2,3</sup> and Jens P. Nielsen<sup>4</sup>**

<sup>1</sup> Department of Economics, University of Oxford, Oxford OX1 2JD, UK

<sup>2</sup> Programme on Economic Modelling, INET, University of Oxford, Oxford OX1 2JD, UK

<sup>3</sup> Nuffield College, Oxford OX1 1NF, UK

<sup>4</sup> Cass Business School, City University London, 106 Bunhill Row, London EC1Y 8TZ, UK

Correspondence should be addressed to Bent Nielsen; bent.nielsen@nuffield.ox.ac.uk

Received 22 January 2014; Accepted 17 April 2014; Published 2 June 2014

Academic Editor: Montserrat Guillén

Copyright © 2014 B. Nielsen and J. P. Nielsen. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Mortality models often have inbuilt identification issues challenging the statistician. The statistician can choose to work with well-defined freely varying parameters, derived as maximal invariants in this paper, or with ad hoc identified parameters which at first glance seem more intuitive, but which can introduce a number of unnecessary challenges. In this paper we describe the methodological advantages from using the maximal invariant parameterisation and we go through the extra methodological challenges a statistician has to deal with when insisting on working with ad hoc identifications. These challenges are broadly similar in frequentist and in Bayesian setups. We also go through a number of examples from the literature where ad hoc identifications have been preferred in the statistical analyses.

## 1. Introduction

Mortality models are commonly used in a wide range of fields such as actuarial sciences, epidemiology, and sociology. They are often used in important decisions such as how to deal with unisex legislation in the pension industry; see Ornelas et al. [1] and Jarner and Kryger [2]. However, such models do often have inbuilt identification issues stemming from overparametrisation. While identification issues are omnipresent in statistical modelling, this paper focuses on mortality modelling, where estimated parameters are treated as time series and extrapolated to give forecasts of future mortality. The underlying theme of this paper is to provide strategies of avoiding arbitrariness resulting from the identification process. We suggest two ways forward. First, we can reparametrise the model in terms of a freely varying parameter, which therefore has to be of lower dimension than the original parameter. Secondly, we can work with an identified version of the original parameter as long as we keep track of the consequences of the identification choice. That way we ensure that two researchers making different identification choices get the same statistical inferences and forecasts.

A simple example is the age-period model for an age-period array of mortality rates. It is well-known that the levels of the age- and period-effects cannot be determined from the likelihood representing the overparametrisation of the model. When the estimated age- and period-effects are treated as time series and subjected to plotting and extrapolation, then our approach ensures that the statistical analysis is the same for two researchers identifying the above model in two different ways. Whereas this issue is relatively simple for the age-period model, identification becomes more tricky for complicated models such as the age-period-cohort model and the model of Lee and Carter [3], let alone two-sample situations.

Mortality models are built as a combination of age, period, and cohort-effects, but the likelihood only varies with a surjective function of these time effects. The time effects can be divided into two parts. One part that moves the likelihood function and another part which does not induce variation in the likelihood function. We will argue that all inferences and forecasts should be concerned primarily with the part of the parameter that moves the likelihood function. This does not preclude the researcher from working with the time effects, but it gives some limitations on what can be done.

This is important because the motivation and the intuition of mortality models typically originate in the time effects. For instance, in the context of an age-period-cohort model linear trends cannot be identified so time series plots of the time effects need to be invariant to linear trends and extrapolations of time effects must preserve the arbitrary linear trend in the time effects. This applies regardless of whether the identification issue is dealt with in a frequentist manner or by Bayesian methods.

To formalise the discussion slightly return to the age-period example. Denote the predictor for the age-period data array by  $\mu$ . The age-period model then determines how the predictor  $\mu$  varies with a vector  $\theta$  summarising age and period effects. That vector is split into two components  $\xi$  and  $\lambda$  so that the predictor only depends on  $\theta$  through  $\xi$  but not on  $\lambda$  which cannot be identified by statistical analysis. In the age-period example  $\xi$  could reflect the contrasts and the overall level of the predictor  $\mu$ , whereas  $\lambda$  reflects the level of the age effect. The more principled solution is then to work exclusively with  $\xi$  and simply consider  $\theta$  as a motivation rather than the objective of the analysis. Another solution is to ad hoc identify  $\lambda$  based on a notion of mathematical convenience or based on a particular purpose given the substantive context.

Once an ad hoc identification of  $\lambda$  is chosen the identification problem appears to go away, because the likelihood analysis can now go through. The reason is that the variation of  $\theta$  is now reduced to the variation of  $\xi$  precisely because  $\lambda$  is fixed. Suppose two researchers choose the same likelihood and the same parametrisation of  $\xi$  but different ad hoc identifications  $\lambda^\dagger$  and  $\lambda^\ddagger$ . Which of their conclusions will be the same and which will be different? As the likelihood only depends on  $\xi$  the fits of the two researchers will be identical. But differences might arise if the statistical inference or forecasting or any other statistical analysis involves  $\lambda$  in some way.

Indeed, with many extrapolation methods forecasts will be invariant to the choice of  $\lambda$ . But, there will also be extrapolation methods where this is not the case. Examples arise in the age-period-cohort model, where linear trends have to be handled with care.

We will start by analysing linearly parametrised models at a rather general level. We do this with two aspects in mind. First, we need to step back to a point in the analysis before ad hoc identification is made. Secondly, we also want to avoid the discussion of how to choose  $\xi$  and  $\lambda$ , which tend to be specific to the mortality model in question. Working at the general level we can focus on the mappings between different parametrisations and the invariance properties coming from these mappings. It is then seen that the parameter  $\xi$  arises as a maximal invariant. The general setting also allows the formulation of a series of results discussing different types of ad hoc identification, first in a frequentist fashion and then in a Bayesian fashion.

Subsequently, we will consider the age-period-cohort model in detail, both for one- and two-sample situations. Using the general results it becomes easier to see that a number of popular methods inadvertently include features that are not invariant to ad hoc identification. These include

the ‘‘intrinsic estimator’’ advocated by Yang et al. [4], the ‘‘mixed model approach’’ by Yang and Land [5], the Bayesian approach by Berzuini and Clayton [6], and the two-sample analysis by Riebler and Held [7]. Finally, we consider the nonlinearly parametrised model of Lee and Carter [3]. The nonlinearity gives a further complication since the mapping from the time effects to the mortality predictor is nondifferentiable. As it turns out the mortality predictor varies in a smooth space, so the nondifferentiability is avoided by working directly with the mortality predictor instead of the original time effects. Instead, a Lee-Carter application should consider whether a certain matrix has rank of unity or zero. Apart from that the analysis is similar to that of linearly parametrised models. Likewise a theory is given for two-sample situations.

Throughout the paper our concern rests exclusively with the identification problem and the consequences of ad hoc identification for estimation, plots, inference, and forecasting. In practice, important additional concerns are how to choose appropriate models and forecasting methods. We would like to refer to Girosi and King [8], Pitacco et al. [9] for general discussions of these issues, and also to Kuang et al. [10] and Coelho and Nunes [11] for discussions of forecast methods in the light of structural breaks. Instead, the aim of the paper is to present an overall framework that can help streamlining the identification discussion that has appeared in so many papers in so many fields over so many years.

Section 2 of this paper considers standard linear statistical models, which lend themselves to a relative straightforward analysis based on linear algebra. Any ad hoc identification splits the time effect into two components. The first component is an arbitrary component, which is not needed for the identification of the likelihood. The other component is necessary and sufficient to identify the model and hence sufficient for statistical analysis. In Section 3 it is outlined how to analyze the statistical model when the latter component is ad hoc identified. It is argued that this can cause difficulties for estimation, interpretation, and forecast. In Section 4 it is shown that Bayesian analysis shares the same challenges as the frequentist approach. In Sections 5 and 6 we study the two particular examples: the omnipresent age-period-cohort and Lee-Carter mortality models. All proofs are collected in the Appendix.

## 2. Statistical Models with Linear Parametrisations

In this section we present the identification problem in a linear framework. The problem is solved by analysing the mapping from the original time effect to the predictor which, in turn, leads to standard statistical analysis. In Section 6 we show how these ideas transfer to a nonlinear context. This contrasts with Section 3 in which we illustrate the analytical challenges and inconveniences arising from ad hoc identification.

In Section 2.1 we present the overparametrized linear model for the mortality predictor. The identification problem

is defined in Section 2.2 via the likelihood. In an overparametrized linear model two different parameters might produce the same likelihood. In Section 2.3 we analyze the mapping from the overparametrised parameter to the predictor. This mapping enables us to split the overparametrised parameter into two. One arbitrary parameter and one parameter identify the model without being overparametrised. This latter parameter is shown to be a maximal invariant parameter. In Section 2.4 it is demonstrated how any statistical analysis can be based on this maximal invariant parameter alone. In particular we comment that visual data representations, hypothesis testing, and forecasting are simple and well defined. This in turn leads to standard statistical analysis.

The analysis of the linearly parametrised involves projections on linear or affine spaces and on their orthogonal complements. It is therefore convenient to introduce the following notation. A matrix  $m$  has full column rank if  $m'm$  is invertible. In this case the orthogonal complement  $m_{\perp}$  is a matrix so  $m'_{\perp}m = 0$  and  $(m, m_{\perp})$  is invertible. Thus, when  $m$  itself is invertible then  $m_{\perp}$  is the empty matrix. It is not difficult to calculate  $m_{\perp}$  in practise, an explicit construction of  $m_{\perp}$  follows from a singular value decomposition of  $mm'$ , choosing  $m_{\perp}$  as the eigenvectors associated with the zero eigenvalues. Moreover, let  $\bar{m} = m(m'm)^{-1}$  so that  $\bar{m}'m$  is the identity matrix, while  $\bar{m}_{\perp} = m_{\perp}(m'_{\perp}m_{\perp})^{-1}$ .

**2.1. The Model.** Think of the time effect  $\theta$  as our preferred intuitive, but unidentified parameter, and think of the predictor  $\mu$  as some function of  $\theta$  specifying the model at hand. In a Poisson type model, where the mean specifies the distribution,  $\mu$  could be the log of that mean. Such Poisson models are omnipresent in mortality models. We will often think of  $\theta$  as containing some time effects. Often forecasting is carried out simply by isolating and extrapolating such a time effect.

Consider a data vector  $Y$  of dimension  $n$ . This could, for instance, be the vector consisting of the stacked mortality rates for a rectangular age-period array of dimension  $I \times J$  in which case  $n = IJ$ . The statistical model for  $Y$  could be a generalized linear model. This involves an appropriately chosen distribution and a link function, which links the expected mortality rate to an  $n$ -dimensional predictor, which is denoted by  $\mu$ . Taken together this defines a likelihood function  $L(\mu; Y)$ .

The model for the predictor  $\mu$  is constructed in terms of, for instance, age, period, and cohort time effects. These time effects are summarized in a vector  $\theta$ , which is of dimension  $q < n$ . Therefore  $\mu$  is a surjective function of  $\theta$ . For the moment the specification of the predictor is assumed linear so that

$$\mu = D\theta \quad \text{for } \theta \in \Theta = \mathbb{R}^q, \tag{1}$$

for some design matrix  $D \in \mathbb{R}^{n \times q}$ . We refer to this specification as the mortality model, while the space  $\Theta$  is the time effect space. The time effect space is chosen as an unrestricted real space in accordance with the starting point of most mortality analyses.

The parameter space for the likelihood function and therefore for the statistical model is given by the range of variation for the predictor  $\mu$ ; that is,

$$M = (\mu \in \mathbb{R}^n : \mu = D\theta \text{ for } \theta \in \Theta = \mathbb{R}^q). \tag{2}$$

The likelihood function is assumed uniquely identified on this space in the sense that for all pairs of predictors so  $\mu^{\dagger} \neq \mu^{\ddagger}$  then the likelihood of  $\mu^{\dagger}, \mu^{\ddagger}$  differ; that is,

$$L(\mu^{\dagger}; Y) \neq L(\mu^{\ddagger}; Y), \tag{3}$$

for  $Y$  in a set with positive probability.

**2.2. The Identification Problem.** The identification problem of mortality models arises when the mapping from the time effect space  $\Theta$  to the parameter space  $M$  is surjective but not injective. With a linear parametrisation this arises when the design matrix  $D$  has reduced column rank  $p < q$  so  $D'D$  is singular. In this situation there exists time effects  $\theta^{\dagger} \neq \theta^{\ddagger}$  with the same likelihood:

$$L(D\theta^{\dagger}; Y) = L(D\theta^{\ddagger}; Y), \tag{4}$$

for all data  $Y$ . Then the time effect space  $\Theta$  is not useful as parameter space for the statistical model.

**2.3. Analysing the Mapping  $\theta \mapsto \mu$ .** When analysing the mapping from our intuitively preferred parametrisation  $\theta$  into the linear predictor  $\mu$ , we will be able to rewrite  $\theta$  as a sum of two components: one is a function of the predictor and the other is the arbitrary part varying with  $\theta$ , but not with the predictor. We provide two methods for analysis.

The first method is to find a basis  $X \in \mathbb{R}^{n \times p}$  with full column rank  $p$  for the design  $D$ . The design matrix of the mortality model can then be expressed as  $D = XA'$  for some matrix  $A \in \mathbb{R}^{q \times p}$  with full column rank  $p$ . Introduce a new  $p$ -dimensional parameter:

$$\xi = A'\theta. \tag{5}$$

The parameter space  $M$  can then be written more parsimoniously as

$$M = (\mu \in \mathbb{R}^n : \mu = X\xi \text{ for } \xi \in \mathbb{R}^p). \tag{6}$$

The mapping from  $\xi$  to  $\mu$  is bijective, so the statistical model can just as well be parametrised in terms of  $\xi \in \Xi = \mathbb{R}^p$ .

Alternatively, the identification problem can be expressed through an invariance argument. This argument relates to the parameterization but resembles the classical invariance argument for reduction of data; see Cox and Hinkley [12, page 157]. With a linear parametrisation the argument involves the orthogonal complement to the matrix  $A$ . That is a matrix  $A_{\perp} \in \mathbb{R}^{q \times (q-p)}$  which has the properties that  $A'_{\perp}A = 0$  and that  $(A, A_{\perp})$  is invertible. The mortality model (1) is defined by the mapping

$$\theta \mapsto \mu = D\theta = XA'\theta, \tag{7}$$

from  $\Theta = \mathbb{R}^q$  to  $M$ . This mapping is surjective in that two different values of  $\theta$  may result in the same  $\mu$  and therefore the same likelihood. These equivalence classes in the time effect space can be described by the group of transformations

$$g : \theta \mapsto \theta + A_{\perp} \zeta, \tag{8}$$

acting on  $\Theta$  for arbitrary  $\zeta \in \mathbb{R}^{q-p}$ . Indeed, it holds that  $\theta$  and  $g(\theta)$  will result in the same  $\mu$ . The mapping (7) is therefore invariant to the group  $g$ . We will argue that the parameter  $\xi = A'\theta$  is a maximal invariant to the group  $g$  acting on  $\Theta$ , which provides a link with (6). It has to be argued that for any  $\theta^{\dagger}, \theta^{\ddagger}$  so that  $\xi^{\dagger} = A'\theta^{\dagger}$  equals  $\xi^{\ddagger} = A'\theta^{\ddagger}$  then  $\theta^{\ddagger} = g(\theta^{\dagger})$ , see Cox and Hinkley [12, page 159]. For this argument use the orthogonal projection identity to write

$$\theta = A(A'A)^{-1}\xi + A_{\perp}(A'_{\perp}A_{\perp})^{-1}\varphi; \tag{9}$$

for unique  $\xi = A'\theta$  and  $\varphi = A'_{\perp}\theta$ . Thus, if  $A'\theta^{\ddagger} = A'\theta^{\dagger}$  then  $\theta^{\ddagger} = g(\theta^{\dagger})$  with  $\zeta = \varphi^{\ddagger} - \varphi^{\dagger} = A'_{\perp}(\theta^{\ddagger} - \theta^{\dagger})$ .

In applications it can be difficult to find a basis  $X$  for the design  $D$ . It can be easier to find a group  $g$  and hence  $A_{\perp}$  and then use this information to construct  $A$  and a candidate basis  $X = D\bar{A}$ , noting that  $D = XA'$ . This argument leaves it to be proven that  $X$  is a basis, or equivalently, that the suggested group  $g$  actually describes the equivalence classes of the mapping from  $\theta$  to  $\mu$ .

It is useful to note that in the choices of  $X$ ,  $A$  only the spaces spanned by them are unique since  $XA' = Xmm^{-1}A'$  for any invertible  $m \in \mathbb{R}^{p \times p}$ . Likewise, the maximal invariant  $\xi$  is only unique up to bijective transformations. This lack of uniqueness has no impact on the analysis of the likelihood albeit it influences interpretations.

**2.4. Statistical Analysis Using the Maximal Invariant Parameter.** The statistical model parametrised with the maximal invariant parameter  $\xi$  can be analysed by standard statistical techniques. This contrasts to a range of problems that arise when working with an ad hoc identified time effect  $\theta$ . In the following the relatively simple standard statistical analysis of the model parametrised by  $\xi$  is discussed with respect to likelihood theory, interpretation, plots, hypothesis testing, forecasting, and Bayesian analysis. In Sections 3 and 4 we give an overview of the much more complicated theory underpinning models parametrised by the ad hoc identified time effect  $\theta$ . Age-period-cohort examples follow in Section 5.

**2.4.1. Exponential Family Theory.** Suppose the likelihood is drawn from a generalized linear model based on an exponential family. Then the model is actually a regular exponential family where the maximal invariant parameter  $\xi$  is the canonical parameter since it is freely varying in a real space; see Barndorff-Nielsen [13, page 116]. This opens up for a wealth of convenient statistical properties such as a likelihood equation with a simple expression and explicit conditions for a unique solution. In contrast, ad hoc identified parameters are based on an injective mapping of the canonical parameter

$\xi$  into  $\theta$ ; see Sections 3.1 and 3.2. It is then more difficult to fully exploit the exponential family theory.

**2.4.2. Interpretation and Plots.** The maximal invariant parameter  $\xi$  varies freely in  $\mathbb{R}^p$ . It can therefore be interpreted as the parameter of any standard statistical model. Since  $\xi$  is freely varying the coordinates of  $\xi$  can be interpreted independently. When  $\theta$  is a collection of time effects then  $\xi$  can be organised as a collection of time series. Since the coordinates of  $\xi$  are freely varying the time series plots of the components of  $\xi$  have the usual interpretation of time series. In contrast, ad hoc identified estimators are constrained to a  $p$ -dimensional subspace  $\Theta_{\lambda}$  of  $\Theta = \mathbb{R}^q$ , which is often affine but can be more complicated. A consequence is that plots are complicated to evaluate; see Section 3.4.1.

**2.4.3. Hypothesis Testing.** Hypotheses are easily formulated and analysed when using the maximal invariant parametrisation. An affine hypothesis that restricts  $\xi$  to vary in a  $p_H$ -dimensional affine subspace can be formulated as  $H'\xi = \eta$  for known matrices  $H \in \mathbb{R}^{p \times (p-p_H)}$ ,  $\eta \in \mathbb{R}^{p-p_H}$ . This implies a restriction on the predictor  $\mu = X\xi$  of (6). Form the orthogonal complement  $H_{\perp}$  and recall the orthogonal projection identity  $I_n = \bar{H}H' + H_{\perp}\bar{H}'_{\perp}$  so that  $\mu = X\bar{H}H'\xi + XH_{\perp}\bar{H}'_{\perp}\xi$ . Introduce a  $p_H$ -dimensional parameter  $\varphi = \bar{H}'_{\perp}\xi$ , a design matrix  $X_H = XH_{\perp}$ , and an offset  $Z_H = X\bar{H}\eta$ . The restricted parameter space is

$$M_H = (\mu \in \mathbb{R}^n : \mu = X_H\varphi + Z_H \text{ for } \varphi \in \mathbb{R}^{p_H}). \tag{10}$$

In an exponential family context both the unrestricted model and the restricted model form regular exponential families. A variety of nice properties then follow for the estimators and the test statistics from the exponential family theory. Examples are given in Sections 5.3 and 5.5.3. In contrast, the hypothesis derived from restrictions on ad hoc identified parameters and the resulting degrees of freedom are complicated to analyse; see Section 3.4.2.

**2.4.4. Forecasting.** Most often the objective of a mortality study is to forecast the future mortality. In the linear context,  $\mu = X\xi$ , this is done by extending the design  $X$  and by extrapolating  $\xi$ .

It is usually easy to extend the design  $X$  into the forecast horizon. This involves the construction of a triangular block matrix with an appropriate number of extra rows corresponding to the data over the forecast horizon as well as extra columns representing the extra parameters that would be needed:

$$X^h = \begin{pmatrix} X & 0 \\ X_1^h & X_2^h \end{pmatrix}. \tag{11}$$

Extrapolating  $\xi$  into a vector  $\tilde{\xi}$  then gives the forecast

$$\tilde{\mu} = (X_1^h, X_2^h) \begin{pmatrix} \xi \\ \tilde{\xi} \end{pmatrix}. \tag{12}$$

The extrapolation of the parameter  $\xi$  can be done as follows. The estimated parameter, or part of it, can be thought of as a time series. Any forecast techniques from the time series literature applied directly to  $\xi$  can be used, subject to the usual contextual considerations.

Ad hoc identified time effects can be extrapolated in a similar way; see Section 3.4.3. This may, however, result in avoidable arbitrary effects in the forecast. Necessary and sufficient conditions for this eventuality are given for age-period-cohort models in Section 5.4.3. The practical examples are mainly Bayesian in nature and are discussed next.

**2.4.5. Bayesian Analysis.** The introduction of the canonical parameter shows that the likelihood, in Bayesian notation, is of the form  $p(y | \theta) = p(y | \xi)$  where  $\xi$  is freely varying. A purist Bayesian analysis can simply introduce a prior on the canonical parameter,  $p(\xi)$ . This is updated in a straight forward way, resulting in the posterior  $p(\xi | y) = p(y | \xi)p(\xi)/p(y)$ .

In contrast, introducing a prior on ad hoc identified parameters gives various difficulties. Only parts of the prior are updated by the likelihood, so that it becomes unclear which information arises from the data and which information arises from the ad hoc identification. Moreover, avoidable arbitrariness is introduced in the forecast; see Section 4. Introduction of hyperparameters exacerbates the issue. Examples are given in Sections 5.4.4, 5.5.2, and 6.1.6.

### 3. Working with the Time Effects

In Section 2 we considered the situations where estimation, hypothesis testing a hypothesis, or forecasting is carried out using the canonical parameter. However, there might be situations, where the original time effect parametrisation is preferred, perhaps because it is felt that this parametrisation is particularly helpful in guiding the intuition. This requires ad hoc identification of the time effect. In this section we will guide the considerations a statistician that has to go through when insisting on an analysis based on some nonunique parametrisations. As in Section 2 we focus on linearly parametrised models. Specific examples follow in Sections 5 and 6.

In Section 3.1 ad hoc identification is defined. As an example we consider a least squares estimation problem with collinear regressors in Section 3.2. For the age-period-cohort model reviewed in Section 5 it is common to ad hoc identify in two steps: first identifying levels then the linear trends. We consider such two-step ad hoc identification in Section 3.3. The consequence of ad hoc identification is considered in Section 3.4. Indeed, when forecasting the time effect, we do not want the forecast to depend on the identification scheme. The same applies to graphical visualisation of our data, where the eye may extract patterns that depend on the identification scheme. Likewise, confusion may arise when formulating a hypothesis directly on the time effect parameters.

**3.1. Ad Hoc Identification.** In this section the time effect parametrisation is considered. An identification scheme has

to be introduced when working with the time effects. This may rest on mathematical convenience or it may be chosen for a particular purpose given the substantive context. We therefore call it ad hoc identification. Here we consider a simple identification scheme but turn to a more common two-step identification scheme in Section 3.3.

Once the canonical parameter  $\xi$  has been estimated there is often a wish to return to the original time effect  $\theta$ . The two are linked through the surjective mapping

$$\theta \mapsto \xi = A'\theta, \tag{13}$$

from  $\Theta = \mathbb{R}^q$  to  $\Xi = \mathbb{R}^p$ . Indeed, since  $\xi$  is constructed as a function of  $\theta$  the notation for  $\xi$  is often chosen to reflect  $\theta$ . The canonical parameter  $\xi$  does, however, only give partial information about  $\theta$ . The remaining part, say  $\lambda$ , of  $\theta$  will have to be chosen by the researcher and combined with  $\xi$ .

A linear ad hoc identification of  $\theta$  comes about by the researcher choosing a constraint

$$L'\theta = \lambda \tag{14}$$

for some known  $\lambda \in \mathbb{R}^{q-p}$  and some matrix  $L \in \mathbb{R}^{q \times (q-p)}$  chosen so the square matrix  $(A, L)$  is invertible. The time effect space  $\Theta$  is now reduced to an affine subspace

$$\Theta_\lambda = (\theta_\lambda \in \Theta : L'\theta_\lambda = \lambda). \tag{15}$$

Given  $\theta$  we can find  $\xi, \lambda$  through (13) and (14) as  $(\xi', \lambda')' = (A, L)'\theta$ . At the same time, given values of  $\xi, \lambda$  and the invertibility of  $(A, L)$ , the ad hoc identified time effect is found through

$$\theta_\lambda = \begin{pmatrix} A' \\ L' \end{pmatrix}^{-1} \begin{pmatrix} \xi \\ \lambda \end{pmatrix} = L_\perp (A' L_\perp)^{-1} \xi + A_\perp (L' A_\perp)^{-1} \lambda. \tag{16}$$

In this notation a subindex  $\lambda$  is introduced to avoid confusion with the time effect  $\theta$  in the original mortality model. Indeed, there are now four different parameters in play, namely, the original time effect  $\theta \in \Theta$ , the predictor  $\mu \in M$ , the maximal invariant parameter  $\xi \in \Xi$  and the ad hoc identified time effect  $\theta_\lambda \in \Theta_\lambda$ , each of which has a different interpretation. The mapping from  $\theta$  to each of  $\mu, \xi$ , and  $\theta_\lambda$  is surjective, while there are bijective mappings between the latter three. The interpretations of the time effect  $\theta$  and the canonical parameter  $\xi$  will inevitably be different. For a start they have different dimensions. Endowing the spaces with Euclidean norms shows that distances in the two spaces  $\Theta$  and  $\Xi$  will be judged differently. The time effect  $\theta$  and the ad hoc identified time effect  $\theta_\lambda$  will similarly have different interpretations. Although they have the same dimensions the Euclidean norms on  $\Theta$  and  $\Theta_\lambda$  will be rather different. Confusion may arise in the interpretation of a mortality analysis if there is no clear distinction between  $\theta$  and  $\theta_\lambda$ . In addition an unnecessary arbitrariness may arise when making inference on  $\theta_\lambda$  or extrapolating  $\theta \in \Theta_\lambda$ . We will return to these issues in Section 3.4.

It is perhaps interesting to note that despite the linear parametrisation the ad hoc identification need not be done

in a linear fashion as in (14). Indeed it is common for Poisson models with a log link to ad hoc identify  $\theta$  through the original multiplicative scale. That means that the ad hoc identification is done nonlinearly through

$$\xi = A'\theta_\lambda, \quad \lambda = f(\theta_\lambda). \tag{17}$$

The fit of the model is unaffected by the ad hoc identification. Indeed the fit is measured in terms of the estimate of the predictor  $\mu = D\hat{\theta}_\lambda$  where  $D = XA'$ . Since the identification is made so  $\xi = A'\hat{\theta}_\lambda$ ; the estimated predictor reduces to

$$\hat{\mu} = D\hat{\theta}_\lambda = XA'\hat{\theta}_\lambda = X\hat{\xi}, \tag{18}$$

regardless of the choice of ad hoc identification.

**3.2. A Least Squares Example.** As an illustration of estimation in the presence of ad hoc identification consider a normal likelihood. Different, but equivalent, expressions can be found depending on the parametrisation. The likelihood of the predictor  $\mu$  is

$$L(\mu, \sigma^2; Y) = (2\pi\sigma^2)^{-n/2} \exp\left\{-\frac{1}{2\sigma^2}(Y - \mu)'(Y - \mu)\right\} \\ \text{for } \mu \in M, \quad \sigma^2 > 0. \tag{19}$$

Rewriting it in terms of the canonical parameter it is

$$L(\xi, \sigma^2; Y) = (2\pi\sigma^2)^{-n/2} \exp\left\{-\frac{1}{2\sigma^2}(Y - X\xi)'(Y - X\xi)\right\} \\ \text{for } \xi \in \Xi = \mathbb{R}^p, \quad \sigma^2 > 0, \tag{20}$$

while introducing the time effect parameter gives

$$L(\theta, \sigma^2; Y) \\ = (2\pi\sigma^2)^{-n/2} \exp\left\{-\frac{1}{2\sigma^2}(Y - XA'\theta)'(Y - XA'\theta)\right\} \\ \text{for } \theta \in \Theta = \mathbb{R}^q, \quad \sigma^2 > 0. \tag{21}$$

The likelihood (20) of the canonical parameter  $\xi$  can be analysed by the least squares method since the design  $X$  has full column rank. The maximum likelihood estimator for  $\xi$  and the predictor for the data are

$$\hat{\xi} = (X'X)^{-1}X'Y, \quad \hat{Y} = X\hat{\xi} = X(X'X)^{-1}X'Y. \tag{22}$$

Along with the residual variance this is all the information that is given by the likelihood.

The likelihood (21) of the time effect  $\theta$  only depends on  $\theta$  through  $\xi = A'\theta$ . The lack of identification means that the maximum likelihood estimator for  $\theta$  has an arbitrary element, so that it is a set valued estimator. Based on (16) this can be expressed by

$$\hat{\Theta}_\lambda = L_\perp(A'L_\perp)^{-1}\hat{\xi} + A_\perp(L'A_\perp)^{-1}\lambda \\ \text{where } \hat{\Theta}_\lambda \in \Theta_\lambda \subset \Theta, \tag{23}$$

for any  $L$  so  $(A, L)$  is invertible and for any  $\lambda \in \mathbb{R}^{q-p}$ . The fit, however, remains the same and (18) becomes

$$\hat{\mu} = D\hat{\theta}_\lambda = XA' \left\{ L_\perp(A'L_\perp)^{-1}\hat{\xi} + A_\perp(L'A_\perp)^{-1}\lambda \right\} \\ = X\hat{\xi} = \hat{Y}. \tag{24}$$

In order to compute actual estimates then  $L, \lambda$  have to be chosen, which amounts to ad hoc identification. For instance, with the ad hoc identifying restrictions  $L = A_\perp$  and  $\lambda = 0$  then  $\hat{\theta}_\lambda$  can be thought of as the least squares estimator of  $Y$  on  $D$  using the Moore-Penrose generalised inverse for the singular matrix  $D'D$ ; see Searle [14, page 212]. See Section 5.4.1 for an example.

**3.3. Step-Wise Identification.** It is common to ad hoc identify parameter in a step-wise fashion. In the first step the time effect parameter is only partially constrained. The full identification then follows in a second step. An example is given in Section 5.4.1 for an age-period-cohort model in which the levels of the time effects are constrained in the first step leaving the ad hoc identification of the linear trends to the second step.

The first step constraints are affine of the type

$$C'\theta_C = \psi, \tag{25}$$

for known matrices  $C \in \mathbb{R}^{q \times (q-q_c)}$ ,  $\psi \in \mathbb{R}^{q-q_c}$ . The constrained time effect space is then

$$\Theta_C = \left( \theta_C \in \Theta : C'\theta_C = \psi \right). \tag{26}$$

Thereby the  $q$ -dimensional time effect space  $\Theta$  is reduced to a  $q_C$ -dimensional variation. The properties of this partially ad hoc identified parameter space depends on the rank of the matrix  $(A, C)$ . If the number of constraints,  $q - q_C$ , is at most equal to the number of unidentified components  $q - p$ , it is possible that  $(A, C)$  has full column rank. In that case the constraint implies a partial ad hoc identification without constraining the parameter space  $M$  of the statistical model. This is shown in Theorem 1; see also Section 5.4.1 for an example, while the proof is given in the Appendix. When  $(A, C)$  has reduced rank the parameter space  $M$  is also constrained; see Section 3.4.2 for a discussion.

**Theorem 1.** Suppose  $(A, C)$  has full column rank. Then the matrix  $m = A'_\perp C \in \mathbb{R}^{(q-p) \times (q-q_c)}$  has full column rank and the constraint (25) does not constrain the canonical parameter  $\xi$  and the predictor  $\mu$ . Hence, the predictor space remains of the form (2). The equivalence classes in  $\Theta_C$  under the mapping  $\theta \mapsto \mu = XA'\theta$  are given by the group

$$g_C : \theta \mapsto \theta + A_\perp m_\perp \zeta, \tag{27}$$

for arbitrary  $\zeta \in \mathbb{R}^{q_c-p}$  where  $m_\perp \in \mathbb{R}^{(q-p) \times (q_c-p)}$  is the orthogonal complement of  $m$ . The maximal invariant remains  $\xi = A'\theta$ .

The partial ad hoc identification by (25) implies that any time series analysis of the time effects has to happen

relative to the constrained space  $\Theta_C$  rather than the space  $\Theta$ . This is awkward as discussed in Section 3.4 below. It is also considerably more complicated than working with the freely varying canonical parameter  $\xi$ ; see Section 2.4.2.

**3.4. Consequences of Ad Hoc Identification.** In the following we will look closer at the consequences of working with the ad hoc identified time effect parameter  $\theta$  in the context of a linear mortality model of the form  $\mu = D\theta$ . We consider the consequences for plotting, hypothesis testing, and forecasting.

**3.4.1. Plots of Time Effects.** In the mortality model (1) the time effect  $\theta$  is the concatenation of age, period, and cohort effects. It seems natural to think of these individual time effects as time series and to plot them against time. As the time effect  $\theta$  varies in the unrestricted space  $\Theta = \mathbb{R}^q$  this maps the  $q$ -vector into unrestricted time series.

Estimates of the time effects are constructed by combining an estimate of  $\xi$  with an ad hoc chosen value for  $\lambda = L'\theta$ , see (14). The resulting estimate  $\hat{\theta}_\lambda$  is therefore constrained to the space  $\Theta_\lambda \subset \Theta$ . The interpretation of the estimate  $\hat{\theta}_\lambda$  is therefore different from the interpretation of the original time effect  $\theta$ . Distances on the spaces  $\Theta$  and  $\Theta_\lambda$  are judged differently and the variability of  $\hat{\theta}_\lambda$  is deduced exclusively from  $\hat{\xi}$  through (16). The time series components of  $\hat{\theta}_\lambda$  are now restricted through  $\lambda = L'\theta_\lambda$ . Plots of the  $\hat{\theta}_\lambda$ -time series are therefore interpreted differently from the imagined plots of the original  $\theta$ -time series and from the plots of the maximal invariant parameter  $\xi$  discussed in Section 2.4.2. Indeed, if one were to analyse the estimated  $\hat{\theta}_\lambda$ -time series statistically the linear constraint should be taken into account. This is a bit complicated as illustrated below, but it is the consequence of working with the ad hoc identified parameter  $\theta_\lambda$  rather than the canonical parameter  $\xi$ .

Attempts to give intrinsic meaning to  $\lambda$  will be specific to the index set for the data set at hand. For instance, the requirement that the age effect should be zero on average does not carry over when looking at a subsample or when forecasting. It is not obvious that such an ad hoc identification is any more or less arbitrary than saying that, for instance, the first or the last age effect should have a particular value.

Adding confidence bands to a plot of  $\hat{\theta}_\lambda$  is in itself not difficult. If  $\hat{\xi}$  is asymptotically normal with mean  $\xi$  and variance  $\Sigma$ , then  $\hat{\theta}_\lambda$  is asymptotically normal with mean  $\theta_\lambda$  and variance  $L_\perp(A'L_\perp)^{-1}\Sigma(L'_\perp A)^{-1}L'_\perp$ . This is a normal distribution on the space  $\Theta_\lambda$ . The interpretation of these standard errors will therefore be similar to that of  $\hat{\theta}_\lambda$  itself.

Finally, it may be of interest to analyse the estimated  $\hat{\theta}_\lambda$ -time series statistically. Denote this time series by  $x_\lambda$ . Its sample space is now  $\Theta_\lambda$ . A statistical model on  $\Theta_\lambda$  can be built as follows. The starting point could be a time series model for unrestricted variables  $x$  on the sample space  $\Theta$ . This gives a joint density for  $x \in \Theta$ , which can be reduced by marginalisation to a density for  $x_\lambda \in \Theta_\lambda$ . Whether one is working with the unrestricted model for  $x \in \Theta$  or the

restricted model for  $x_\lambda \in \Theta_\lambda$  inferences that are invariant to  $g$  must be based on those statistics of  $x$  or  $x_\lambda$  that are invariant to  $g$ . Thus, inferences must be based on the maximal invariant under  $g$ . For a general overview of invariant reduction see Cox and Hinkley [12, page 175f], whereas Nielsen [15] gives the argument in some detail for an autoregression with a linear trend.

**3.4.2. Hypothesis Testing.** Having formulated the model in terms of time effects it may be of interest to test the hypothesis that one of these time effects is absent. No identification issues arise when the hypothesis is formulated as a restriction on the canonical parameter  $\xi$  as discussed in Section 2.4.3. But one has to be careful when formulating hypotheses in terms of the original time effect. See Sections 5.4.5, 5.5.3, and 5.5.4 for examples.

Affine hypotheses on the time effect are of the form

$$R'\theta_R = \rho, \tag{28}$$

for known matrices  $R \in \mathbb{R}^{q \times (q - q_R)}$ ,  $\rho \in \mathbb{R}^{q - q_R}$ . The constrained time effect space is then

$$\Theta_R = \{\theta_R \in \mathbb{R}^q : R'\theta_R = \rho\}. \tag{29}$$

To see how the restriction (28) restricts the predictor space  $M \subset \mathbb{R}^n$  recall that the predictor  $\mu$  only depends on  $\theta$  through  $\xi = A'\theta$ . Thus, the analysis of the restriction (28) depends on the interplay between the matrices  $A$ ,  $R$ . Theorem A.3 in Appendix A.3 gives a general result to that effect. It shows that the hypothesis (28) restricts the predictor space  $M$  to a  $p_R$ -dimensional affine subspace of  $\mathbb{R}^n$  in so far as it restricts the canonical parameter  $\xi$ . In particular, the degrees of freedom of the hypothesis,  $p - p_R$ , may in general be different from the dimension reduction of the time effect parameter,  $q - q_R$ . When this is the case the restriction (28) has an element of ad hoc identifying the time effect.

**3.4.3. Forecasts.** Forecasts can be made by extrapolating the ad hoc identified time effects  $\theta_\lambda$ . Two researchers choosing different ad hoc identification schemes, but otherwise making the same analysis, may make different forecasts. This can be avoided if the extrapolation method is chosen with some care.

Following the linear approach outlined in Section 2.4.4 the predictor  $\mu = D\theta = XA'\theta$  is forecasted by extending the design  $D$  into

$$D^h = \begin{pmatrix} D & 0 \\ D_1^h & D_2^h \end{pmatrix}. \tag{30}$$

Extrapolating the ad hoc identified  $\theta_\lambda$  into a vector  $(\theta'_\lambda, \tilde{\theta}'_\lambda)'$  then gives the forecast

$$\tilde{\mu} = (D_1^h, D_2^h) \begin{pmatrix} \theta_\lambda \\ \tilde{\theta}_\lambda \end{pmatrix} = D_1^h \theta_\lambda + D_2^h \tilde{\theta}_\lambda. \tag{31}$$

Often both components  $D_1^h \theta_\lambda$  and  $D_2^h \tilde{\theta}_\lambda$  depend on the ad hoc identification. Nonetheless, these dependencies of



ad hoc identification may cancel each other so that the overall forecast  $\bar{\mu}$  is invariant to the ad hoc identification. Such invariance would seem desirable in most applications unless there is strong substantial reason for the ad hoc identification scheme. Necessary and sufficient conditions for invariance are presented for the age-period-cohort model in Section 5.4.3 and for a nonlinear model in Section 6.1.5.

In contrast, these considerations are redundant when working with the canonical parameter,  $\xi$ ; see Section 2.4.4.

#### 4. Bayesian Models and Random Effects Models

Mortality analysis is often carried out using either Bayesian methods or random effects methods. The mortality model is then altered through the introduction of a prior distribution on the parameters. One might think that the identification problems become less of an issue or even disappear. This is not the case since the Bayesian method and the random effects method is based on the mortality likelihood which only depends on the time effect  $\theta$  through the maximal invariant parameter  $\xi$ . Thus, the identification challenges remain. The issue is that a prior on the unidentified part, say  $\lambda$ , of the time effect amounts to an ad hoc identification. Indeed, the conditional prior of  $\lambda$  given  $\xi$  is not updated by the mortality likelihood. A main difference is that a maximum likelihood analysis of the original mortality likelihood usually prompts the researcher when there is an identification issue, whereas both Bayesian methods and random effects methods allow computations to go through despite an identification issue.

In Section 4.1 it is seen that introduction of a conditional prior on  $\lambda$  given  $\xi$  is the Bayesian analogue of ad hoc identification. This leads to the same type of forecasting challenges as in the frequentist settings as is seen in Section 4.2. In Section 4.3 we show how the Bayesian identification issues transfer to random effects models.

*4.1. Bayesian Estimation.* For Bayesian and random effects models we formulate a likelihood and a prior. Thus, consider a likelihood  $p(y | \theta) = L(\theta; y)$ . Replacing  $\theta$  by  $\xi, \lambda$  the identification problem implies that

$$p(y | \xi, \lambda) = p(y | \xi) \quad \text{for all outcomes } y. \quad (32)$$

The prior on  $\theta$  is factorised as  $p(\theta) = p(\xi, \lambda) = p(\xi)p(\lambda | \xi)$ . In the case of Bayesian estimation the following result emerges.

**Theorem 2.** *Suppose the likelihood satisfies (32). Then*

- (i) *the predictive distribution does not depend on the conditional prior for  $\lambda$ :*

$$p(y) = \int p(y | \xi) p(\xi) d\xi; \quad (33)$$

- (ii) *the posterior satisfies*

$$p(\xi | y) = \frac{p(y | \xi) p(\xi)}{p(y)}, \quad p(\lambda | \xi, y) = p(\lambda | \xi); \quad (34)$$

- (iii) *the posterior means satisfy*

$$E(\xi | y) = \int \xi p(\xi | y) d\xi, \\ E(\lambda | \xi, y) = E(\lambda | \xi), \quad (35)$$

$$E(\lambda | y) = \int E(\lambda | \xi) p(\xi | y) d\xi.$$

Theorem 2 shows that it suffices to give a prior to  $\xi$  and ignore  $\lambda$  as advocated in Section 2.4.5. Indeed the conditional prior for  $\lambda$  given  $\xi$  is not updated. Theorem 2 appears to be well-known; see Poirier [16, Proposition 2] or Smith [17, Section B].

Due to Theorem 2 the Bayesian analyst faces the complications outlined in Section 3.4. Indeed, suppose that two Bayesian researchers choose the same likelihood  $p(x | \xi, \lambda) = p(x | \xi)$  and the same prior  $p(\xi)$  for  $\xi$ , but different conditional priors for  $\lambda$  given  $\xi$ . Their marginal distributions for the data are identical, but any inferences regarding interpretation or forecasting will differ in so far as they involve the unidentified parameter  $\lambda$ . A Bayesian researcher should therefore be cautious with inference related to  $\lambda$ . There will of course be situations where the prior knowledge of  $\lambda$  given  $\xi$  is found to be of substantive relevance. In such situations it seems more fruitful to change the likelihood to include that information.

*4.2. Forecasting.* Bayesian forecasts involve integrating an extrapolative distribution. This can be done in two ways, either working exclusively with the identified, maximal invariant parameter  $\xi$  as in Section 2.4.4, or working with the time effect  $\theta = (\xi, \lambda)$  as in Section 3.4.3.

*4.2.1. Forecasting Using the Maximal Invariant Parameter.* Consider first the case where only the maximal invariant parameter  $\xi$  is used. In that case the forecast is computed by sampling from the posterior  $p(\xi | y)$  and then extrapolating  $\bar{\mu}$  using the sampled value  $\xi$  using some extrapolative methods, say  $p(\bar{\mu} | \xi, y)$ . In combination this gives the forecast

$$p(\bar{\mu} | y) = \int p(\bar{\mu} | \xi, y) p(\xi | y) d\xi. \quad (36)$$

*4.2.2. Forecasting Using the Ad Hoc Identified Time Effect.* Consider now forecasts involving the full time effect  $\theta = (\xi, \lambda)$ . Theorem 2(ii) shows that the posterior satisfies  $p(\theta | y) = p(\xi | y)p(\lambda | \xi)$ . The distribution forecast with extrapolation  $p(\bar{\mu} | \xi, \lambda, y)$  is then

$$p(\bar{\mu} | y) = \iint p(\bar{\mu} | \xi, \lambda, y) p(\xi | y) p(\lambda | \xi) d\lambda d\xi. \quad (37)$$

The concern is now as follows. Suppose a second researcher chooses the same extrapolative method, likelihood, and prior for  $\xi$ , but different conditional priors  $p^\dagger(\lambda | \xi)$ . In general, this will result in a different distribution forecast:

$$p^\dagger(\bar{\mu} | y) = \iint p(\bar{\mu} | \xi, \lambda, y) p(\xi | y) p^\dagger(\lambda | \xi) d\lambda d\xi. \quad (38)$$

The question is then under which conditions will  $p(\bar{\mu} | y) = p^\dagger(\bar{\mu} | y)$  so that the distribution forecasts are invariant to the choice of conditional prior for  $\lambda$  given  $\xi$ ? A sufficient condition is that the extrapolation method does not depend on  $\lambda$  so

$$p(\bar{\mu} | \xi, \lambda, y) = p(\bar{\mu} | \xi, y). \tag{39}$$

Condition (39) could alternatively be expressed as requiring that the forecast  $p(\bar{\mu} | \theta, y) = p(\bar{\mu} | \xi, \lambda, y)$  is invariant to the group  $g$  acting on the time effect space  $\Theta$  so that  $p(\bar{\mu} | \xi, \lambda, y) = p\{\bar{\mu} | \xi, g(\lambda), y\}$ .

**Theorem 3.** *Suppose that the likelihood satisfies (32) and the priors are probabilities. If the extrapolative distribution does not depend on  $\lambda$  so (39) holds; then the forecast distribution  $p(\bar{\mu} | y)$  computed in (37) is invariant to the choice of conditional prior for  $\lambda$  given  $\xi$ . The forecast then reduces to (36).*

To summarise, the identification issues surrounding Bayesian analysis are similar to those outlined in the previous sections. Examples of the problems that can arise are discussed in Sections 5.4.4, 5.5.2, and 6.1.6. There are two solutions to the identification problem. The first is only to formulate a prior on  $\xi$ ; see Section 2.4.5. Incidentally, this is what Bernardo and Smith [18, page 218] do in their discussion of the two-way analysis of variance, albeit without linking it to the considerations of Smith [17]. The prior  $p(\xi)$  can of course be constructed by formulating a prior on  $\theta$  and then reduce it to a prior on  $\xi$  by marginalisation so  $p(\xi) = \int p(\xi, \lambda) d\lambda$ . The other solution is to work with a prior on  $\theta$  but avoid those parts of the posterior that depend on  $\lambda$ .

**4.3. Random Effects Models.** It is common to combine mortality models with a random effects approach, which effectively forms a new model. An example is given in Section 5.4.6. We consider the consequence of the lack of identification.

The random effect models are typically constructed as follows. Suppose the density of the data  $y$  given the time effects  $\theta = (\xi, \lambda)$  is of the form  $p(y | \xi, \lambda) = p(y | \xi)$  as before; see (32). A prior  $p(\theta | \psi)$  is chosen that now depends on a parameter  $\psi$ . The prior can be decomposed as  $p(\theta | \psi) = p(\xi | \psi)p(\lambda | \xi, \psi)$ . Theorem 2 implies that the density of the data  $y$  given  $\psi$  is

$$p(y | \psi) = \int p(y | \xi) p(\xi | \psi) d\xi. \tag{40}$$

This in turn is used to form the random effects likelihood of  $\psi$  as

$$L_{RE}(\psi | y) = p(y | \psi). \tag{41}$$

This, effectively, defines a new model. The random effects likelihood only depends on the prior  $p(\theta | \psi)$  through  $p(\xi | \psi)$ . Two researchers choosing the same prior  $p(\xi | \psi)$  but different conditional priors  $p(\lambda | \xi, \psi)$  will then get the same random effects likelihood and the same maximum likelihood estimator  $\hat{\psi}$ .

In mortality modelling it is common to go one step further and estimate the time effects  $\theta$  through the mean of the posterior  $p(\theta | y, \psi)$  evaluated at  $\psi = \hat{\psi}$ . Then the identification problem may show up. Theorem 2 shows that

$$p(\xi | \hat{\psi}, y) = \frac{p(y | \xi) p(\xi | \hat{\psi})}{p(y | \hat{\psi})}, \tag{42}$$

$$p(\lambda | \xi, \hat{\psi}, y) = p(\lambda | \xi, \hat{\psi}),$$

so that the prior for  $\xi$  is updated, while the conditional posterior for  $\lambda$  given  $\xi$  is not updated by the data. Thus, in general the estimate for  $\theta$  is based, in part, on a prior which is not updated by the data.

### 5. Age-Period-Cohort Models

We will now apply the theoretical considerations to analyse the age-period-cohort model. The methodological literature on this model is large and the consequences of the above theory are wide ranging.

In Section 5.1 we present the age-period-cohort model along with the maximal invariant parameter. This maximal invariant parameter is also called the canonical parameter because the age-period-cohort model is usually implemented as an exponential family; see Section 2.4.1. When formulating the model we choose a notation matching the age-period-cohort literature rather than the reserving literature. At the same time the exposition takes it starting point in Kuang et al. [19], but the notation deviates.

The implementation of the canonical parameter depends on the type of data array. In Section 5.2 design matrices are given for age-cohort, age-period, and period-cohort data arrays. While they illustrate interesting differences in the structure for these data arrays, they also provide the basis for an immediate implementation via any generalised linear model software. The age-cohort model is expressed as a hypothesis of the age-period-cohort model in Section 5.3. Time effects and forecasting are considered in Section 5.4, while the two-sample age-period-cohort model is discussed in Section 5.5.

**5.1. The Model and the Canonical Parameter.** Here the age-period model is set up and a quite general identification result is reported.

Consider data  $Y_{ij}$  indexed by  $(i, j) \in \mathcal{J}$  where  $i$  is the age and  $j$  is the period. The index set may be a rectangle given by  $i = 1, \dots, I$  and  $j = 1, \dots, J$  so that the cohort  $k = I - i + j$  runs from 1 to  $K = I + J - 1$ . More generally, the index set could be a generalized trapezoid where two corners are cut off the rectangle so that the cohort  $k$  runs from  $1 + h_1$  to  $I + J - 1 - h_2$  for some  $h_1, h_2 \geq 0$ . The class of generalized trapezoids includes the three types of Lexis diagrams discussed by Keiding [20]. We will return to those special cases below.

The statistical model is defined by the assumption that the variables  $Y_{ij}$  are independent with an exponential family distribution with predictor  $\mu_{ij}$  given by

$$\mu_{ij} = \alpha_i + \beta_j + \gamma_k + \delta \quad \text{for } i, j \in \mathcal{J}. \tag{43}$$

The time effect  $\theta = (\alpha_1, \dots, \alpha_I, \beta_1, \dots, \beta_j, \gamma_{h_1+1}, \dots, \gamma_{I+J-1-h_2}, \delta)'$  now varies in some time effect space  $\Theta \in \mathbb{R}^q$  where  $q = I + J + K + 1 - h_1 - h_2$ .

The model (43) is of the form (1) discussed in Section 2. Specifically, the predictors  $\mu_{ij}$  can be stacked in a vector  $\mu$  of dimension  $n = \dim \mathcal{S}$  and written as  $\mu = D\theta$ . Thus, the parameter space for the model is of the form  $M = (\mu \in \mathbb{R}^n : \mu = D\theta \text{ for } \theta \in \Theta)$  as outlined in (2). The mapping  $\theta \mapsto \mu$  from  $\Theta$  to  $M$  is surjective and the equivalence classes in the time effect space can be described by a group of transformations that are discussed in (8). This group can be represented as

$$g : \begin{pmatrix} \alpha_i \\ \beta_j \\ \gamma_k \\ \delta \end{pmatrix} \mapsto \begin{pmatrix} \alpha_i + a + (i - 1)d \\ \beta_j + b - (j - 1)d \\ \gamma_k + c + (k - 1)d \\ \delta - a - b - c - (I - 1)d \end{pmatrix} \text{ for } \theta \in \Theta, \tag{44}$$

for any  $a, b, c$ , and  $d$ . This is of the form (8) with  $\zeta = (a, b, c, d)'$  although the definition of the matrix  $A$  depends on the structure of the index set  $\mathcal{S}$ .

A first clue for the canonical parametrisation is given by Fienberg and Mason [21] and Clayton and Schiffler [22] who pointed out that, on the multiplicative scale, ratios of relative risks are invariant. On the additive scale this amounts to looking at second differences, such as  $\Delta^2 \alpha_i = \alpha_i - 2\alpha_{i-1} + \alpha_{i-2}$ . A graphical illustration of the double differences is given in Figure 1 (graphics were done using R 3.0.2, see R Development Core Team [23]), which is taken from Miranda et al. [24]. Panel (a) illustrates the interpretations of the formula for  $\Delta^2 \alpha_i$  as follows. Consider the 1970 and 1971 cohorts. In 2010 these have ages 40 and 39, while in 2011 these have ages 41 and 40. Thus,  $\Delta^2 \alpha_{41}$  represents the increase in mortality from ages 40 to 41 in 2011 relative to the increase from ages 39 to age 40 in 2010. An equivalent interpretation is that which represents the increase in mortality from ages 40 to 41 for the 1970 cohort relative to the increase from ages 39 to 40 for the 1971 cohort. In a similar way panels (b) and (c) illustrate the formulas for  $\Delta^2 \beta_{2012}$  and  $\Delta^2 \gamma_{1972}$ .

Kuang et al. [19] introduces a parameter formed by these second differences as well as three entries of the predictor; that is,

$$\xi = (\mu_{i_1 j_1}, \mu_{i_2 j_2}, \mu_{i_3 j_3}, \Delta^2 \alpha_3, \dots, \Delta^2 \alpha_I, \Delta^2 \beta_3, \dots, \Delta^2 \beta_J, \Delta^2 \gamma_{h_1+3}, \dots, \Delta^2 \gamma_{K-h_2}). \tag{45}$$

The parameter  $\xi$  varies in the space  $\Xi = \mathbb{R}^p$  where  $p = q - 4$ . If the three points  $\mu_{i_1 j_1}$ ,  $\mu_{i_2 j_2}$ , and  $\mu_{i_3 j_3}$  are chosen not to be linearly related then they define the levels and the linear trends in the predictor. The formal condition is that a certain determinant defined from the indices is nonzero; that is,

$$i_2 j_3 - i_3 j_2 + i_3 j_1 - i_1 j_3 + i_1 k_2 - i_2 k_1 \neq 0. \tag{46}$$

**Theorem 4** (see [19], [25, Corollary 2]). *Let  $\mu$  satisfy (43). If the condition (46) is satisfied then the parameter  $\xi$  of (45) satisfies the following:*

- (i)  $\xi$  is a function of  $\theta$  which is invariant to the group  $g$  in (44);
- (ii)  $\mu$  is a function of  $\xi$ ;
- (iii) the parametrisation of  $\mu$  by  $\xi$  is exactly identified in that  $\xi^\dagger \neq \xi^\ddagger \Rightarrow \mu(\xi^\dagger) \neq \mu(\xi^\ddagger)$ .

Theorem 4 therefore shows that  $\xi$  varies freely in  $\Xi = \mathbb{R}^p$ . Moreover,  $\xi$  is a maximal invariant of the mapping  $m$  from  $\theta$  to  $\mu$  under the transformations  $g$ . It should be noted that the choice of maximal invariant is not unique. Indeed, any bijective mapping of  $\xi$  can serve as maximal invariant. The choice of  $\xi$  is convenient since it becomes the canonical parameter in generalized linear models of the exponential family type.

In itself this theorem does not tell how to express the predictor  $\mu$  in terms of the canonical parameter  $\xi$ . The link depends on the structure of the index set  $\mathcal{S}$ . The above mentioned paper gives implicit expressions for generalized trapezoid index sets. In the following we report explicit expressions for the 3 principal Lexis diagrams.

**5.2. Design Matrices for Lexis Diagrams.** The link between the canonical parameter  $\xi$  and the predictor  $\mu$  is analysed for the 3 principal Lexis diagrams. We start with age-cohort data arrays, which were the focus of attention in Kuang et al. [19]. Such arrays are easiest to analyse because all three time scales increase from the point where  $i = j = k = 1$ . As a consequence the results are relatively easier for these arrays.

**5.2.1. Age-Cohort Data Arrays.** Age-cohort data arrays are rectangular in the age and cohort indices and given by

$$\mathcal{S}_{ac} = \{(i, k) : i = 1, \dots, I, k = 1, \dots, K\}. \tag{47}$$

Consequently, the period index  $j = i + k - 1$  varies over  $j = 1, \dots, J = I + K - 1$ . Keiding [20] refers to this Lexis diagram as the first principal set of death.

Age-cohort arrays are in particular used for reserving in general insurance. In that situation, only the triangle  $1 \leq i, j, k \leq I$  is observed. The issue is to forecast the other triangle in the square  $1 \leq i, k \leq I$ . In the reserving literature these triangles are referred to as the upper and lower triangles, since the cohort axis has reverse order. The two-factor age-cohort model for triangular age-cohort arrays is known as the chain-ladder model; see England and Verrall [26] for an overview. Zehnwirth [27] introduced an age-period-cohort model for such triangular arrays. The identification issue is analysed in detail in Kuang et al. [19, 25]. Subsequently, Kuang et al. [28] analysed the Poisson likelihood, while Kuang et al. [10] give an empirical analysis focusing on forecasting.

The age-period-cohort model for the age-cohort arrays is parametrised by

$$\mu_{ik} = \alpha_i + \beta_{i+k-1} + \gamma_k + \delta \text{ for } i, k \in \mathcal{S}_{ac}. \tag{48}$$

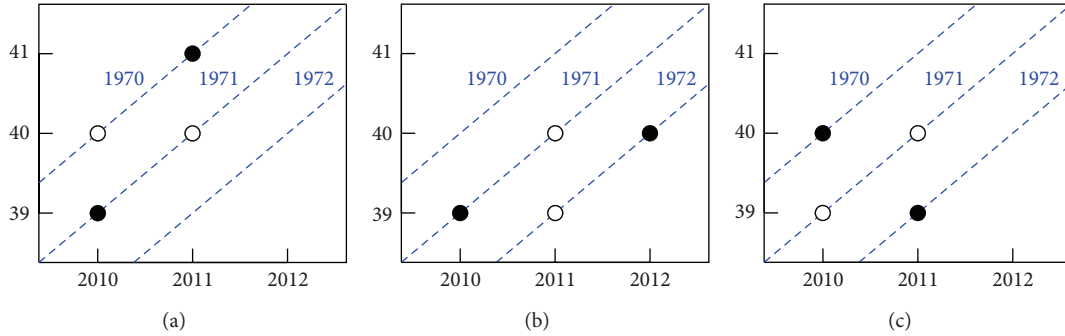


FIGURE 1: Illustration of interpretation of  $\Delta^2\alpha_{41}$ ,  $\Delta^2\beta_{2012}$ , and  $\Delta^2\gamma_{1972}$ .

The time effect  $\theta = (\alpha_1, \dots, \alpha_I, \beta_1, \dots, \beta_J, \gamma_1, \dots, \gamma_K, \delta)'$  now varies in  $\Theta = \mathbb{R}^{2(I+K)}$ .

The design matrix linking the canonical parameter  $\xi$  in (45) and the predictor  $\mu$  is essentially an identity linking the two parameters. A natural choice of the three levels points to the predictors that are  $\mu_{11}$ ,  $\mu_{12}$ , and  $\mu_{21}$ . We then get the representation

$$\begin{aligned} \mu_{ik} &= \mu_{11} + (i - 1)(\mu_{21} - \mu_{11}) + (k - 1)(\mu_{12} - \mu_{11}) \\ &+ \sum_{\ell=3}^i \sum_{h=3}^{\ell} \Delta^2\alpha_h + \sum_{\ell=3}^j \sum_{h=3}^{\ell} \Delta^2\beta_h + \sum_{\ell=3}^k \sum_{h=3}^{\ell} \Delta^2\gamma_h, \end{aligned} \quad (49)$$

with the convention that empty sums are zero, and recalling that second differences are defined as  $\Delta^2\alpha_i = \alpha_i - 2\alpha_{i-1} + \alpha_{i-2}$  so that  $\sum_{h=3}^i \Delta^2\alpha_h = \Delta\alpha_i - \Delta\alpha_2$  and  $\sum_{\ell=3}^i \sum_{h=3}^{\ell} \Delta^2\alpha_h = \alpha_i - \alpha_1 - (i - 1)\Delta\alpha_2$ .

The identity (49) is crucial to the understanding of the age-period-cohort model. It shows that the predictor has a single level expressed as  $\mu_{11}$ , which in turn satisfies  $\mu_{11} = \alpha_1 + \beta_1 + \gamma_1 + \delta$ . The level  $\mu_{11}$  is therefore estimable, but the individual levels  $\alpha_1$ ,  $\beta_1$ ,  $\gamma_1$ , and  $\delta$  are not identifiable from the model. Further, the model has two linear trends, here expressed with slopes  $\mu_{21} - \mu_{11}$  and  $\mu_{12} - \mu_{11}$  in terms of the age and cohort indices. These slopes can be expressed as  $\mu_{21} - \mu_{11} = \Delta\alpha_2 + \Delta\beta_2$  and  $\mu_{12} - \mu_{11} = \Delta\beta_2 + \Delta\gamma_2$ . They are estimable, but the individual slopes  $\Delta\alpha_2$ ,  $\Delta\beta_2$ , and  $\Delta\gamma_2$  are not identifiable.

The design matrix now follows from the identity (49) so that the predictor satisfies  $\mu = X\xi$ , where

$$\begin{aligned} \xi &= (\mu_{11}, \mu_{21} - \mu_{11}, \mu_{12} - \mu_{11}, \Delta^2\alpha_3, \dots, \Delta^2\alpha_I, \\ &\Delta^2\beta_3, \dots, \Delta^2\beta_J, \Delta^2\gamma_3, \dots, \Delta^2\gamma_K)' \end{aligned} \quad (50)$$

$$\begin{aligned} X_{ik} &= \{1, (i - 1), (k - 1), h(i, 3), \dots, h(i, I), \\ &h(j, 3), \dots, h(j, J), h(k, 3), \dots, h(k, K)\}' \end{aligned} \quad (51)$$

where  $\xi \in \mathbb{R}^p$ , where  $p = 2(I + K - 2)$  and  $h(t, s) = \max(t - s + 1, 0)$ .

The identification relies on Theorem 4, which can be specialised to age-cohort arrays as follows.

**Theorem 5** (see [19, Theorem 1]). *Let  $\mu$  satisfy (48). The parameter  $\xi$  of (50) satisfies the following:*

- (i)  $\xi$  is a function of  $\theta$  which is invariant to the group  $g$  in (44);
- (ii)  $\mu$  is a function of  $\xi$ , because of (49);
- (iii) the parametrisation of  $\mu$  by  $\xi$  is exactly identified in that  $\xi^\dagger \neq \xi^\ddagger \Rightarrow \mu(\xi^\dagger) \neq \mu(\xi^\ddagger)$ .

Theorem 5 in turn implies that the parameter  $\xi$  varies freely in  $\Xi = \mathbb{R}^p$ , while the design matrix  $X$  given by (51) has full column rank. Originally, the more general Theorem 4 was proved as a corollary to Theorem 5.

**5.2.2. Age-Period Arrays.** An age-period data array is rectangular in the age and cohort indices and given by

$$\mathcal{F}_{ap} = \{(i, j) : i = 1, \dots, I, j = 1, \dots, J\}. \quad (52)$$

Consequently, the cohort index  $k = j - i + I$  varies over  $k = 1, \dots, K = I + J - 1$ . Keiding [20] refers to this Lexis diagram as the third principal set of death.

Age-period arrays are commonly used in epidemiology, in mortality analysis, and in sociology. The analysis of identification issue is largely similar to that of age-cohort arrays. However, the representation of the predictor  $\mu$  in terms of  $\xi$  differs in an intriguing way, because the third time index, the cohort  $k$ , is the difference of the other two indices.

The age-period-cohort model for the age-period arrays is parametrised by

$$\mu_{ij} = \alpha_i + \beta_j + \gamma_{j-i+I} + \delta \quad \text{for } i, j \in \mathcal{F}_{ap}. \quad (53)$$

The time effect  $\theta = (\alpha_1, \dots, \alpha_I, \beta_1, \dots, \beta_J, \gamma_1, \dots, \gamma_K, \delta)'$  now varies in  $\Theta = \mathbb{R}^{2(I+J)}$ . A representation of the predictor  $\mu$  in terms of the canonical parameter  $\xi$  is now

$$\begin{aligned} \mu_{ij} &= \mu_{11} + (i - I)(\mu_{11} - \mu_{I-1,1}) + (j - 1)(\mu_{12} - \mu_{11}) \\ &+ \sum_{\ell=i}^{I-2} \sum_{h=\ell}^{I-2} \Delta^2\alpha_h + \sum_{\ell=3}^j \sum_{h=3}^{\ell} \Delta^2\beta_h \\ &+ \sum_{\ell=3}^{j-i+I} \sum_{h=3}^{\ell} \Delta^2\gamma_{h+2}. \end{aligned} \quad (54)$$

The representation (54) differs from that of (49) in a subtle way. The three reference points for the levels of the predictor are chosen in the corner  $i = I, j = 1$ . From this corner period and cohort indices increase, while age decreases. Hence, the age double differences  $\Delta^2\alpha_i$  are now cumulated backwards. This phenomenon arises because the cohort index is the difference of the principal indices of age and period, whereas for the age-cohort array the period index is the sum of the principal indices of age and cohort. The predictor is now  $\mu = X\xi$  where, with  $h(t, s) = \max(t - s + 1, 0)$ ,

$$\begin{aligned} \xi &= (\mu_{I1}, \mu_{I1} - \mu_{I-1,1}, \mu_{I2} - \mu_{I1}, \Delta^2\alpha_3, \dots, \Delta^2\alpha_I, \\ &\quad \Delta^2\beta_3, \dots, \Delta^2\beta_J, \Delta^2\gamma_3, \dots, \Delta^2\gamma_K)' \end{aligned} \tag{55}$$

$$X_{ij} = \{1, i - I, j - 1, h(1, i), \dots, h(I - 2, i),$$

$$(j, 3), \dots, h(j, J), h(j - i + I, 3), \dots,$$

$$h(j - i + I, K)\}' \tag{56}$$

The identification relies on Theorem 4. It is specialised to age-period arrays as follows.

**Theorem 6** (see [24, Theorem 4.1]). *Let  $\mu$  satisfy (53). The parameter  $\xi$  of (55) satisfies the following:*

- (i)  $\xi$  is a function of  $\theta$  which is invariant to the group  $g$  in (44);
- (ii)  $\mu$  is a function of  $\xi$ , because of (54);
- (iii) the parametrisation of  $\mu$  by  $\xi$  is exactly identified in that  $\xi^\dagger \neq \xi^\ddagger \Rightarrow \mu(\xi^\dagger) \neq \mu(\xi^\ddagger)$ .

The group of transformations in (44) can be specialised as

$$g : \begin{pmatrix} \alpha_i \\ \beta_j \\ \gamma_{i-i+I} \\ \delta \end{pmatrix} \mapsto \begin{pmatrix} \alpha_i + a + di \\ \beta_j + b - dj \\ \gamma_{j-i+I} + c + d(j - i + I) \\ \delta - a - b - c - dI \end{pmatrix} \tag{57}$$

for  $\theta \in \Theta = \mathbb{R}^{2(I+J)}$ ;

see, for instance, Carstensen [29]. This is of the form (8) with  $\zeta = (a, b, c, d)'$  and

$$A'_\perp = \begin{pmatrix} 1 & 1 & \dots & 1 & & & & & & & -1 \\ & & & & 1 & 1 & \dots & 1 & & & -1 \\ 1 & 2 & \dots & I & -1 & -2 & \dots & -J & 1 & 1 & \dots & 1 & -1 \\ & & & & & & & & & & & & -1 \end{pmatrix} \tag{58}$$

5.2.3. *Period-Cohort Arrays.* An age-period data arrays is rectangular in the age and cohort indices and given by

$$\mathcal{F}_{pc} = \{(j, k) : j = 1, \dots, J, k = 1, \dots, K\}. \tag{59}$$

Consequently, the age index  $i = j - k + K$  varies over  $i = 1, \dots, I = J + K - 1$ . Keiding [20] refers to this Lexis diagram as the second principal set of death. Age-period arrays are commonly used in prospective cohort studies in

epidemiology and in sociology. The analysis is similar to that of age-period arrays when swapping the role of age and cohort.

The age-period-cohort model for the age-cohort arrays is parametrised by

$$\mu_{jk} = \alpha_{j-k+1} + \beta_j + \gamma_k + \delta \quad \text{for } j, k \in \mathcal{F}_{ap}. \tag{60}$$

The time effect  $\theta = (\alpha_1, \dots, \alpha_I, \beta_1, \dots, \beta_J, \gamma_1, \dots, \gamma_K, \delta)'$  now varies in  $\Theta = \mathbb{R}^{2(J+K)}$ . A representation of the predictor  $\mu$  in terms of the canonical parameter  $\xi$  is now

$$\begin{aligned} \mu_{jk} &= \mu_{1K} + (j - 1)(\mu_{2K} - \mu_{1K}) + (k - K)(\mu_{1K} - \mu_{1,K-1}) \\ &+ \sum_{\ell=3}^{j-k+1} \sum_{h=3}^{\ell} \Delta^2\alpha_h + \sum_{\ell=3}^j \sum_{h=3}^{\ell} \Delta^2\beta_h \\ &+ \sum_{\ell=i}^{K-2} \sum_{h=\ell}^{K-2} \Delta^2\gamma_{h+2}. \end{aligned} \tag{61}$$

Thus, the canonical parameter and the design matrix are given by

$$\begin{aligned} \xi &= (\mu_{1K}, \mu_{2K} - \mu_{1K}, \mu_{1K} - \mu_{1,K-1}, \Delta^2\alpha_3, \dots, \Delta^2\alpha_I, \\ &\quad \Delta^2\beta_3, \dots, \Delta^2\beta_J, \Delta^2\gamma_3, \dots, \Delta^2\gamma_K)' \end{aligned} \tag{62}$$

$$X_{jk} = \{1, j - 1, k - K, h(j - k + 1, 3), \dots, h(j - k + 1, I),$$

$$h(j, 3), \dots, h(j, J), h(1, k), \dots, h(K - 2, k)\}' \tag{63}$$

In parallel with Theorem 6 we then have the following identification result.

**Theorem 7.** *Let  $\mu$  satisfy (60). The parameter  $\xi$  of (62) satisfies the following:*

- (i)  $\xi$  is a function of  $\theta$  which is invariant to the group  $g$  in (44);
- (ii)  $\mu$  is a function of  $\xi$ , because of (61);
- (iii) the parametrisation of  $\mu$  by  $\xi$  is exactly identified in that  $\xi^\dagger \neq \xi^\ddagger \Rightarrow \mu(\xi^\dagger) \neq \mu(\xi^\ddagger)$ .

5.3. *Expressing the Age-Cohort Model as a Hypothesis.* It is often of interest to test the absence of the period effect. An application to analysing asbestos related mortality can be found in Miranda et al. [24].

The hypothesis is that  $\beta_1 = \dots = \beta_J$ , when expressed in terms of the time effect parameters. The restricted model is given by, with  $k = j - i + I$ ,

$$\mu_{ij}^{ac} = \alpha_i + \gamma_{j-i+I} + \delta \quad \text{for } i, j \in \mathcal{F}_{ap}. \tag{64}$$

The identification problem simplifies to a question of determining the levels of  $\alpha_i$  and  $\gamma_k$ . Therefore the (log) relative risk parameters  $\Delta\alpha_i$  are identified as pointed out by Clayton

and Schifflers [30]. In this model the cohort index is present and keeps the difference of the principal age and period indices. Therefore the representation of the predictor involves backward cumulated age differences as before but with a subtle change of sign, so that (54) reduces to

$$\mu_{ij}^{ac} = \mu_{I1} - \sum_{\ell=i}^{I-1} \Delta\alpha_{\ell+1} + \sum_{\ell=2}^k \Delta\gamma_{\ell}. \quad (65)$$

As a consequence the canonical parameter and the design reduce to  $\mu_{ij}^{ac} = X_{ij}^{ac} \xi_{ij}^{ac}$ , where

$$X_{ij}^{ac} = \{1, -1_{(1 \geq i)}, \dots, -1_{(I-1 \geq i)}, 1_{(k \geq 2)}, \dots, 1_{(k \geq K)}\}, \quad (66)$$

$$\xi^{ac} = (\mu_{I1}, \Delta\alpha_2, \dots, \Delta\alpha_I, \Delta\gamma_2, \dots, \Delta\gamma_K)'. \quad (67)$$

Miranda et al. [24, Theorem 4.2] establish an identification result similar to Theorem 6.

The age-cohort model can also be formulated as a hypothesis on the maximal invariant  $\xi$  in the age-period-cohort model following Section 2.4.3. The period effects  $\Delta^2\beta_j$  are set to zero through  $H'\xi = 0$ , where  $H' = (0, I_{J-2}, 0)$ . Applying this to the expression for  $\xi$  in (55) gives

$$\xi_H = \overline{H}'_1 \xi = (\mu_{I1}, \Delta\alpha_1 - \Delta\gamma_2, \Delta\gamma_2, \Delta^2\alpha_3, \dots, \Delta^2\alpha_I, \Delta^2\gamma_3, \dots, \Delta^2\gamma_K), \quad (68)$$

since in the absence of period effects; then  $\mu_{I1} - \mu_{I-1,1} = \Delta\alpha_1 - \Delta\gamma_2$  and  $\mu_{I2} - \mu_{I1} = \Delta\gamma_2$ . The double differences cumulate to first differences through  $\sum_{i=3}^I \Delta^2\alpha_i = \Delta\alpha_1 - \Delta\alpha_2$ , so the above expression  $\xi_H$  is seen to be a linear transformation of  $\xi^{ac}$  in (67). In other words the age-cohort model arises from the age-period-cohort model by restricting the maximal invariant parameter.

**5.4. Working with the Time Effect.** There is a large literature seeking to identify the original time effects  $\alpha_i, \beta_j$ , and  $\gamma_k$  of the age-period-cohort model from the predictor. Here we look closer at some of those ad hoc identification proposals.

**5.4.1. Ad Hoc Identification of Levels.** For the age-period-cohort model it is popular to impose ad hoc identifications in two steps of the type discussed in Section 3.3. Here the first step is concerned with the level of the time effects and the second step is concerned with the linear trend. Examples are given in Sections 5.4.2 and 5.5.4.

A common first step ad hoc identification is to require that

$$\sum_{i=1}^I \alpha_i = \sum_{j=1}^J \beta_j = \sum_{k=1}^K \gamma_k = 0. \quad (69)$$

This ad hoc identification is specific to the chosen data range. For instance, the constraint  $\sum_{i=1}^I \alpha_i = 0$  is not easily transferable to a different data set drawn from the same population but with a different set of age groups. This aspect would have to be kept in mind if a substantive motivation was

to be found for this constraint. Other ad hoc identification schemes such as  $\alpha_I = \beta_J = \gamma_K = 0$  have similar problems.

The constraint (69) is a special case of affine constraints of the form  $C'\theta_C = \psi$  discussed in Section 3.3. The involved dimensions are  $q = 2(I + J)$  and  $p = q - 4$ , while the number of constrains is  $q - q_C = 3$ . The matrix  $C' \in \mathbb{R}^{(q-q_C) \times q}$  is given by the top left  $\{3 \times (q - 1)\}$ -block of  $A'_\perp$  in (58) padded with a column of zeros, while  $\psi \in \mathbb{R}^3$  is given by  $\psi = 0$ . Theorem 1 shows that  $m = A'_\perp C \in \mathbb{R}^{(q-p) \times (q-q_C)}$  has full rank. Indeed,  $m$  and its orthogonal complement are given by

$$m = A'_\perp C = \begin{pmatrix} I & 0 & 0 \\ 0 & J & 0 \\ 0 & 0 & K \\ I\check{i} & -J\check{j} & K\check{k} \end{pmatrix}, \quad m_\perp = \begin{pmatrix} -\check{i} \\ \check{j} \\ -\check{k} \\ 1 \end{pmatrix}, \quad (70)$$

where, for instance,  $\check{i} = I^{-1} \sum_{i=1}^I i = (I + 1)/2$ . Thus, the constrained group of equivalence classes (27) is

$$g_C : \begin{pmatrix} \alpha_i \\ \beta_j \\ \gamma_k \\ \delta \end{pmatrix} \mapsto \begin{Bmatrix} \alpha_i + d(i - \check{i}) \\ \beta_j - d(j - \check{j}) \\ \gamma_k + d(k - \check{k}) \\ \delta \end{Bmatrix} \quad (71)$$

for  $\theta \in \Theta_C$ .

**5.4.2. Ad Hoc Identification of Slopes: The “Intrinsic” Estimator.** The “intrinsic” estimator is a popular estimator in the sociology literature; see Yang et al. [4] and see also O’Brien [31, 32] and Fu et al. [33] for a recent discussion of its merits. It has its roots in a suggestion by Kupper et al. [34], with an early critique given by Holford [35].

The “intrinsic” estimator is defined in two steps. In the first step, the levels are identified by the ad hoc constraint (69). Three of the  $\theta$ -coordinates are then dropped; that is  $\alpha_I, \beta_J$ , and  $\gamma_K$  are dropped. In a second step the linear trend is ad hoc identified using a Moore-Penrose inverse as in (23).

We can analyse these steps using the developed framework. The first step identifies the levels by the ad hoc constraint (69), which is a constraint of the form  $C'\theta = 0$  for the  $C$  discussed in Section 5.4.1. This  $\theta$  is defined on  $\Theta_C$  which is a linear subspace with a dimension deficiency of 3. Introduce a selection matrix  $S_\perp \in \mathbb{R}^{q \times (q-3)}$  that selects all coordinates of  $\theta$  except  $\alpha_I, \beta_J$ , and  $\gamma_K$ . Thus  $S_\perp$  arises as a  $q$ -dimensional with 3 columns deleted corresponding to  $\alpha_I, \beta_J$ , and  $\gamma_K$ . This is chosen so that  $(C, S_\perp)$  is invertible. Then  $S'_\perp \theta$  is freely varying in that  $S'_\perp \Theta_C = \mathbb{R}^{q-3}$ . The skew projection identity  $I_q = S(C'S)^{-1}C' + C_\perp(S'_\perp C_\perp)^{-1}S'_\perp$  and the constraint  $C'\theta = 0$  then implies that  $\theta = C_S \vartheta$  where  $C_{S_\perp} = C_\perp(S'_\perp C_\perp)^{-1}$  and  $\vartheta = S'_\perp \theta \in \mathbb{R}^{q-3}$ . Note that while  $C_{S_\perp}$  depends on  $S_\perp$  and  $C_\perp$ , it does not depend on the normalisation of  $C_\perp$ , since we can replace  $C_\perp$  by  $C_\perp m$  for arbitrary invertible matrices  $m \in \mathbb{R}^{(q-3) \times (q-3)}$ . This implies that  $C_S$  is a function of  $S_\perp$  and  $C$ . The predictor  $\mu$  is now parametrised by  $\mu = XA'\theta = XA'_C \vartheta$  with  $A'_C = A'_C S$ . This corresponds to equation 5 of Yang et al. [4] who use the notation  $X$  and  $b$  for  $XA'_C$  and  $\vartheta$ , respectively.

In the second step the linear trend is ad hoc identified through a time effect parameter of the form (23) with  $A, \theta$  replaced by  $A_C, \vartheta$  so that  $\theta_{ad.hoc} = C_{\perp} \vartheta_{ad.hoc}$  where  $\vartheta_{ad.hoc} = L_{\perp} (A_C' L_{\perp})^{-1} \xi + (A_C)_{\perp} \{L' (A_C)_{\perp}\}^{-1} \lambda$  for some scalar  $\lambda$  and some matrix  $L_{\perp} \in \mathbb{R}^{(q-3) \times p}$ .

The “intrinsic” estimator is ad hoc identified through the choices  $\lambda = 0$  and  $L_{\perp} = A_C$ , while  $C$  is chosen by (69). It therefore estimates an “intrinsic” parameter:

$$\theta_{intrinsic} = C_{S_{\perp}} C'_{S_{\perp}} A (A' C_{S_{\perp}} C'_{S_{\perp}} A)^{-1} \xi, \tag{72}$$

which depends on the choices of  $S'_{\perp}, C$ , and  $A_{\perp}$ . However, since we can replace  $C_{\perp}$  by  $C_{\perp} m$  for arbitrary invertible matrices  $m \in \mathbb{R}^{(q-3) \times (q-3)}$  without changing  $\theta_{intrinsic}$  the expression  $\theta_{intrinsic}$  does not depend on the normalisation of  $C_{\perp}$ . The “intrinsic” parameter satisfies the following result.

**Theorem 8.** *The “intrinsic” parameter is an injective mapping of the canonical parameter  $\xi \in \mathbb{R}^p$  into a  $p = q - 4$  dimensional linear subspace  $\Theta_{intrinsic}$  of  $\Theta = \mathbb{R}^q$ . The “intrinsic” time effect space is a  $p$ -dimensional linear subspace of  $\mathbb{R}^q$  of the form*

$$\begin{aligned} \Theta_{intrinsic} &= \left\{ \theta \in \mathbb{R}^q : \theta = C_{S_{\perp}} C'_{S_{\perp}} A (A' C_{S_{\perp}} C'_{S_{\perp}} A)^{-1} \xi \text{ for } \xi \in \mathbb{R}^p \right\}, \\ &= \left\{ \theta \in \mathbb{R}^q : C' \theta = 0, w' (C'_{\perp} S_{\perp}) (S'_{\perp} C_{\perp}) \bar{C}'_{\perp} \theta = 0 \right\}, \end{aligned} \tag{73}$$

where  $w \in \mathbb{R}^{q-3}$  is uniquely defined up to a scale by  $w' C'_{\perp} A = 0$ .

Theorem 8 implies that the “intrinsic” parameter should be interpreted as an object varying in the linear subspace  $\Theta_{intrinsic}$  rather than in the unrestricted time effect space  $\Theta = \mathbb{R}^q$ . As outlined in Section 3.4 this has consequences for the interpretation of plots of the time effects, hypothesis testing, and forecasts. A consequence of this argument is that different choices of  $C, S_{\perp}, L$ , and  $\lambda$  would lead to other ad hoc identified parameters varying in other affine subspaces of  $\Theta$ . In other words, the “intrinsic” estimator carries the cost of working with the somewhat complicated linear subspace  $\Theta_{intrinsic}$ . This effort may be worthwhile if the particular choice of  $C, S_{\perp}, L$ , and  $\lambda$  can be made on substantive grounds.

**5.4.3. Forecasting.** Forecasting of future mortality rates involves an extrapolation of the time parameters. In Section 2.4.4 it was argued that ad hoc identification may introduce an undesirable arbitrariness in the forecast. When working exclusively with the canonical parameter  $\xi$  this arbitrariness is avoided. It is, however, also possible to work with ad hoc identified time effects under specific circumstances that we characterise here for age-period arrays. This builds on the theory developed in Kuang et al. [25] for age-cohort data arrays.

In the context of an age-period data array  $\mathcal{F}_{ap}$  it is often of interest to forecast  $h$  periods ahead. Suppose it is of interest to forecast the mortality at age  $i$  in period  $J + h$ , so that the

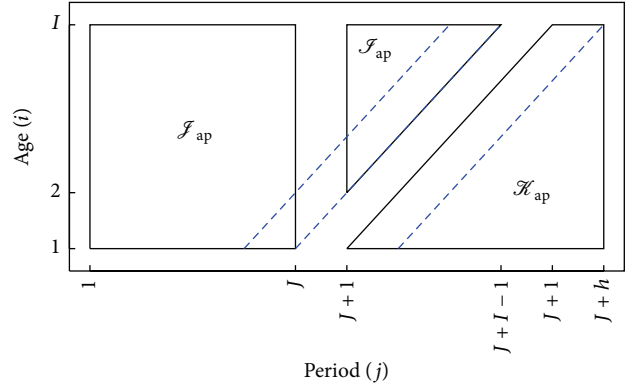


FIGURE 2:  $\mathcal{F}_{ap}$  is the data array.  $\mathcal{F}_{ap,1}$  is the forecast array where only period parameters need to be extrapolated.  $\mathcal{F}_2$  is the forecast array where both period and cohort parameters need to be extrapolated. Cohorts are indicated by dashed lines.

cohort is  $k = I + J + h - i$ . This requires an extrapolation of the period effect. If the cohort index is sufficiently large, that is,  $k > K$ , then the cohort effect needs to be extrapolated too. Thus, there are two forecast index arrays of interest:

$$\begin{aligned} \mathcal{F}_{ap} &= \{(i, j) : i = 1, \dots, I; j = J + 1, \dots, J + h; k \leq K\}, \\ \mathcal{H}_{ap} &= \{(i, j) : i = 1, \dots, I; j = J + 1, \dots, J + h; k > K\}. \end{aligned} \tag{74}$$

Figure 2 illustrates these forecast index arrays.

Identification plays a role when extrapolating the estimates obtained on the data array  $\mathcal{F}_{ap}$ . The identification issues can be ignored if the investigator simply extrapolates  $\Delta^2 \beta_j$  and  $\Delta^2 \gamma_k$ . In the context of ad hoc identified time effects arbitrary linear trends are introduced in the model. The forecast of the predictor  $\mu_{i,J+h}$  is invariant to these if and only if the chosen extrapolation method for  $\beta_j, \gamma_k$  preserves these linear trends so that they can cancel with the arbitrary linear trend in  $\alpha_i$ . The next result gives a precise formulation of this statement. It applies both to point forecasts and distribution forecasts.

**Theorem 9.** *Consider the predictor  $\mu_{ij}$  for  $i, j \in \mathcal{F}_{ap}$  as given in (53). Suppose the time effects  $\alpha_i, \beta_j$ , and  $\gamma_k$  are ad hoc identified. Consider the class of  $h$  periods-ahead forecasts over  $\mathcal{F}_{ap}$  constructed as  $\tilde{\mu}_{i,J+h} = \tilde{\alpha}_i + \tilde{\beta}_{J+h} + \tilde{\gamma}_{I+J+h-i} + \tilde{\delta}$ , where  $\tilde{\beta}_{J+h} + \tilde{\gamma}_{I+J+h-i}$  is a function of the ad hoc identified estimate  $\hat{\theta}$ . Let  $g$  be the group (57). Invariance of the forecast  $\tilde{\mu}_{i,J+h}$  with respect to the ad hoc identification is equivalent to either of the following:*

- (i) the extrapolation method for period and cohort effects is linear trend-preserving:

$$\begin{aligned} &\tilde{\beta}_{J+h}(\hat{\theta}) + \tilde{\gamma}_{I+J+h-i}(\hat{\theta}) \\ &= [\tilde{\beta}_{J+h} \{g(\hat{\theta})\} - b + d(J + h)] \\ &\quad + [\tilde{\gamma}_{I+J+h-i} \{g(\hat{\theta})\} - c - d(I + J + h - i)] \end{aligned} \tag{75}$$

$$\forall b, c, d \in \mathbb{R};$$

(ii) functions  $f_\beta, f_\gamma$  exist so that with  $\hat{\xi}_{\beta,\gamma} = (\Delta^2 \hat{\beta}_3, \dots, \Delta^2 \hat{\beta}_J, \Delta^2 \hat{\gamma}_3, \dots, \Delta^2 \hat{\gamma}_K)'$ ; then

$$\begin{aligned} \tilde{\beta}_{J+h}(\hat{\theta}) + \tilde{\gamma}_{I+J+h-i}(\hat{\theta}) &= \{\hat{\beta}_J + h\Delta\hat{\beta}_J + f_\beta(\hat{\xi})\} \\ &+ \{\hat{\gamma}_K + (h-i+1)\Delta\hat{\gamma}_K + f_\gamma(\hat{\xi})\}. \end{aligned} \tag{76}$$

To illustrate the use of Theorem 9 consider the extrapolation methods  $\tilde{\beta}_{J+h} = \hat{\beta}_J$  and  $\Delta\tilde{\beta}_{J+h} = \Delta\hat{\beta}_J$ . The first forecast is a random walk forecast and it is seen to violate (ii). The second forecast is a cumulated random walk and satisfies (ii). The reason is that  $\beta_{J+h} = \beta_J + \sum_{\ell=1}^h \Delta\beta_{J+\ell}$ . Since  $\Delta\tilde{\beta}_{J+\ell} = \Delta\hat{\beta}_J$ , then  $\tilde{\beta}_{J+h} = \hat{\beta}_J + h\Delta\hat{\beta}_J$ . Further examples of forecasts that are linear trend-preserving as well as some which are not are given Kuang et al. [25, Table 1].

Kuang, Nielsen, and Nielsen [10] apply this to reserving data organised in an age-cohort array  $\mathcal{S}_{ac}$  and discuss the issue of robustification of forecast with respect to structural breaks at the forecast origin. Miranda et al. [24] give an application to asbestos related mortality using an age-period array  $\mathcal{S}_{ap}$ .

**5.4.4. Bayesian Ad Hoc Identification Using a Dynamic Prior.** A Bayesian ad hoc identification using a dynamic prior does not solve the identification problem as discussed in Section 4 and the same care has to be exercised to avoid the problems outlined in Section 3.4. Berzuini and Clayton [6] suggest such an ad hoc identification approach. On page 831 they write “*Identifiability problems may be solved by imposing an arbitrary linear constraint on the log-linear trend components of age, period and cohort effects. Happily, such an arbitrary constraint has no effect on the predictions of the model.*” The previous analysis suggests that this is far from innocent.

The Berzuini-Clayton suggestion is to ad hoc identify the model (53) through

$$\begin{aligned} \alpha_i &= \alpha_1 + \alpha_2 i + \sum_{\ell=3}^i \sum_{h=3}^{\ell} \Delta^2 \alpha_h, \\ \beta_j &= \beta_1 + \beta_2 j + \sum_{\ell=3}^j \sum_{h=3}^{\ell} \Delta^2 \beta_h, \\ \gamma_k &= \gamma_1 + \gamma_2 k + \sum_{\ell=3}^k \sum_{h=3}^{\ell} \Delta^2 \gamma_h, \end{aligned} \tag{77}$$

$$\delta = 0.$$

A dynamic prior is chosen so that the double differences  $\Delta^2 \alpha_i, \Delta^2 \beta_j,$  and  $\Delta^2 \gamma_k$  are independent zero mean normal with variances  $\phi = (\sigma_\alpha^2, \sigma_\beta^2, \sigma_\gamma^2)$  that have  $\chi^2$ -type prior. The purpose of this is in part to facilitate extrapolations  $\Delta^2 \alpha_i, \Delta^2 \beta_j,$  and  $\Delta^2 \gamma_k$  for  $i > I, j > J,$  and  $k > K,$  which is done through further draws from normal distributions. The level/trend effects  $\theta_{level} = (\alpha_1, \alpha_2, \beta_1, \beta_2, \gamma_1, \gamma_2)'$  have independent uniform priors on some large intervals.

We will analyse the Berzuini-Clayton model as applied to an age-period data array  $\mathcal{S}_{ap}$ . Decompose the canonical parameter  $\xi$  from (54) into two parts: the slope and level parameters, say  $\xi_\mu = (\mu_{T1}, \mu_{T1} - \mu_{T-1,1}, \mu_{T2} - \mu_{T1})'$ , and the collection of double differences, say  $\xi_\Delta$ . The assumed prior for  $\xi_\Delta$  is a simple collection of independent normal distributions with variances  $\phi$ . The assumed prior for  $\xi_\mu$  is a linear combination of not only the independent uniform variables  $\theta_{level}$ , but also on  $\xi_\Delta$ , since the age double differences  $\Delta^2 \alpha_i$  are cumulated backwards in (54), but forwards in (77). Thus, the prior for  $\xi = (\xi'_\mu, \xi'_\Delta)'$  depends on the  $\theta_{level}$  construction.

We get a hyper-parameter  $\lambda_{hyper} = (\lambda, \phi)$ , where  $\lambda$  is some three-dimensional ad hoc identified level/trend effect dependent on  $\theta_{level}, \xi_\Delta$ . We will argue that the ad hoc identified level/trend effect  $\lambda$  will wash out in the Berzuini-Clayton model. However, the level/trend parameter  $\xi_\mu$  is a function of the  $\theta_{level}$  construction that is tailored to the ad hoc identification. That construction remains in the analysis.

In the presentation of the posterior Berzuini and Clayton are careful only to consider the double differences  $\xi_\Delta$  and stay clear of the ad hoc identified level/trend effect  $\theta_{level}$ . Theorem 2 yields the posterior  $p(\xi | y) = p(y | \xi)p(\xi)/p(y)$ . Thus, the marginal posterior for the double differences is  $p(\xi_\Delta | y) = \int p(y | \xi_\Delta, \xi_\mu)p(\xi_\Delta, \xi_\mu)d\xi_\mu/p(y)$ . This links  $\xi_\Delta$  to  $\xi_\mu$  and in turn to the  $\theta_{level}$  construction.

The extrapolative method is based on double differences so it only depends on  $\lambda_{hyper}$  through  $\phi$  due to Theorem 9 and the subsequent discussion. Thus, the extrapolative method is of the form  $p(\tilde{\mu} | \xi, \lambda_{hyper}, y) = p(\tilde{\mu} | \xi, \phi, y)$ . By construction it does not reduce to  $p(\tilde{\mu} | \xi, y)$  so that condition (39) for Theorem 3 is not satisfied. The distribution forecast is of the form

$$p(\tilde{\mu} | y) = \iint p(\tilde{\mu} | \xi, \phi, y) p(\xi | y) p(\phi | \xi) d\phi d\xi, \tag{78}$$

which, apart from depending on the  $\theta_{level}$  construction, also depends on the conditional prior  $p(\phi | \xi)$ , which is not updated by the likelihood.

In summary, it appears that the Berzuini-Clayton analysis depends on the  $\theta_{level}$  construction as well as the conditional prior  $p(\phi | \xi)$ . The dependence on the  $\theta_{level}$  construction could of course be addressed by introducing priors directly on  $\xi_\mu$ , which in turn would be updated by the likelihood. Since the conditional prior  $p(\phi | \xi)$  cannot be updated by the likelihood that its sole justification rests on the substantial context.

**5.4.5. A Functional Form Hypothesis.** It is instructive to consider functional form restrictions on the time effects. Such hypotheses can be analysed using the results outlined in Section 3.4.2. As an example restrict the age effect to be quadratic in a similar way to Yang and Land [5] so that

$$\alpha_i = \sigma_0 + \sigma_1 i + \sigma_2 i^2 \quad \text{for } i = 1, \dots, I. \tag{79}$$

This restriction on the time effect can be analysed by writing it on the form  $R'\theta = \rho$ , see (28), and then applying



Theorem A.3. Alternatively, in this particular case, we can show that the restriction actually only affects the ad hoc identified time effect through the canonical parameter, so a simpler analysis can be made.

A quadratic polynomial has constant second order derivative. Therefore the restriction (79) implies

$$\Delta^2 \alpha_i = 2\sigma_2 \quad \text{for } i = 3, \dots, I. \quad (80)$$

This expression has one free parameter. Thus, it is useful to consider the third order difference:

$$\Delta^3 \alpha_i = \Delta^2 \alpha_i - \Delta^2 \alpha_{i-1} = 0 \quad \text{for } i = 4, \dots, I. \quad (81)$$

This gives  $I - 3$  linear restrictions on the canonical parameter. The age time effect  $\alpha_i$  then has three remaining parameters, say  $\alpha_1, \alpha_2$ , and  $\alpha_3$ . These are freely varying since the parameters  $\sigma_0, \sigma_1$ , and  $\sigma_2$  are freely varying.

If the constraint is imposed directly on the canonical parameter, the restricted model is a regular exponential family with the advantages outlined in Section 2.4. However, if the analysis is done with the time effect the levels and trend will have to be ad hoc identified while bearing in mind the issues discussed above.

5.4.6. *The “Hierarchical Age-Period Cohort Regression Model”.* In some cases a random effects approach can be used to get an overview of the many parameters of the age-period model. When applied to the time effects this implies an ad hoc identification. An example is the “hierarchical age-period cohort regression model” by Yang and Land [5]. In that paper the age effect is given a quadratic structure, but that does not have to be the case. The model is then given by

$$\begin{aligned} \alpha_i &= \sigma_0 + \sigma_1 i + \sigma_2 i^2, & \beta_j &\stackrel{D}{=} \mathbf{N}(0, \sigma_\beta^2), \\ \gamma_k &\stackrel{D}{=} \mathbf{N}(0, \sigma_\gamma^2), & \delta &= 0. \end{aligned} \quad (82)$$

Since random effects are only introduced for some of the time effects, the analysis of Section 4.3 has to be modified in a similar way to the analysis in Section 5.4.4.

From (80) it is seen that the model restricts  $\Delta^2 \alpha_i = 2\sigma_2$ . Thus, divide the canonical parameter  $\xi$  into three elements: the slope and level parameters, say  $\xi_\mu = (\mu_{11}, \mu_{11} - \mu_{I-1,1}, \mu_{12} - \mu_{I1})'$ , the age-double differences  $\xi_\alpha = (\Delta^2 \alpha_3, \dots, \Delta^2 \alpha_I)$ , and the remaining double differences  $\xi_{\beta,\gamma}$ . Here  $\xi_\alpha$  is restricted by the hypothesis  $2\sigma_2$  and  $\xi_{\beta,\gamma}$  is linear function of the normal random effects, while  $\xi_\mu$  is a three-dimensional linear function of  $\sigma_2$  and of the six-dimensional object  $\nu = (\sigma_0, \sigma_1, \beta_1, \beta_2, \gamma_1, \gamma_2)'$ . This leaves a three-dimensional ad hoc identified level/slope parameter  $\lambda$  which is also a function of  $\nu$  but not entering the likelihood. Let  $\psi = (\sigma_0, \sigma_1, \sigma_2, \sigma_\beta^2, \sigma_\gamma^2)$ .

The random effects likelihood are constructed in three steps. First, we have the usual age-period-cohort likelihood  $p(y | \xi)$ . Secondly, the random effects distribution for  $\xi_\mu, \xi_{\beta,\gamma}$ , and  $\lambda$  is multivariate normal, while  $\xi_\alpha$  is deterministic function of  $\psi$ . Thus, decompose the prior as  $p(\xi_\mu, \xi_{\beta,\gamma}, \lambda | \psi) = p(\xi_\mu, \xi_{\beta,\gamma} | \psi)p(\lambda | \xi_\mu, \xi_{\beta,\gamma}, \psi)$ . Thirdly, following

Section 4.3 the random effects likelihood will not depend on  $p(\lambda | \xi_\mu, \xi_{\beta,\gamma}, \psi)$  and it is given by

$$\begin{aligned} p(y | \psi) &= \int p(y | \xi_\alpha, \xi_\mu, \xi_{\beta,\gamma}) p(\xi_\alpha, \xi_\mu, \xi_{\beta,\gamma} | \psi) d(\xi_\mu, \xi_{\beta,\gamma}). \end{aligned} \quad (83)$$

The prior  $p(\lambda | \xi_\mu, \xi_{\beta,\gamma}, \psi)$  is not updated by the data. Plots and inferences based on the posterior  $p(\theta | \psi, y)$  will then suffer from the ad hoc identification issues outlined in Section 3.4.

5.5. *A Two-Sample Age-Period-Cohort Model.* When confronted with two samples for women and for men it may be of interest to apply the age-period-cohort model (43) to each of the samples and impose that some of the time effects are the same across samples. The models for samples  $r = 1, 2$  are

$$\mu_{ijr} = \alpha_{ir} + \beta_{jr} + \gamma_{kr} + \delta_r \quad \text{for } i, j \in \mathcal{J}, r = 1, 2. \quad (84)$$

The time effect  $\theta = (\dots, \alpha_{ir}, \beta_{jr}, \gamma_{kr}, \delta_r, \dots)'$  now varies in  $\Theta = \mathbb{R}^q$  where  $q = 4(I + J)$ .

5.5.1. *Analysis of the Unrestricted Two-Sample Model.* The unrestricted two-sample model is simply analysed as two copies of the one sample model of Section 5.1. The time effects of each copy are only defined up to linear trends. The group of transformations characterizing the identification problem combines two copies of the one sample group (44). The maximal invariant parameter is  $\xi = (\xi'_1, \xi'_2)' \in \mathbb{R}^p$  where  $p = 4(I + J - 2)$  and each of  $\xi_r$  are of the form (45). The benefits of Section 2 hold when working with that parameter.

5.5.2. *Bayesian Ad Hoc Identification Using a Dynamic Model.* An application of the unrestricted two-sample model can be found in Cairns et al. [36]. The two samples are the population of England and Wales and the subpopulation of assured lives, so the substantive question is whether there is a selection effect for the assured lives. A Bayesian model with dynamic prior is used. It shares some features with the Berzuini and Clayton [6] model discussed in Section 5.4.4 although the details of the ad hoc identification of the levels and slopes are slightly different. When it comes to forecasting the extrapolative method appears to depend on the ad hoc identified parameter as well as the hyperparameters. This complicates the analysis of the forecast relatively the discussion in Section 5.4.4.

5.5.3. *The Hypothesis of Common Period Parameters.* The two-sample model allows the possibility for adding cross-sample restrictions on the parameters. As an example we consider the hypothesis of common period parameters.

Working with the canonical parameter the hypothesis is

$$\Delta^2 \beta_{i1} = \Delta^2 \beta_{i2} \quad \text{for } j = 3, \dots, J. \quad (85)$$

This is a simple linear restriction as that discussed in Section 2.4.3. It is readily seen that the degrees of freedom of

the hypothesis are  $p - p_H = J - 2$  so the dimension of the restricted model is  $p_H = 4I + 3J - 6$ . The canonical parameter under the hypothesis is then

$$\xi_H = (\dots, \mu_{11r}, \mu_{21r} - \mu_{11r}, \mu_{12r} - \mu_{11r}, \Delta^2 \alpha_{ir}, \Delta^2 \beta_j, \Delta^2 \gamma_{kr}, \dots)'. \tag{86}$$

The same result arises when writing the hypothesis in terms of time effects so that

$$\beta_{j1} = \beta_{j2} \quad \text{for } j = 1, \dots, J. \tag{87}$$

Such hypotheses on the time effect were discussed in Section 3.4.2. It can be analysed using the general result in Theorem A.3. However, we will take the simpler route of arguing that this only restricts the canonical parameter given a hypothesis of the type (85). The argument relies on noting that analysing the restriction for the predictors  $\mu_{ij1}$  and  $\mu_{ij2}$  is equivalent to analysing the restriction for the predictors  $\mu_{ij1}$  and  $\mu_{ij2} - \mu_{ij1}$ , where the cross-sample differenced predictor is of the form

$$\mu_{ij2} - \mu_{ij1} = (\alpha_{i2} - \alpha_{i1}) + (\beta_{j2} - \beta_{j1}) + (\gamma_{k2} - \gamma_{k1}) + (\delta_2 - \delta_1). \tag{88}$$

Now, the restricted model for the cross-sample differenced predictor  $\mu_{ij2} - \mu_{ij1}$  is an age-cohort model:

$$\mu_{ij2} - \mu_{ij1} = (\alpha_{i2} - \alpha_{i1}) + (\gamma_{k2} - \gamma_{k1}) + (\delta_2 - \delta_1). \tag{89}$$

Following the analysis of Section 5.3 the (87) therefore implies the  $J - 2$  linear restrictions given by (85). At the same time the predictor for the first sample  $\mu_{ij1}$  is left unrestricted by (87). In summary, the restrictions (85) and (87) are equivalent.

The restriction has an interesting implication for the interpretation of the involved double differences. For the unrestricted model it was found that only plain double differences, such as  $\Delta^2 \alpha_{jr}$ , are identified. Under the restriction the cross-sample differenced predictor is of age-cohort form (89) so also the cross-double differences  $\Delta(\alpha_{i2} - \alpha_{i1})$  and  $\Delta(\gamma_{k2} - \gamma_{k1})$  are identified.

**5.5.4. Step-Wise Ad Hoc Identification under the Hypothesis.** The analysis of Riebler and Held [7] finds that the difference  $\alpha_{i2} - \alpha_{i1}$  is identified under the hypothesis (85). This is not consistent with the above analysis showing that the cross-sample differenced predictor is an age-cohort model under the hypothesis, for which levels such as  $\alpha_{i2} - \alpha_{i1}$  are identified.

The apparent difference comes about because Riebler and Held follow a step-wise identification approach along the lines of Sections 3.3 and 5.4.1. In a first step the time effects  $\alpha_{ir}$ ,  $\beta_{jr}$ , and  $\gamma_{kr}$  are constrained to have zero-sums as in (69). In a second step the slopes are ad hoc identified using a Bayesian approach similar to that of Berzuini and Clayton [6]; see Sections 4 and 5.4.4 for a discussion of the consequences.

The identification in the first step implies that  $\alpha_{i2} - \alpha_{i1}$  has a zero sum. Under the hypothesis (85) this is exactly what is

needed to ad hoc identify the levels in the age-cohort model (89). In other words a different level identification in the first step leads to a different level for the difference  $\alpha_{i2} - \alpha_{i1}$ .

## 6. Models with Nonlinear Parametrisations

Some additional issues arise when looking at models with nonlinear parametrisations. A prominent example is the mortality model proposed by Lee and Carter [3] and which is the current benchmark in mortality studies done by government agencies and pension funds. For this model the time effect space  $\Theta$  has a nondifferentiability which can actually be avoided by working directly with the parameter space  $M$ .

We analyze the Lee-Carter model in Section 6.1. In Section 6.2 we turn to a two-sample problem where some additional difficulties can arise when forecasting.

**6.1. The Lee-Carter Model.** The mortality model proposed by Lee and Carter [3] has predictor of the form

$$\mu_{ij} = \alpha_i + \beta_j \kappa_j \quad \text{for } i, j \in \mathcal{F}_{\text{ap}}. \tag{90}$$

The time effects  $\theta = (\alpha_1, \dots, \alpha_I, \beta_1, \dots, \beta_I, \kappa_1, \dots, \kappa_J)$  vary in  $\Theta = \mathbb{R}^{2I+J}$ .

Lee and Carter pointed towards two identification issues of the model. If  $\alpha$ ,  $\beta$ , and  $\kappa$  are one solution to (90), then  $\alpha - \beta c$ ,  $\beta$ ,  $\kappa + c$  is also a solution for any scalar  $c$ , just as  $\alpha$ ,  $\beta/d$ , and  $\kappa d$  are a solution for any  $d \neq 0$ . Consequently, they proposed the ad hoc identification:

$$\sum_i \beta_i = 1, \quad \sum_j \kappa_j = 0. \tag{91}$$

This is, however, not the full story about the identification issues. To get at this we follow the outline from the linear parametrised models and start by finding the parameter space for the predictor  $\mu$ .

**6.1.1. The Parameter Space.** We start by finding the predictor space  $M$ . Write the model in matrix form. Let  $\underline{\mu}$  denote the  $I \times J$ -matrix of  $\mu_{ij}$ . Then

$$\underline{\mu} = \alpha \iota' + \beta \kappa', \tag{92}$$

where  $\alpha$ ,  $\beta$ , and  $\kappa$  are vectors concatenating  $\alpha_i$ ,  $\beta_i$ , and  $\kappa_j$  and where  $\iota = (1, \dots, 1)' \in \mathbb{R}^I$ . Postmultiply by the projection identity  $I_J = \bar{u}' + \bar{t}_\perp \bar{t}'_\perp$  to get

$$\underline{\mu} = \alpha \iota' + \beta \kappa' (\bar{u}' + \bar{t}_\perp \bar{t}'_\perp) = (\alpha + \beta \kappa' \bar{t}_\perp) \iota' + \beta (\kappa' \bar{t}_\perp) \bar{t}'_\perp, \tag{93}$$

where the orthogonal complement  $\iota_\perp$  can be chosen so that  $\bar{t}'_\perp \kappa = (\Delta \kappa_2, \dots, \Delta \kappa_J)'$  but could also be chosen otherwise. Equation (93) shows that the model is composed of two matrices with rank one. Thus, the parameter space is given by

$$M = \left\{ \underline{\mu} \in \mathbb{R}^{I \times J} : \underline{\mu} = \gamma \iota' + \delta \bar{t}'_\perp \right. \\ \left. \text{for } (\gamma, \delta) \in \mathbb{R}^I \times \mathbb{R}^{I \times (J-1)} \right. \\ \left. \text{so } \text{rank}(\delta) \leq 1 \right\}. \tag{94}$$

Note that  $M$  does not depend on the normalisation of  $\iota_{\perp}$  since  $\delta$  is freely varying. The space  $M$  is a manifold since the space of matrices  $\delta$  with an upper bound to the rank is a manifold as opposed to the space where  $\delta$  has rank of unity. This space can be parametrised parsimoniously by

$$\xi = (\gamma, \delta) \quad \text{where } \gamma = \alpha + \beta\kappa' \iota, \quad \delta = \beta\kappa' \iota_{\perp}, \quad (95)$$

varying in the manifold

$$\Xi = \{(\gamma, \delta) \in \mathbb{R}^I \times \mathbb{R}^{I \times (J-1)} : \text{rank}(\delta) \leq 1\}. \quad (96)$$

The  $\xi$  is the candidate for the maximal invariant describing the equivalence classes of the mapping from the time effect  $\theta$  to the predictor  $\mu$ .

The next step is to analyse the time effect space  $\Theta$ . It is convenient to decompose  $M$  into two disjoint sets depending on the rank of  $\delta$ . These sets are

$$M_1 = \{\underline{\mu} \in \mathbb{R}^{I \times J} : \underline{\mu} = \gamma \iota' + \delta \iota_{\perp}' \text{ for } (\gamma, \delta) \in \mathbb{R}^I \times \mathbb{R}^{I \times (J-1)} \\ \text{so } \text{rank}(\delta) = 1\},$$

$$M_0 = \{\underline{\mu} \in \mathbb{R}^{I \times J} : \underline{\mu} = \gamma \iota' \text{ for } \gamma \in \mathbb{R}^I\}. \quad (97)$$

Correspondingly, the time effect space  $\Theta$  can be decomposed into two disjoint sets:

$$\Theta_1 = \{\theta \in \Theta : \exists i, j \text{ so } \beta_i \kappa_j \neq \beta_i \kappa_j\},$$

$$\Theta_0 = \{\theta \in \Theta : \forall i, j \text{ so } \beta_i \kappa_j = \beta_i \kappa_j\}. \quad (98)$$

Note that  $\delta = 0$  if and only if  $\theta \in \Theta_0$ . Consider first the time effect space  $\Theta_1$ , which is implicitly what Lee and Carter had in mind. The mapping  $\theta \mapsto \mu$  on  $\Theta_1$  to  $M$  is invariant to the group of transformations:

$$g_1 : \begin{pmatrix} \alpha_i \\ \beta_i \\ \kappa_j \end{pmatrix} \mapsto \begin{pmatrix} \alpha_i + \beta_i c \\ \beta_i \\ d \\ (\kappa_j - c) d \end{pmatrix}, \quad (99)$$

acting on  $\Theta_1$  for all  $c \in \mathbb{R}$  and all  $d \neq 0$ . The parameter  $\xi = (\gamma', \delta)'$  is invariant under  $g_1$  acting on  $\Theta_1$ . Now, consider the space  $\Theta_0$  with deficient rank. Then  $\alpha_i, \beta_i,$  and  $\kappa_j$  map into  $\alpha_i + \varphi_i$  where  $\varphi_i = \beta_i \kappa_j$  is constant in  $j$ , so that  $\delta = \beta \kappa' \iota_{\perp} = 0$ . This mapping is invariant to the group of transformations:

$$g_0 : \begin{pmatrix} \alpha_i \\ \beta_i \kappa_j \end{pmatrix} \mapsto \begin{pmatrix} \alpha_i + a_i \\ \beta_i \kappa_j - a_i \end{pmatrix}, \quad (100)$$

acting on  $\Theta_0$  for all  $(a_1, \dots, a_I)' \in \mathbb{R}^I$ .

**Theorem 10.** Let  $\underline{\mu} \in M$ . The parameter  $\xi \in \Xi$  of (95) satisfies the following:

- (i)  $\xi$  is a function of  $\theta \in \Theta$  which is invariant to the groups  $g_0, g_1$  in (99) and (100);

- (ii)  $\underline{\mu}$  is a function of  $\xi$ ;
- (iii) the parametrisation of  $\underline{\mu}$  by  $\xi$  is exactly identified in the sense that  $\xi^{\dagger} \neq \xi^{\ddagger} \Rightarrow \underline{\mu}(\xi^{\dagger}) \neq \underline{\mu}(\xi^{\ddagger})$ .

Theorem 10 shows that  $\xi$  varies freely on the space  $\Xi$  and it gives a unique parametrisation of  $\mu$ . As a function of  $\theta$  it is invariant to  $g_0, g_1$ ; hence it is a maximal invariant.

It is interesting to compare the properties of the spaces  $M, \Xi,$  and  $\Theta$ . The spaces  $M$  and  $\Xi$  are spaces of matrices with deficient rank. These are smooth spaces, but they are not vector spaces since the sum of matrices with rank one may have rank larger than one. In contrast  $\Theta$  is a vector space. The mapping from  $\Theta$  to  $M$  will inevitably be nondifferentiable. This nondifferentiability is avoided by working directly with  $M$ . Likewise, in a Bayesian setting it would seem more difficult to introduce a meaningful prior of  $\Theta$  with its nondifferentiability than on  $M$ .

**6.1.2. Maximum Likelihood Estimation.** The maximum likelihood estimator for  $\xi$  can be derived analytically in the normal case.

Consider a situation where the data array is of age-period form so  $Y_{ij}$  for  $(i, j) \in \mathcal{F}_{ap}$ . Suppose  $Y_{ij}$  are independent normal with mean  $\mu_{ij}$  and variance  $\sigma^2$ . Organise the data in a matrix  $\underline{Y}$ . Then the log likelihood is of the form

$$\ell(\underline{\mu}, \sigma^2; \underline{Y}) = -\frac{IJ}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \text{tr} \left\{ (\underline{Y} - \underline{\mu})(\underline{Y} - \underline{\mu})' \right\}. \quad (101)$$

The maximum likelihood estimator is of the following form. Subsequently, this is related to the estimator suggested by Lee and Carter.

**Theorem 11.** For a normal age-period array parametrised by (94) the maximum likelihood estimators are

$$\hat{\gamma} = \underline{Y} \iota (\iota' \iota)^{-1}, \quad \hat{\delta} = \left[ \text{svd}_1 \left\{ \underline{Y} \iota_{\perp} (\iota'_{\perp} \iota_{\perp})^{-1} \iota'_{\perp} \right\} \right] \iota_{\perp} (\iota'_{\perp} \iota_{\perp})^{-1}, \quad (102)$$

where  $\text{svd}_1(\cdot)$  is the singular value decomposition truncated to one factor.

Thus,  $\gamma$  is estimated by the row-averages of the data matrix, while  $\delta$  is estimated by the singular value decomposition of the row-wise demeaned data matrix.

**6.1.3. Estimation of Ad Hoc Identified Time Effects.** The ad hoc identification (91) gives a time effect  $\theta_{\lambda}$  varying in a  $2I + J - 2$  dimensional affine subspace of  $\Theta = \mathbb{R}^{2I+J}$ . The ad hoc identified  $\theta_{\lambda}$  can now be expressed in terms of the maximal invariant parameter  $\xi$  using (95). In the case where  $\delta \neq 0$  then it has singular value decomposition  $\delta = \delta_L \delta_S \delta_R'$  for two vectors  $\delta_L \in \mathbb{R}^I$  and  $\delta_R \in \mathbb{R}^{J-1}$  so  $\delta_L' \delta_L = 1$  and  $\delta_R \delta_R' = 1$ , while  $\delta_S > 0$  is a positive scale. The ad hoc identification of Lee and Carter then gives

$$\alpha_{\lambda} = \gamma, \quad \beta_{\lambda} = \delta_L (\iota' \delta_L)^{-1}, \quad \kappa_{\lambda} = \iota_{\perp} \delta_R \iota' \delta_L \delta_S. \quad (103)$$

Inserting the maximum likelihood estimators from Theorem 11 yields the estimators proposed by Lee and Carter. However, the disentangling of the singular values and singular vectors of  $\hat{\delta}$  is done by the ad hoc identification  $\beta'_{i'} = 1$  and  $\kappa'_{i'} = 0$ . These estimators are therefore specific to the considered data array and data set in parallel with the discussion in Sections 3.2 and 5.4.1.

**6.1.4. Consequences of the Possible Rank Deficiency.** The parameter space  $M$  was split into spaces  $M_1$  and  $M_0$  depending on the rank of  $\delta$ . The space  $M_0$  is a Lebesgue null set relative to  $M$ . Broadly speaking, there are two consequences of the possible rank deficiency. The first consequence is an estimation problem arising in the vicinity of  $M_0$ . The second consequence is that the usual normal asymptotic distribution theory does not apply in the vicinity of  $M_1$ . Whether this becomes a problem in practice depends on the data. One solution is to ensure that the time effect really is present when using the Lee-Carter model.

Investigate whether the time effects are present amounts to estimating the rank of  $\delta$ . For a given data set two Lee-Carter models can be estimated. The first model with predictor space  $M$  is the unrestricted model in which  $\text{rank}(\delta) \leq 1$ . The second model has predictor space  $M_0$  so  $\delta = 0$ . Twice the difference of the likelihood values gives a likelihood ratio test statistic which is asymptotically  $\chi^2$ . If the smaller model,  $M_0$ , is accepted this is used in subsequent analysis. However, if the smaller model  $M_0$  is rejected then it is likely that the predictor is not located in the vicinity of  $M_0$  and it is then safe to work with the predictor space  $M_1$ .

The consistency of this step-wise procedure is discussed in a cointegration context by Johansen [37, Section 12]. Even when this procedure points towards working with the parameter space  $M_1$  the rank deficiency may still affect inference under  $M_1$ . Analysis of simple canonical correlation models suggests that inference under  $M_1$  will be nearly similar if the distance to  $M_0$  is sufficiently large. A problem is that the distribution for the test statistic will have poor finite sample properties when the parameter value is close to  $M_0$ . A simple way to get around this problem is to test for  $M_0$  using a test with lower level than the conventional level. A more complicated way to address this is to employ a finite sample correction when seeking to test for  $M_0$ . See Nielsen [38, 39] for further discussion of this issue in the context of simple canonical correlation models.

The rank deficiency issue is typically not encountered in a standard Lee-Carter analysis. The reason is that the analysis is typically applied to data where there is a marked improvement in mortality rates over time. A Lee-Carter analysis could however run into trouble if it were applied to data without a strong calendar effect. The issue becomes more pertinent when extending Lee-Carter model with a cohort component such as

$$\mu_{ij} = \beta_i^{(1)} + \beta_i^{(2)} \kappa_j^{(2)} + \beta_i^{(3)} \gamma_{j-i+1}^{(3)}; \tag{104}$$

see Renshaw and Haberman [40]. If the cohort effect is modest the latter matrix is nearly rank deficient and the likelihood will be nearly flat in certain directions. This is

presumably the reason for the estimation problem noted by Cairns et al. [41].

**6.1.5. Forecasting.** The purpose of Lee-Carter model is usually to forecast future mortality. This issue is considered for the model with parameter space  $M_1$ . The standard approach is to extrapolate  $\kappa$ , ad hoc identified through, for instance,  $\kappa'_{i'} = 0$ . The  $h$ -step ahead extrapolation of  $\kappa_{J+h}$  based on some forecast methods is denoted by  $\tilde{\kappa}_{J+h}(\hat{\kappa})$ . Combined with the estimates  $\hat{\alpha}_i, \hat{\beta}_i$  this gives the overall forecast

$$\tilde{\mu}_{i,J+h}(\hat{\theta}) = \hat{\alpha}_i + \hat{\beta}_i \tilde{\kappa}_{J+h}(\hat{\kappa}). \tag{105}$$

The identification question is then for which extrapolation methods this equals

$$\tilde{\mu}_{i,J+h} \{g_1(\hat{\theta})\} = (\hat{\alpha}_i + \hat{\beta}_i c) + \left(\frac{\hat{\beta}_i}{d}\right) \tilde{\kappa}_{J+h} \{(\hat{\kappa} - c) d\}. \tag{106}$$

The condition for avoiding adverse impact of the ad hoc identification is as follows.

**Theorem 12.** *Let  $\mu \in M_2$ . The forecast  $\tilde{\mu}_{i,J+h}$  in (105) is invariant to ad hoc identification if and only if the extrapolation method for the period effect is location-scale preserving:*

$$\tilde{\kappa}_{J+h} \{(\hat{\kappa} - c) d\} = \{\tilde{\kappa}_{J+h}(\hat{\kappa}) - c\} d \quad \forall c \in \mathbb{R}, \forall d \neq 0. \tag{107}$$

The default forecast method in the literature is a random walk with a drift, which was the preferred forecast of Lee and Carter [3]. This is given by

$$\tilde{\kappa}_{J+h} = \tilde{\kappa}_{J+h-1} + \nu_c + \varepsilon_h, \tag{108}$$

with estimates  $\hat{\nu}_c = (J - 1)^{-1} \sum_{j=2}^J (\tilde{\kappa}_j - \tilde{\kappa}_{j-1})$  and normal errors  $\varepsilon_h$  with mean zero and estimated variance  $\hat{\sigma}^2(\hat{\kappa}) = (J - 2)^{-1} \sum_{j=2}^J (\tilde{\kappa}_j - \tilde{\kappa}_{j-1} - \hat{\nu}_c)^2$ . This extrapolation method is location-scale preserving as required in Theorem 12. It is even linear trend preserving. Other valid forecasts are a random walk without intercept as given by the equation  $\tilde{\kappa}_{J+h} = \tilde{\kappa}_{J+h-1} + \varepsilon_h$ , or an autoregression given by  $\tilde{\kappa}_{J+h} = \rho \tilde{\kappa}_{J+h-1} + \nu_c + \varepsilon_h$ .

An alternative approach to forecasting would consider the predictor of the model for a particular age ground, say  $i$ . That predictor is  $\hat{\mu}_i = e'_i(\hat{\gamma}'_{i'} + \hat{\delta}'_{i'})$ , where  $e_i$  is the  $i$ th unit vector. From this we can generate forecasts  $\tilde{\mu}_{i,J+h}$  using any time series method. The resulting forecast will in general depend on  $\hat{\kappa}$  as well as  $\hat{\alpha}_i, \hat{\beta}_i$  and it is therefore more general than the forecasts discussed in Theorem 12, which only depends on  $\hat{\kappa}$ . The forecast for another age group, say  $i'$ , should be the same up to a linear transformation dictated by the Lee-Carter structure. Thus, the  $h$ -step ahead forecasts for the entire array are

$$\tilde{\underline{\mu}}_{-J+h} = \frac{\hat{\delta}'_{i'} e'_i}{e'_i \hat{\delta}'_{i'} e'_i} (\tilde{\mu}_{i,J+h} - e'_i \hat{\gamma}'_{i'}) + \hat{\gamma}'_{i'}, \tag{109}$$

for an index  $i'$  is chosen so that  $e'_i \hat{\delta}'_{i'} e'_i \neq 0$ .

6.1.6. *Bayesian Ad Hoc Identification Using a Dynamic Model.*

A Bayesian model with dynamic specification of the prior has been suggested by Pedroza [42]. Dynamic priors are presented for the time effects  $\theta = (\xi, \lambda)$  involving a hyper parameter  $\phi$ . The ad hoc identification (91) is imposed so that analysis is made for an ad hoc identified time effect  $\theta_\lambda$ .

Pedroza presents posteriors for  $\theta_\lambda$ . When evaluating this posterior one should bear in mind that the conditional prior  $p(\lambda | \xi)$  is not updated by the data; see Theorem 2. The presented extrapolative method does not depend on  $\lambda$ . Even so, the forecast will depend on conditional prior  $p(\phi | \xi)$  which is not updated by the data; see Theorem 3.

6.2. *The Two-Sample Lee-Carter Model.* We now turn to applications of the Lee-Carter model in two-sample problems. Suppose two samples are for women and men. One approach would be to fit separate Lee-Carter models to the two datasets. These Lee-Carter models are of the form

$$\mu_{ijr} = \alpha_{ir} + \beta_{ir}\kappa_{jr} \quad \text{for } i, j \in \mathcal{J}_{ap}, r = 1, 2. \quad (110)$$

The objective is now to extrapolate the period effects  $\kappa_{jr}$ . Extrapolating the two models separately using separate random walks is often seen to be volatile, so methods that seek to combine information from both estimated series  $\hat{\kappa}_{jr}$  are sought after. The next result describes the invariance problem in forecasting.

**Theorem 13.** *Let  $\mu_r \in M_{2r}$  for  $r = 1, 2$ . The forecast  $\tilde{\mu}_{i,j+h,1}$  for sample  $r = 1$  is invariant to ad hoc identification if the extrapolation method  $\tilde{\kappa}_{j+h,1}$  preserves location/scale for sample 1, but is invariant to location and scale for sample 2. That is for all  $c_1, c_2 \in \mathbb{R}$  and all  $d_1, d_2 \neq 0$ ; then*

$$\tilde{\kappa}_{j+h,1} \{d_1(\hat{\kappa}_1 - c_1), d_2(\hat{\kappa}_2 - c_2)\} = d_1 \{\tilde{\kappa}_{j+h,1}(\hat{\kappa}_1, \hat{\kappa}_2) - c_1\}. \quad (111)$$

For one sample the standard forecasting technique appears to be the random walk with a drift as in (108). For the two-sample problem a suggestion could be that women and men should share a common random walk with a drift but deviate from this by a stationary process. In econometrics this idea is referred to as cointegration as proposed by Engle and Granger [43]; see also Johansen [37] for a likelihood based vector autoregressive approach. It is tempting to require that the calendar effects should cointegrate with coefficients of unity, so  $\kappa_{j1} - \kappa_{j2}$  should be stationary. However, that apparently intuitive choice violates Theorem 13 because the locations and scales of  $\kappa_{jr}$  are different and arbitrary.

There are two fixes to this problem. The first solution is to work directly with the mortality predictors  $\mu_{ijr}$  for an arbitrary age group  $i$  as outlined for the one-sample case in connection with (109). Since no identification is involved it is permitted to impose that  $\mu_{ij1}$  and  $\mu_{ij2}$  cointegrate with coefficients of unity. The forecast for age group  $i$  is then carried over to other age groups. The second solution is to work with the estimated series  $\hat{\kappa}_{jr}$  but estimate the cointegrating coefficients from the data. In other words, the cointegrating relation  $\hat{\kappa}_{j1} - \phi\hat{\kappa}_{j2} - \psi$  should be zero mean,

stationary, with coefficients  $\phi, \psi$  estimated from the data. This can, for instance, be done by Johansen’s approach for a bivariate vector autoregression; see Hendry and Nielsen [44, Section 17].

**7. Conclusion**

Ad hoc identification is intimately linked to interpretation, inference, numerical analysis, and forecasting. The ad hoc identification will often introduce an arbitrary element in the statistical analysis, whether it is based on frequentist or Bayesian methods. This arbitrary element is entirely avoidable and is in our view best avoided unless there is a clear substantial motivation for ad hoc identification. For decades there has been a debate over how it is best to ad hoc identify mortality models. Our proposal is to bypass this discussion by analysing the surjective mapping between the unidentified time effect parameter and the predictor of the model and then deduce a maximal invariant parametrisation. In our experience there are typically two substantial benefits. First, it simplifies estimation and other statistical computations which is what we have focused on here. Secondly and perhaps more importantly, it helps to focus the substantial question that gives rise to the analysis in the first place.

The issue of dealing with two time scales also occurs in other statistical models, such as the Cox regression model; see Cabrera et al. [45] for a recent application. In future research it would be interesting to consider whether the analysis presented here has any bearing on that problem.

**Appendix**

**A. Proofs**

*A.1. Some Linear Algebra Results*

**Lemma A.1.** *Consider  $A \in \mathbb{R}^{q \times p}$  and  $C \in \mathbb{R}^{q \times (q-q_C)}$  so  $p, q_C \leq q$ . Suppose they have full column rank. Then the following statements are equivalent for some  $p_C \leq p$ :*

- (i)  $(A, C) \in \mathbb{R}^{q \times (p+q-q_C)}$  has rank  $p_C + q - q_C$ ;
- (ii)  $A'_\perp C \in \mathbb{R}^{(q-p) \times (q-q_C)}$  has rank  $p_C + q - q_C - p$ ;
- (iii)  $C'_\perp A \in \mathbb{R}^{q_C \times p}$  has rank  $p_C$ .

*Proof of Lemma A.1.* (i)  $\Leftrightarrow$  (ii) Premultiply the matrix  $(A, C)$  with the invertible matrix  $(\bar{A}, A'_\perp)'$  to get the identity

$$\begin{aligned} \begin{pmatrix} \bar{A}' \\ A'_\perp \end{pmatrix} (A, C) &= \begin{pmatrix} I_p & \bar{A}' C \\ 0 & A'_\perp C \end{pmatrix} = \begin{pmatrix} I_p & 0 \\ 0 & A'_\perp C \end{pmatrix} M \\ &\text{where } M = \begin{pmatrix} I_p & \bar{A}' C \\ 0 & I_{p+q-q_C} \end{pmatrix}. \end{aligned} \quad (A.1)$$

Since the first matrix  $(\bar{A}, A'_\perp)'$  and the last matrix  $M$  have full rank then

$$\text{rank}(A, C) = \text{rank}(I_p) + \text{rank}(A'_\perp C) = p + \text{rank}(A'_\perp C). \quad (A.2)$$

(i)⇔(iii) Swap the roles of  $A, C$  so  $\text{rank}(A, C) = q - q_C + \text{rank}(C'_\perp A)$ .  $\square$

**Lemma A.2.** Consider  $A \in \mathbb{R}^{q \times p}$  and  $C \in \mathbb{R}^{q \times (q-q_C)}$  so  $p \leq q_C \leq q$ . Suppose  $(A, C) \in \mathbb{R}^{q \times (p+q-q_C)}$  has full column rank  $p + q - q_C$ . Then  $m = A'_\perp C \in \mathbb{R}^{(q-p) \times (q-q_C)}$  has full column rank and  $(A, C)$  has orthogonal complement given by  $(A, C)_\perp = A_\perp m_\perp$  where  $m_\perp \in \mathbb{R}^{(q-p) \times (q_C-p)}$ .

*Proof of Lemma A.2.* Since  $(A, C)$  has full column rank Lemma A.1 (i, ii) implies that  $m$  has full column rank. Since  $m'_\perp A'_\perp (A, C) = (m'_\perp A'_\perp A, m'_\perp m) = 0$  we argue that  $(A, C, A_\perp m_\perp)$  is invertible. Premultiply by the invertible matrix  $(\overline{A}, A_\perp)'$  to get

$$\begin{pmatrix} \overline{A} \\ A'_\perp \end{pmatrix} (A, C, A_\perp m_\perp) = \begin{pmatrix} I_p & \overline{A}' C & 0 \\ 0 & m & A'_\perp A_\perp m_\perp \end{pmatrix}. \quad (\text{A.3})$$

The matrix  $(m, A'_\perp A_\perp m_\perp)$  has full rank. Indeed, its inverse is

$$\begin{aligned} (m, A'_\perp A_\perp m_\perp)^{-1} &= \left[ (A'_\perp A_\perp)^{-1} m \{ m' (A'_\perp A_\perp)^{-1} m \}^{-1}, \right. \\ &\quad \left. m_\perp (m'_\perp A'_\perp A_\perp m_\perp)^{-1} \right]'. \end{aligned} \quad (\text{A.4})$$

Thus, the block triangular matrix (A.3) and, hence,  $(A, C, A_\perp m_\perp)$ , have full rank.  $\square$

A.2. Proofs of Main Theorems

*Proof of Theorem 1.* Since  $(A, C)$  has full column rank then Lemma A.2 shows it has orthogonal complement  $A_\perp m_\perp$  so that  $(A, C, A_\perp m_\perp)$  is invertible. The orthogonal projection identity shows

$$\begin{aligned} \theta &= (\overline{A, C}) \begin{pmatrix} A' \\ C' \end{pmatrix} \theta + A_\perp m_\perp (\overline{A_\perp m_\perp})' \theta \\ &= (\overline{A, C}) \begin{pmatrix} \xi \\ \psi \end{pmatrix} + A_\perp m_\perp \zeta, \end{aligned} \quad (\text{A.5})$$

by the constraint  $C'\theta = \psi$  and the definitions  $A'\theta = \xi$  and  $(\overline{A_\perp m_\perp})'\theta = \zeta$ . This defines the constrained time effect space  $\Theta_C$ . Consider now the mapping  $\theta \mapsto \mu = XA'\theta$ . Premultiply the above expression for  $\theta$  by  $A' = (I_p, 0)(A, C)'$  to get  $A'\theta = \xi$  so  $\mu = X\xi$ . Thus, since  $\psi$  is fixed, the equivalence classes in  $\Theta_C$  are given by  $g_C : \theta \mapsto A_\perp m_\perp \zeta$ , with  $\xi = A'\theta$  as a maximal invariant.  $\square$

*Proof of Theorem 2.* (i) With the likelihood (32) so  $p(y | \xi, \lambda) = p(y | \xi)$  then

$$\begin{aligned} p(y) &= \int \left\{ \int p(y | \xi, \lambda) p(\xi, \lambda) d\lambda \right\} d\xi \\ &= \int p(y | \xi) \left\{ \int p(\xi, \lambda) d\lambda \right\} d\xi \\ &= \int p(y | \xi) p(\xi) d\xi. \end{aligned} \quad (\text{A.6})$$

(ii) By Bayes formula and the likelihood (32) then

$$\begin{aligned} p(\xi | y) &= \frac{p(y | \xi) p(\xi)}{p(y)}, \\ p(\lambda | \xi, y) &= \frac{p(y | \xi, \lambda) p(\lambda | \xi)}{p(y | \xi)} = p(\lambda | \xi). \end{aligned} \quad (\text{A.7})$$

(iii) The posterior means are

$$\begin{aligned} E(\xi | y) &= \int \xi p(\xi | y) d\xi, \\ E(\lambda | \xi, y) &= \int \lambda p(\lambda | \xi, y) d\lambda = \int \lambda p(\lambda | \xi) d\lambda, \\ E(\lambda | y) &= \int \lambda p(\lambda | y) d\lambda = \int \lambda \left\{ \int p(\lambda, \xi | y) d\xi \right\} d\lambda \\ &= \int p(\xi | y) \left\{ \int \lambda p(\lambda | \xi) d\lambda \right\} d\xi, \end{aligned} \quad (\text{A.8})$$

noting that  $p(\lambda | \xi, y) = p(\lambda | \xi)$ .  $\square$

*Proof of Theorem 3.* Consider the expressions in (37) and (38); that is,

$$p(\bar{\mu} | y) = \iint p(\bar{\mu} | \xi, \lambda, y) p(\xi | y) p(\lambda | \xi) d\lambda d\xi \quad (\text{A.9})$$

and a similar expression involving  $p^\dagger$ . The question is when they are identical. Assuming  $p(\bar{\mu} | \xi, \lambda, y) = p(\bar{\mu} | \xi, y)$  as in (39) the expression reduces to

$$\begin{aligned} p(\bar{\mu} | y) &= \int p(\bar{\mu} | \xi, y) p(\xi | y) \left\{ \int p(\lambda | \xi) d\lambda \right\} d\xi \\ &= \int p(\bar{\mu} | \xi, y) p(\xi | y) d\xi, \end{aligned} \quad (\text{A.10})$$

since the conditional prior integrates to unity. The same applies for the expression involving  $p^\dagger(\lambda | \xi)$ .  $\square$

*Proof of Theorem 5.* Similar to the proof of Kuang et al. [19, Theorem 1], albeit for a rectangular instead of a triangular data array.  $\square$

*Proof of Theorem 8.* Recall  $\theta_{\text{intrinsic}} = C_{S_\perp} C'_{S_\perp} A (A' C_{S_\perp} C'_{S_\perp} A)^{-1} \xi$  where  $C_{S_\perp} = C_\perp (S'_\perp C_\perp)^{-1}$  as defined in (72). Premultiply by  $A'$  and  $C'$  to see that  $A'\theta_{\text{intrinsic}} = \xi$  and  $C'\theta_{\text{intrinsic}} = \psi$ . This does, however, not describe the full variation of  $\theta_{\text{intrinsic}}$  since  $(A, C)$  is not a square matrix. We must extend the matrix  $(A, C)$  with columns so that it is square and invertible.

We find a vector  $w \in \mathbb{R}^{q-3}$  so that  $(A, C, C_\perp w)$  is invertible. Recall  $A \in \mathbb{R}^{q \times p}$  where  $q = p + 4$ . The matrix  $C \in \mathbb{R}^{q \times 3}$  is chosen so that  $(A, C)$  has full column rank  $q - 1 = p + 3$  as discussed in Section 5.4.1. Apply Lemma A.2, swapping the role of  $A, C$ , to see  $w_\perp = C'_\perp A$ , say, has full column rank. Then  $(A, C)$  has orthogonal complement  $C_\perp w$ .

We show that for any invertible matrix  $M \in \mathbb{R}^{q \times q}$  then the  $1 \times 2$  block matrix  $\{(A, C), MC_\perp w\}$  is invertible.

To see this hold, premultiply by  $\{(A, C), C_{\perp} w\}'$  to see that an invertible upper triangular matrix arises.

To analyse the properties of  $\theta_{\text{intrinsic}}$  it suffices to analyse  $\{(A, C), MC_{\perp} w\}'\theta_{\text{intrinsic}}$ , since there is a bijective mapping between the two. Choose  $M = \overline{C}_{\perp}(C'_{\perp}S_{\perp})(S'_{\perp}C_{\perp})\overline{C}'_{\perp} + CC'$ . Then it holds that  $C'_{\perp}M' = (C'_{\perp}S_{\perp})(S'_{\perp}C_{\perp})\overline{C}'_{\perp}$  so that  $C'_{\perp}M'C'_{S_{\perp}}C'_{S_{\perp}} = C'_{\perp}$  and therefore

$$w'C'_{\perp}M'\theta_{\text{intrinsic}} = w'C'_{\perp}A(A'C'_{S_{\perp}}C'_{S_{\perp}}A)^{-1}\xi = 0, \quad (\text{A.11})$$

since  $w'C'_{\perp}A = 0$  by construction. Thus, it holds that  $w'(C'_{\perp}S_{\perp})(S'_{\perp}C_{\perp})\overline{C}'_{\perp}\theta_{\text{intrinsic}} = 0$  as required.  $\square$

*Proof of Theorem 9.* This is a generalisation of the proof of Kuang et al. [25, Theorems 1, 2]. Let  $k_h = J + h + I - i - K = h - i - 1$ .

(i) Recall the group  $g$  in (57). Then (i) follows by comparing the equations

$$\begin{aligned} \tilde{\mu}_{i,J+h}(\hat{\theta}) &= \hat{\alpha}_i + \tilde{\beta}_{J+h}(\hat{\theta}) + \tilde{\gamma}_{K+k_h}(\hat{\theta}) + \hat{\delta}, \\ \tilde{\mu}_{i,J+h}\{g(\hat{\theta})\} &= \hat{\alpha}_i + a + di + \tilde{\beta}_{J+h}\{g(\hat{\theta})\} + \tilde{\gamma}_{K+k_h}\{g(\hat{\theta})\} \\ &\quad + \hat{\delta} - a - b - c - dI. \end{aligned} \quad (\text{A.12})$$

(ii) As in Section 3.1 there is a bijective mapping from  $\theta$  to  $\lambda$ ,  $\xi$ , where  $\xi$  is invariant to  $g$ , but  $\lambda$  is not. The choice of  $\lambda$  is not important. Any extrapolations of  $\tilde{\beta}_j, \tilde{\gamma}_k$  can then be written in the form

$$\begin{aligned} \tilde{\beta}_{J+h}(\hat{\theta}) &= \hat{\beta}_J + h\Delta\hat{\beta}_J + F_{\beta}(\hat{\lambda}, \hat{\xi}), \\ \tilde{\gamma}_{K+k_h}(\hat{\theta}) &= \hat{\gamma}_K + k_h\Delta\hat{\gamma}_K + F_{\gamma}(\hat{\lambda}, \hat{\xi}), \end{aligned} \quad (\text{A.13})$$

for some functions  $F_{\beta}, F_{\gamma}$ . Applying the group  $g$  it follows

$$\begin{aligned} \tilde{\beta}_{J+h}\{g(\hat{\theta})\} &= \hat{\beta}_J + b - dJ + h(\Delta\hat{\beta}_J - d) + F_{\beta}\{g(\hat{\lambda}), \hat{\xi}\}, \\ \tilde{\gamma}_{K+k_h}\{g(\hat{\theta})\} &= \hat{\gamma}_K + c + dK + k_h(\Delta\hat{\gamma}_K + d) + F_{\gamma}\{g(\hat{\lambda}), \hat{\xi}\}. \end{aligned} \quad (\text{A.14})$$

Due to (i) it must hold  $F(\lambda, \xi) = F_{\beta}(\lambda, \xi) + F_{\gamma}(\lambda, \xi)$  and must equal  $F\{g(\lambda), \xi\} = F_{\beta}\{g(\lambda), \xi\} + F_{\gamma}\{g(\lambda), \xi\}$ . The function  $F$  must then be constant in the first argument. This must apply in the forecast region  $\mathcal{F}_{\text{ap}}$  where  $\gamma$  is not extrapolated. Therefore, it must also hold that the  $F_{\beta}(\lambda, \xi) = F_{\beta}\{g(\lambda), \xi\}$ , and in turn that  $F_{\gamma}(\lambda, \xi) = F_{\gamma}\{g(\lambda), \xi\}$ . Conversely, if the functions  $F_{\beta}, F_{\gamma}$  are constant in  $\lambda$  then the forecast is invariant to  $g$ .  $\square$

*Proof of Theorem 10.* (i) Equation (95) shows that  $\xi$  is a function of  $\theta$ .

(ii) Equation (93) shows that  $\underline{\mu} = \gamma' + \delta'_{\perp}$  is a function of  $\xi$ .

(iii) The decomposition  $\underline{\mu} = \underline{\mu}\underline{u}' + \underline{\mu}'_{\perp}t_{\perp}$  shows that there is a one-one mapping between  $\underline{\mu}$  and  $(\underline{\mu}, \underline{\mu}'_{\perp})$ . In turn, (93) shows that  $(\underline{\mu}, \underline{\mu}'_{\perp}) = (\gamma, \delta)$ . Thus if  $(\gamma^{\ddagger}, \delta^{\ddagger}) \neq (\gamma^{\ddagger}, \delta^{\ddagger})$  then  $\underline{\mu}^{\ddagger} \neq \underline{\mu}^{\ddagger}$ .  $\square$

*Proof of Theorem 11.* Rewrite the trace term using the identity  $I_J = \underline{u}' + \underline{t}_{\perp}t'_{\perp}$  to get

$$\begin{aligned} \mathcal{T} &= \text{tr} \left\{ (\underline{Y} - \underline{\mu})(\underline{Y} - \underline{\mu})' \right\} = \text{tr} \left\{ (\underline{Y} - \underline{\mu})\underline{u}'(\underline{Y} - \underline{\mu})' \right\} \\ &\quad + \text{tr} \left\{ (\underline{Y} - \underline{\mu})\underline{t}_{\perp}t'_{\perp}(\underline{Y} - \underline{\mu})' \right\}. \end{aligned} \quad (\text{A.15})$$

By (94) then  $\underline{\mu} = \gamma' + \delta'_{\perp}$  so that  $\mathcal{T} = \mathcal{T}_1 + \mathcal{T}_2$  where

$$\begin{aligned} \mathcal{T}_1 &= \text{tr} \left\{ (\underline{Y} - \gamma')\underline{u}'(\underline{Y} - \gamma')' \right\} = t' \text{tr} \left\{ (\underline{Y}I - \gamma)(\underline{Y}I - \gamma)' \right\}, \\ \mathcal{T}_2 &= \text{tr} \left\{ (\underline{Y} - \delta'_{\perp})\underline{t}_{\perp}t'_{\perp}(\underline{Y} - \delta'_{\perp})' \right\} \\ &= \text{tr} \left\{ (\underline{Y}I_{\perp}t'_{\perp} - \delta'_{\perp})(\underline{Y}I_{\perp}t'_{\perp} - \delta'_{\perp})' \right\}. \end{aligned} \quad (\text{A.16})$$

The term  $\mathcal{T}_1$  has a minimum of zero if and only if  $\underline{Y}I = \gamma$ . In  $\mathcal{T}_2$  replace for a moment  $\delta'_{\perp}$  by a matrix  $\phi$  with rank of at most one. The altered  $\mathcal{T}_2$  is minimised when  $\phi$  is the singular value decomposition of  $\underline{Y}I_{\perp}t'_{\perp}$  truncated to rank one due; see Golub and van Loan [46, Theorem 2.5.2]. That singular value decomposition has the property that it is zero when multiplied by  $t$ . Therefore it is also the minimiser of the original problem.  $\square$

*Proof of Theorem 12.* It follows by comparing (105) and (106).  $\square$

*Proof of Theorem 13.* Write  $\tilde{\mu}_{i,J+h,1}(\hat{\theta}_1, \hat{\theta}_2) = \hat{\alpha}_{i1} + \hat{\beta}_{i1}\tilde{\kappa}_{J+h,1}(\hat{\kappa}_1, \hat{\kappa}_2)$ . This is equals to  $\tilde{\mu}_{i,J+h,1}\{g_2(\hat{\theta}_1), g_2(\hat{\theta}_2)\} = (\hat{\alpha}_{i1} + \hat{\beta}_{i1}c_1) + (\hat{\beta}_{i1}/d_1)\tilde{\kappa}_{J+h,1}\{d_1(\hat{\kappa}_1 - c_1), d_2(\hat{\kappa}_2 - c_2)\}$  under the given condition.  $\square$

**A.3. A Further Result on Time Effect Hypotheses.** Consider the restriction  $R'\theta_R = \rho$  of (28). The following result holds.

**Theorem A.3.** Consider the restriction

$$R'\theta_R = \rho \quad (\text{A.17})$$

of (28) where  $A \in \mathbb{R}^{q \times p}$  and  $R \in \mathbb{R}^{q \times (q - q_R)}$  so  $p, q_R \leq q$ . Suppose  $A$  and  $R$  have full column rank. Then, for some  $p_R \leq \min(p, q_R)$  it holds  $p_R = \text{rank}(R'_{\perp}A) = \text{rank}(A, R) - (q - q_R)$ . Then write  $R'_{\perp}A = ab'$  for some matrices  $a \in \mathbb{R}^{q_R \times p_R}$  and  $b \in \mathbb{R}^{p \times p_R}$  with full column rank. The hypothesis (28) restricts the canonical parameter affinely through

$$b'_{\perp}\xi = b'_{\perp}A'\bar{R}\rho, \quad (\text{A.18})$$

so that the degrees of freedom of the restriction is  $p - p_R$ . Introduce the parameters

$$\varphi_1 = a'\bar{R}'_{\perp}\theta_R, \quad \varphi_2 = \bar{b}'A'\xi = \varphi_1 + \bar{b}'A'\bar{R}\rho. \quad (\text{A.19})$$

Then, the predictor  $\mu$  can be written as the  $p_R$ -dimensional affine subspace of the form

$$M_R = \left\{ \mu \in \mathbb{R}^n : \mu = Xb\varphi_1 + XA'\bar{R}\rho \text{ for } \varphi_1 \in \mathbb{R}^{p_R} \right\} \quad (\text{A.20})$$

$$= \left\{ \mu \in \mathbb{R}^n : \mu = Xb\varphi_2 + X\bar{b}_\perp b'_\perp A'\bar{R}\rho \text{ for } \varphi_2 \in \mathbb{R}^{p_R} \right\}. \quad (\text{A.21})$$

The mapping  $\theta_R \mapsto \mu = XA'\theta_R$  on  $\theta_R \in \Theta_R$  is invariant with respect to the group  $g_R : \theta_R \mapsto \theta_R + R_\perp a_\perp \zeta_3$  with  $\varphi_1$  and  $\varphi_2$  as maximal invariants.

A special case arises if the restriction combines a restriction on  $A$  with ad hoc identification. That is, if  $R = (Ad, C)$  where  $d \in \mathbb{R}^{p \times (p-p_R)}$  and  $C \in \mathbb{R}^{q \times (q-p-(q_R-p_R))}$  so  $(A, C)$  has full column rank. Then  $b = d_\perp$  so that the restriction (A.18) reduces to  $d'\xi = (I_{p-p_R}, 0)\rho$ .

*Proof of Theorem A.3.* Apply Lemma A.1 (i, iii) to see that  $\text{rank}(A, R) = p_R + q - q_R$  is equivalent to  $\text{rank}(R'_\perp A) = p_R$ . We can then write  $R'_\perp A = ab'$  for some matrices  $a \in \mathbb{R}^{q_R \times p_R}$  and  $b \in \mathbb{R}^{p \times p_R}$  with full column rank.

Exploit the projection identity  $I_q = R_\perp \bar{R}'_\perp + \bar{R}R'$  to get  $\xi = A'\theta_R = A'R_\perp \bar{R}'_\perp \theta_R + A'\bar{R}R'\theta_R$ . Insert  $A'R_\perp = ba'$  and  $R'\theta_R = \rho$  to get  $\xi = ba'\bar{R}'_\perp \theta_R + A'\bar{R}\rho$  and therefore

$$\theta \mapsto \mu = XA'\theta = Xba'\bar{R}'_\perp \theta_R + XA'\bar{R}\rho. \quad (\text{A.22})$$

Inserting the orthogonal projection  $I_p = b\bar{b}' + \bar{b}_\perp b'_\perp$  this can also be written as

$$\theta \mapsto \mu = XA'\theta = Xb \left( a'\bar{R}'_\perp \theta_R + \bar{b}' A'\bar{R}\rho \right) + X\bar{b}_\perp b'_\perp A'\bar{R}\rho. \quad (\text{A.23})$$

Noting that  $\varphi_1 = a'\bar{R}'_\perp \theta_R$  is freely varying so is  $\varphi_2 = a'\bar{R}'_\perp \theta_R + \bar{b}' A'\bar{R}\rho$ . To rewrite  $\varphi_2$  premultiply the first term by  $I_{p_R} = \bar{b}' b$  to get  $\varphi_2 = \bar{b}' (ba'\bar{R}'_\perp \theta_R + A'\bar{R}\rho)$ . Then insert  $\rho = R'\theta_R$  and  $ba' = A'R_\perp$  to get  $\varphi_2 = \bar{b}' A' (R\bar{R}'_\perp + \bar{R}R') \theta_R = \bar{b}' A'\theta_R = \bar{b}' A'\xi$ . This gives the space  $M_R$  in (A.21).

To derive the reduced group  $g_R$  choose a  $\theta_R \in \Theta_R$  so that  $R'\theta_R = \rho$ . Any  $\theta_R^\dagger \in \Theta$  can be written as  $\theta_R^\dagger = \theta_R + \bar{R}\zeta_1 + R_\perp \bar{a}\zeta_2 + R_\perp a_\perp \zeta_3$  since  $(\bar{R}, R_\perp)$  and  $(\bar{a}, a_\perp)$  have full rank. Since  $\theta_R \in \Theta_R$  then  $\theta_R^\dagger \in \Theta_R$  if and only if  $\zeta_1 = 0$ . We now consider whether  $\theta_R, \theta_R^\dagger$  are equivalent with respect to the restricted mapping (A.22). It holds

$$\theta_R^\dagger \mapsto \mu^\dagger = Xba'\bar{R}'_\perp (\theta_R + R_\perp \bar{a}\zeta_2 + R_\perp a_\perp \zeta_3) + XA'\bar{R}\rho. \quad (\text{A.24})$$

Noting that  $\mu = Xba'\bar{R}'_\perp \theta_R + XA'\bar{R}\rho$  and  $a'\bar{R}'_\perp R_\perp \bar{a} = I_{p_R}$ , while  $a'\bar{R}'_\perp R_\perp a_\perp = 0$ , then  $\mu^\dagger = \mu + Xb\zeta_2$  which reduces to  $\mu$  if and only if  $\zeta_2 = 0$ . Thus, the mapping  $\theta_R \mapsto \mu$  on  $\theta_R \in \Theta_R$  is invariant with respect to the group  $g_R : \theta_R \mapsto \theta_R + R_\perp a_\perp \zeta_3$  with  $\varphi_1$  and  $\varphi_2$  as maximal invariants.

Now, suppose  $R = (Ad, C)$ . Then Lemma A.2 shows that  $m = A'_\perp C$  has full column rank, while  $(A, C)$  has orthogonal complement  $A_\perp m_\perp$ . Thus,  $R$  has orthogonal complement

$$R_\perp = \left\{ \left( \overline{(A, C)} \begin{pmatrix} d_\perp \\ 0 \end{pmatrix}, (A, C)_\perp \right) \right\} \quad (\text{A.25})$$

$$= \left\{ \left( \overline{(A, C)} \begin{pmatrix} I_p \\ 0 \end{pmatrix} d_\perp, A_\perp m_\perp \right) \right\}.$$

In particular, it holds  $A'R_\perp = (I_p, 0)(A, C)'R_\perp = (d_\perp, 0) = d_\perp(I_{p_R}, 0)$ , which is denoted  $ba'$  in the general case, so that  $b = d_\perp$ . Consider the restriction (A.18). Here,  $b'_\perp \xi = d'\xi$ , while  $b'_\perp A' = d'A' = (I_{p-p_R}, 0)(Ad, C)' = (I_{p-p_R}, 0)R'$  so that  $b'_\perp A'\bar{R}\rho = (I_{p-p_R}, 0)\rho$ .  $\square$

### Conflict of Interests

The authors do not have any conflict of interests.

### Acknowledgment

Comments from Andrew Hunt and María Dolores Martínez Miranda are gratefully acknowledged.

### References

- [1] A. Ornelas, M. Guillén, and M. Alcañiz, "Implications of unisex assumptions in the analysis of longevity for insurance portfolios," in *Modeling and Simulation in Engineering, Economics, and Management*, vol. 145 of *Lecture Notes in Business Information Processing*, pp. 99–107, 2011.
- [2] S. F. Jarner and E. M. Kryger, "Modelling adult mortality in small populations: the SAINT model," *ASTIN Bulletin*, vol. 41, no. 2, pp. 377–418, 2011.
- [3] R. D. Lee and L. R. Carter, "Modeling and forecasting U.S. mortality," *Journal of the American Statistical Association*, vol. 87, pp. 659–671, 1992.
- [4] Y. Yang, W. J. Fu, and K. C. Land, "A methodological comparison of age-period-cohort models: the intrinsic estimator and conventional generalized linear models," *Sociological Methodology*, vol. 34, pp. 75–110, 2004.
- [5] Y. Yang and K. C. Land, "Age-period-cohort analysis of repeated cross-section surveys: fixed or random effects?" *Sociological Methods & Research*, vol. 36, no. 3, pp. 297–326, 2008.
- [6] C. Berzuini and D. Clayton, "Bayesian analysis of survival on multiple time scales," *Statistics in Medicine*, vol. 13, no. 8, pp. 823–838, 1994.
- [7] A. Riebler and L. Held, "The analysis of heterogeneous time trends in multivariate age-period-cohort models," *Biostatistics*, vol. 11, no. 1, pp. 57–69, 2010.
- [8] F. Girosi and G. King, *Demographic Forecasting*, Princeton University Press, Princeton, NJ, USA, 2008.
- [9] E. Pitacco, M. Denuit, S. Haberman, and A. Olivieri, *Modelling Longevity Dynamics for Pensions and Annuity Business*, Oxford University Press, Oxford, UK, 2009.
- [10] D. Kuang, B. Nielsen, and J. P. Nielsen, "Forecasting in an extended chain-ladder-type model," *Journal of Risk and Insurance*, vol. 78, no. 2, pp. 345–359, 2011.



- [11] E. Coelho and L. C. Nunes, "Forecasting mortality in the event of a structural change," *Journal of the Royal Statistical Society A: Statistics in Society*, vol. 174, no. 3, pp. 713–736, 2011.
- [12] D. R. Cox and D. V. Hinkley, *Theoretical Statistics*, Chapman & Hall, London, UK, 1974.
- [13] O. Barndorff-Nielsen, *Information and Exponential Families*, John Wiley & Sons, Chichester, UK, 1978.
- [14] S. R. Searle, *Matrix Algebra Useful for Statistics*, John Wiley & Sons, New York, NY, USA, 1982.
- [15] B. Nielsen, "Power of tests for unit roots in the presence of a linear trend," *Oxford Bulletin of Economics and Statistics*, vol. 70, no. 5, pp. 619–644, 2008.
- [16] D. J. Poirier, "Revising beliefs in nonidentified models," *Econometric Theory*, vol. 14, no. 4, pp. 483–509, 1998.
- [17] A. F. M. Smith, "Bayes estimates in one way and two way models," *Biometrika*, vol. 60, no. 2, pp. 319–329, 1973.
- [18] J. M. Bernardo and A. F. M. Smith, *Bayesian Theory*, John Wiley & Sons, Chichester, UK, 2000.
- [19] D. Kuang, B. Nielsen, and J. P. Nielsen, "Identification of the age-period-cohort model and the extended chain-ladder model," *Biometrika*, vol. 95, no. 4, pp. 979–986, 2008.
- [20] N. Keiding, "Statistical inference in the Lexis diagram," *Philosophical Transactions of the Royal Society of London A*, vol. 332, no. 1627, pp. 487–509, 1990.
- [21] S. E. Fienberg and W. M. Mason, "Identification and estimation of age-period-cohort models in the analysis of discrete archival data," *Sociological Methodology*, vol. 10, pp. 1–67, 1979.
- [22] D. Clayton and E. Schifflers, "Models for temporal variation in cancer rates. II: age-period-cohort models," *Statistics in Medicine*, vol. 6, no. 4, pp. 469–481, 1987.
- [23] R Development Core Team, *R: A Language and Environment for Statistical Computing*, R Foundation for Statistical Computing, Vienna, Austria, 2013.
- [24] M. D.M. Miranda, B. Nielsen, and J. P. Nielsen, "Inference and forecasting in the age-period-cohort model with unknown exposure with an application to mesothelioma mortality," *Journal of the Royal Statistical Society A*, 2014.
- [25] D. Kuang, B. Nielsen, and J. P. Nielsen, "Forecasting with the age-period-cohort model and the extended chain-ladder model," *Biometrika*, vol. 95, no. 4, pp. 987–991, 2008.
- [26] P. D. England and R. J. Verrall, "Stochastic claims reserving in general insurance," *British Actuarial Journal*, vol. 8, pp. 519–544, 2002.
- [27] B. Zehnwirth, "Probabilistic development factor models with applications to loss reserve variability, prediction intervals, and risk based capital," in *Proceedings of the Casualty Actuarial Society Forum*, pp. 447–605, Arlington, Va, USA, 1994.
- [28] D. Kuang, B. Nielsen, and J. P. Nielsen, "Chain-ladder as maximum likelihood revisited," *Annals of Actuarial Science*, vol. 4, pp. 105–121, 2009.
- [29] B. Carstensen, "Age-period-cohort models for the Lexis diagram," *Statistics in Medicine*, vol. 26, no. 15, pp. 3018–3045, 2007.
- [30] D. Clayton and E. Schifflers, "Models for temporal variation in cancer rates. I: age-period and age-cohort models," *Statistics in Medicine*, vol. 6, no. 4, pp. 449–467, 1987.
- [31] R. M. O'Brien, "Constrained estimators and age-period-cohort models," *Sociological Methods & Research*, vol. 40, no. 3, pp. 419–452, 2011.
- [32] R. M. O'Brien, "Intrinsic estimators as constrained estimators in age-period-cohort accounting models," *Sociological Methods & Research*, vol. 40, no. 3, pp. 467–470, 2011.
- [33] W. J. Fu, K. C. Land, and Y. Yang, "On the intrinsic estimator and constrained estimators in age-period-cohort models," *Sociological Methods & Research*, vol. 40, no. 3, pp. 453–466, 2011.
- [34] L. L. Kupper, J. M. Janis, A. Karmous, and B. G. Greenberg, "Statistical age-period-cohort analysis: a review and critique," *Journal of Chronic Diseases*, vol. 38, no. 10, pp. 811–830, 1985.
- [35] T. R. Holford, "An alternative approach to statistical age-period-cohort analysis," *Journal of Chronic Diseases*, vol. 38, no. 10, pp. 831–836, 1985.
- [36] A. J. G. Cairns, D. Blake, K. Dowd, G. D. Coughlan, and M. Khalaf-Allah, "Bayesian stochastic mortality modeling for two populations," *ASTIN Bulletin*, vol. 41, no. 1, pp. 29–59, 2011.
- [37] S. Johansen, *Likelihood-Based Inference in Cointegrated Vector Autoregressive Models*, Oxford University Press, Oxford, UK, 1995.
- [38] B. Nielsen, "The likelihood-ratio test for rank in bivariate canonical correlation analysis," *Biometrika*, vol. 86, no. 2, pp. 279–288, 1999.
- [39] B. Nielsen, "Conditional test for rank in bivariate canonical correlation analysis," *Biometrika*, vol. 88, no. 3, pp. 874–880, 2001.
- [40] A. E. Renshaw and S. Haberman, "A cohort-based extension to the Lee-Carter model for mortality reduction factors," *Insurance: Mathematics and Economics*, vol. 38, no. 3, pp. 556–570, 2006.
- [41] A. J. G. Cairns, B. David, K. Dowd et al., "A quantitative comparison of stochastic mortality models using data from England and Wales and the United States," *North American Actuarial Journal*, vol. 13, no. 1, pp. 1–35, 2009.
- [42] C. Pedroza, "A Bayesian forecasting model: predicting U.S. male mortality," *Biostatistics*, vol. 7, no. 4, pp. 530–550, 2006.
- [43] R. F. Engle and C. W. Granger, "Co-integration and error correction: representation, estimation, and testing," *Econometrica*, vol. 55, pp. 251–276, 1987.
- [44] D. F. Hendry and B. Nielsen, *Econometric Modeling*, Princeton University Press, Princeton, NJ, USA, 2007.
- [45] M. A. S. Cabrera, S. M. de Andrade, and R. M. Dip, "Lipids and all-cause mortality among older adults: a 12-year follow-up study," *TheScientificWorldJOURNAL*, vol. 2012, Article ID 930139, 5 pages, 2012.
- [46] G. H. Golub and C. F. van Loan, *Matrix Computations*, Johns Hopkins University Press, Baltimore, Md, USA, 1989.