

Analysis of Mortality Models with Covariates Missing at Random

Shirin Shoaee¹, Reyhaneh Fathi²

Received: 2021/23/04

Accepted: 2021/14/07

Abstract

Objective: Demographic indicators such as mortality rates play a very important role in health, financial and pension policies. Therefore, the accuracy of mathematical models in estimating mortality rates is an important challenge. One of the tasks of actuaries is to construct a suitable mortality model for the available data so that these mortality models can calculate mortality for different ages and longevity, as well as the different information available to individuals on retirement plans. Missing data is a problem that may be faced by actuaries when they are analyzing the real data. Missing data can occur for a variety of reasons, such as unanswered or censored. The presence of missing data can pose a threat to the accuracy of the data analysis results. The purpose of this study is to model the mortality in a retirement plan. In this regard, it is assumed that data are available at the individual level, including date of birth, date of joining the retirement plan, date of completion of the observation, and reason for discontinuation (usually death or right censoring). Information on covariate variables such as gender, benefits or size of pension, demographic geography or health status will also be available. More precisely, this study aims to model the mortality in a retirement plan based on missing data and access to information from various covariate variables, to carefully analyze the structure of different models, to estimate and finally to investigate the financial implications for different mortality experiences containing missing data.

Methodology: In this article, we deal with a pension plan in which each member's future life expectancy is modeled using parametric survival models incorporating covariates which may be missing for some individuals. Likelihood-based techniques estimate parameters, and in this regard, an algorithm is proposed that can perform the estimation task in the best possible way. One of the necessary features to check the adequacy of the statistical model, especially when the data contains missing values, is identifiable. If not identifiable, it can be claimed that the statistical model is not a full rank and is not a suitable model for the data. It is worth noting that the Jacobin matrix needs to be calculated to verify identifiability. As mentioned, in the analysis of

1. Assistant Professor of Department of Actuarial Science, Faculty of Mathematical Sciences, Shahid Beheshti University, Tehran, Iran. (**Corresponding Author**). sh_shoaee@sbu.ac.ir.

2. MSc in Actuarial Science, Department of Actuarial Science, Faculty of Mathematical Sciences, Shahid Beheshti University, Tehran, Iran. reyhane.fathi1@gmail.com

mortality models with the presence of missing values, the maximum likelihood method can be used. In such cases, an estimation error may often occur when fitting the model, which can be reduced by modeling from a larger population. For this reason, hybrid retirement plans that remain homogeneous are often used. This proposed method can also be useful for calculating financial quantities based on pension factors. In fact, in this proposed method, different data sets with equal or similar death experiences are combined, sample size increases and risk of parameter decreases, which also leads to a reduction in capital requirement. Socio-economic variables such as the level of benefits and geographical characteristics of the population are also considered more if interest rates are low.

Finding: First, complete data are analyzed and modeled for observations of members of a retirement plan, which includes survival time and ancillary variables for each individual. Estimation of parameters is obtained using the maximum likelihood method. However, when the data is missing, it is not easy to estimate the parameters with the maximum likelihood method. In this case, the model parameters are estimated by the maximum likelihood method which are calculated using the proposed algorithm; then, statistical indicators such as identifiability of parameters are calculated to evaluate the performance of the proposed structure and algorithm. Furthermore, the financial effects, in particular the annuity factors, and the misestimation risk capital requirements for the mortality experience which includes the maximum covariates variables are calculated and compared with the individual segments when the data are missing. In addition, it can be seen that when the two statistical variables are not observed together, the model is not identifiable according to the data.

Conclusion: It was found that if the data are missing, the statistical model is not always identifiable using the maximum likelihood, and data combination from two or more experiments can avoid identifiable barriers. The methods proposed in this paper can be useful for actuaries when calculating financial committees based on annuity factors. These methods may combine different datasets with equal or similar mortality experiences, increase sample size, and reduce parameter risk, thus, reducing capital requirements. Socio-economic variables such as the level of benefits and geographical characteristics of the population are given more attention if the interest rate is low.

Keywords: Capital Requirement, Parameter Redundant, Full Rank, Likelihood Contribution, Identifiability, Missing At Random, Mortality Model.

JEL-Classification: C13, C24, C51

تجزیه و تحلیل مدل‌های مرگ‌ومیر با متغیرهای کمکی دارای گمشدگی تصادفی

شیرین شعاعی^۱، ریحانه فتحی^۲

تاریخ پذیرش: ۱۴۰۰/۰۳/۲۴

تاریخ دریافت: ۱۴۰۰/۰۲/۰۳

چکیده

هدف: این پژوهش با هدف مدل‌بندی مرگ‌ومیر در یک طرح بازنشتگی بر اساس داده‌های گمشده و دسترسی به اطلاعات مختلف از متغیرهای کمکی، تجزیه و تحلیل دقیق ساختار مدل‌های مختلف، برآوردیابی و در نهایت بررسی تأثیرات مالی برای تجربه‌های مختلف مرگ‌ومیر و حاوی داده‌های گمشده انجام شده است.

روش‌شناسی: این مقاله با یک طرح بازنشتگی سروکار دارد که در آن طول عمر آتی هر فرد با مدل‌های بقای پارامتری با ترکیب متغیرهای کمکی که ممکن است برای برخی از افراد گمشده باشند، مدل‌سازی شده است. پارامترها با روش ماکسیمم درست‌نمایی برآورد شده و الگوریتمی پیشنهاد گردیده که بتواند وظیفه برآوردیابی را به بهترین شکل ممکن انجام دهد.

یافته‌ها: نتایج نشان داد در صورتی که داده‌ها گمشده باشند، مدل آماری همیشه با استفاده از ماکسیمم درست‌نمایی شناسایی‌پذیر نیست و تلفیق داده‌های حاصل از دو یا چند تجربه می‌تواند از موانع شناسایی‌پذیری جلوگیری نماید.

نتیجه‌گیری: روش‌های پیشنهادی این مقاله هنگام محاسبه کمیت‌های مالی مورد علاقه براساس عامل‌های مستمری، برای بیم‌سازان می‌تواند مفید باشد. این روش‌ها ممکن است مجموعه داده‌های مختلف با تجربه مرگ‌ومیر برابر یا مشابه را با هم ترکیب کنند، اندازه نمونه را افزایش دهند و ریسک پارامتر را کاهش دهند، بنابراین منجر به کاهش الزام سرمایه شوند. متغیرهای اقتصادی-اجتماعی از جمله سطح مزایا و مشخصات جغرافیایی جمعیتی در صورت پایین بودن نرخ بهره بیشتر مورد توجه قرار می‌گیرند.

واژگان کلیدی: الزام سرمایه، پارامتر افزونه، رتبه کامل، سهم درست‌نمایی، شناسایی‌پذیری، گمشدگی تصادفی؛ مدل مرگ‌ومیر.

طبقه‌بندی موضوعی: C51. C24. C13

۱. استادیار گروه بیم‌سنجی، دانشکده علوم ریاضی، دانشگاه شهید بهشتی، تهران، ایران. (نویسنده مسئول).

sh_shoae@sbu.ac.ir

۲. کارشناسی ارشد اکچوئری، گروه بیم‌سنجی، دانشکده علوم ریاضی، دانشگاه شهید بهشتی، تهران، ایران.

reyhane.fathi1@gmail.com

مقدمه

شاخص‌های جمعیتی همواره نقش مهمی در تصمیم‌گیری‌های جمعیتی داشته‌اند و نرخ مرگ‌ومیر^۱ مهم‌ترین عامل تغییر الگوی جمعیت به‌شمار می‌آید. داشتن اطلاعات درست از آینده جمعیت برای برآورد دقیق جمعیت سالخورده به‌منظور برنامه‌ریزی‌های دقیق مالی جهت تأمین آینده افراد، ساخت جداول آتی عمر مطابق با محتمل‌ترین الگوی مرگ‌ومیر و قیمت‌گذاری صحیح سالانه‌های عمر، تنها بخشی از کاربردهای پیش‌بینی مقادیر آتی نرخ مرگ‌ومیر در جمعیت‌شناسی، بیمه و آمار است.

مدل‌بندی نرخ مرگ‌ومیر پیشینه‌ای طولانی دارد و جداول مرگ‌ومیر برای بیم‌سنگان، یکی از متداول‌ترین و قدیمی‌ترین روش‌های بررسی میزان مرگ‌ومیر و ابزار اصلی توصیف الگوی مرگ‌ومیر تلقی می‌شوند. رویکرد یک جدول مرگ‌ومیر بر این اساس است که آزمایش گذشته بدون هیچ دگرگونی در آینده تکرار شود. لذا تلاش‌های بسیاری برای دستیابی به مدلی که بتواند این جداول را تبیین کند، صورت گرفته است. یافتن مدل ریاضی برای احتمال مرگ‌ومیر از نظر کاربردی مفیدتر است. روشن است که استفاده از تابعی که پارامتر کمی دارد و از نظر هندسی دارای منحنی همواری است به‌مراتب آسان‌تر از کار با جدول عمری است که پارامترهای فراوانی دارد و در موارد بسیاری نیز احتیاج به درون‌یابی، برون‌یابی، هموارسازی و تعدیل‌های متعدد آماری دیگر، در جهت یافتن احتمال‌های فوت دارد که در محاسبه حق‌بیمه فنی استفاده می‌شوند. لذا بیم‌سنگ‌ها علاقه‌مند هستند که مدل مرگ‌ومیری را ارائه دهند که بتواند برای سنین و طول عمرهای متفاوت و همچنین اطلاعات متفاوتی که برای افراد در طرح‌های بازنشستگی موجود است، شدت مرگ‌ومیر را محاسبه نماید.

بنابراین، مطالعه حاضر با هدف مدل‌بندی مرگ‌ومیر در یک طرح بازنشستگی انجام شده است. می‌توان فرض نمود که داده‌ها در سطح فردی شامل تاریخ تولد، تاریخ پیوستن به طرح، تاریخ اتمام مشاهده و دلیل قطع مشاهده (معمولاً مرگ یا سانسور از راست) در

دسترس هستند. همچنین سایر اطلاعات به صورت متغیرهای کمکی در دسترس خواهند بود. از جمله متغیرهای کمکی، جنسیت است که معمولاً شناخته شده است. متغیرهای کمکی دیگری که شامل اطلاعات غنی تری هستند، از قبیل مزایا یا اندازه حقوق بازنشستگی، مشخصات جغرافیایی جمعیتی یا وضعیت سلامتی نیز در دسترس هستند.

در طرح‌های بازنشستگی اطلاعات متفاوتی برای هر فرد موجود است و وابستگی مرگ‌ومیر به این اطلاعات به خصوص زمانی که ریسک روی زیرمجموعه‌هایی از اعضا متمرکز است، بسیار اهمیت دارد. در بررسی داده‌ها در طرح‌های بازنشستگی با مواردی مواجه می‌شویم که برخی از اطلاعات در مورد هر فرد موجود نیست. در تحلیل مدل‌های مرگ‌ومیر با حضور مقادیر گمشده^۱ می‌توان از روش ماکسیمم درستنمایی^۲ استفاده نمود. در این گونه حالت‌ها اغلب هنگام برازش مدل ممکن است خطای برآورد حاصل شود، که می‌توان با مدل‌سازی از یک جمعیت بزرگ‌تر آن را کاهش داد. به همین علت معمولاً از طرح‌های بازنشستگی ترکیبی که همگن باقی بمانند، استفاده می‌شود. این امر باعث افزایش اندازه نمونه، کاهش ریسک پارامتر و در نتیجه، کاهش الزام سرمایه^۳ می‌گردد.

در این پژوهش بعد از بیان مقدمه، در ادامه، مبانی نظری و پیشینه پژوهش به اختصار مرور شده است. بخش بعدی شامل نمادگذاری مورد نیاز برای توزیع احتمال متغیرهای کمکی، توصیف تابع درستنمایی^۴ و توسیع آن برای زمانی که داده‌ها گمشده هستند، است. همچنین در ابتدای این بخش، به تحلیل و مدل‌بندی داده‌های کامل برای مشاهدات اعضای یک طرح بازنشستگی که شامل زمان بقا و متغیرهای کمکی برای هر فرد است، پرداخته شده است. همچنین در این بخش براساس مشاهدات کامل از داده‌های صندوق بازنشستگی، ابتدا مجموعه داده‌هایی حاوی داده‌های گمشده با حذف متغیرهای کمکی برای برخی واحدها شبیه‌سازی شده‌اند. سپس مسئله برآورد، شرایط ریاضی برای شناسایی پذیری^۵ برآورد پارامترها و الگوریتم برازش توصیف شده است. به علاوه مشاهده

1. Missing Values
2. Maximum Likelihood
3. Capital Requirement
4. Likelihood Function
5. Identifiability

می‌شود در حالتی که دو متغیر آماری توأمآ مشاهده نمی‌شوند، مدل با توجه به داده‌ها شناسایی پذیر نیست. برای غلبه بر عدم شناسایی پذیری، نمونه بسیار کوچکی در نظر گرفته شده که متغیرهای کمکی آن کاملاً مشاهده شده است. در پایان درباره تأثیر مالی ترکیب مجموعه داده‌های مختلف برای برآورد شدت مرگ‌ومیر مطالعه شده و در نهایت نتیجه‌گیری ارائه گردیده است.

۱. مروری بر پیشینه پژوهش

تدوین مدل‌های مرگ‌ومیر از ساختن جداول مرگ‌ومیر آغاز و اولین دستاورد در این زمینه در سال ۱۶۹۳ به وسیله ستاره‌شناس مشهور، هالی^۱ انجام گرفته است که جدول عمری براساس تعداد فوتی‌های مشاهده شده به دست آورد. در سال ۱۷۴۰، اولین جدول عمر به تفکیک مردان و زنان توسط استرویک^۲ منتشر شد. اولین مدل تحلیلی مرگ‌ومیر که احتمال بقا را به صورت یک تابع افزایشی از سن کنونی فرد در نظر گرفت، توسط دمورآور^۳ معرفی شد. سپس گومپرتز^۴ در سال ۱۸۲۵ مدل کامل‌تری از مدل دمورآور ارائه نمود. در سال ۱۸۶۰ مکهام^۵، مدل گومپرتز را که برای سنین بالاتر مناسب نبود اصلاح و مدل کامل‌تری ارائه کرد. بعد از او پرکس^۶ مدل ریاضی ارائه شده گومپرتز و مکهام را کامل نمود. مدل‌های گومپرتز و مکهام از مدل‌های رایج مرگ‌ومیر باقی ماندند تا این که در اوایل قرن بیستم، اقتصاددان و جامعه‌شناس ایتالیایی پارتو^۷ نظریات خود که بر مبنای مشکلات جامعه بود را در مدل‌بندی مرگ‌ومیر دخیل نمود. در اواخر قرن بیستم مدل‌های مرگ‌ومیر پیچیده‌ای بررسی شدند که اکثر آن‌ها تعدیل شده یا تعمیم یافته مدل گومپرتز

1. Halley
2. Struyck
3. De Moivre
4. Gompertz
5. Makeham
6. Perkes
7. Pareto

بودند. از جمله می‌توان به مدل ذکایی و مقصودی (۱۳۸۹) که به بازسازی مدل مرگ‌ومیر بر پایه شکنندگی با استفاده از تعمیم مدل گومپرتز می‌پرداخت، اشاره نمود. برخلاف مدل‌بندی نرخ مرگ‌ومیر، پیش‌بینی نرخ مرگ‌ومیر پیشینه کوتاه‌تری دارد. تا سه دهه قبل، روش‌های مورد استفاده برای پیش‌بینی مرگ‌ومیر نسبتاً ساده و براساس قضاوت‌های ذهنی انجام می‌شدند. اما در سال‌های گذشته، روش‌های پیشرفته‌تری ارائه شده که کارشناسان بیمه، آمار و جمعیت‌شناسان را به استفاده بیشتر از آنها سوق داده است. از میان این روش‌ها، روش لی-کارتز^۱ به دلیل ساختار ساده و کاربرد موفق آن در بسیاری از کشورها، مورد استفاده محققین قرار گرفته است. پیدایش روش لی-کارتز به تغییر الگوی امید به زندگی از سال ۱۹۰۰ در آمریکا برمی‌گردد. گرچه این روش بر اساس داده‌های مرگ‌ومیر آمریکا در سال‌های ۱۹۰۰ تا ۱۹۸۷، طراحی شده بود، اما عملکرد بسیار خوبی در مدل‌بندی و پیش‌بینی بلندمدت نرخ مرگ‌ومیر در بسیاری از کشورهای توسعه‌یافته داشته است. به طوری که کمیجانی و همکاران (۱۳۹۲) به برآورد و پیش‌بینی نرخ مرگ‌ومیر در ایران با استفاده از مدل لی-کارتز پرداخته و عملکرد آن را شرح داده‌اند.

جدیدترین ایده در مطالعه بقا، مرگ‌ومیر تصادفی است که مرگ‌ومیر همانند یک فرایند تصادفی در نظر گرفته می‌شود (یاشین^۲، ۲۰۰۱). اما رویکرد زنجیر مارکوف در مدل‌بندی مرگ‌ومیر به طوری که توزیع زمان مرگ مشخص و دارای توزیع فاز-نوع^۳ باشد به وسیله به‌وسیله لین و لیو^۴ (۲۰۰۷) مطرح شده است. همچنین شجاعی آذر و حسن‌زاده (۱۳۹۳) (۱۳۹۳) براساس این ایده به مدل‌بندی مرگ‌ومیر پرداخته‌اند.

در ادامه پژوهش‌های صورت گرفته در این حوزه، مهدوی و همکاران (۱۳۹۰) به کاربرد یک مدل مرگ‌ومیر با چند عامل ریسک پرداختند. آن‌ها به این موضوع اشاره نمودند که مدل‌های مرگ‌ومیر معرفی شده عمدتاً معطوف به سن و جنس افراد است، اما عوامل

1. Lee - Carter
2. Yashin
3. Phase-Type
4. Lin & Liu

دیگری شامل وضعیت اجتماعی - اقتصادی افراد مانند سطح تحصیلات، میزان درآمد، محل زندگی، شغل، شرایط تأهل و عوامل رفتاری می‌تواند در میزان مرگ و میر تأثیر داشته باشند.

۲. مبانی نظری

در بررسی بسیاری از مجموعه داده‌ها، با داده‌های گمشده روبه‌رو می‌شویم. گمشدگی داده‌ها می‌تواند به‌علل مختلفی مانند پاسخ داده‌نشده یا سانسور شده به‌وجود آید. در مجموعه داده‌ها در موارد پاسخ داده‌نشده، برخی از اطلاعات هر فرد موجود نیست. موارد سانسور شده از جایی به بعد تمامی اطلاعات درباره فرد موجود نیست که می‌تواند به‌دلیل فوت یا هر علت دیگری فرد از مشاهده خارج شود. وجود داده گمشده می‌تواند منجر به تهدیدی برای صحت نتایج حاصل از تحلیل داده‌ها شود. در نظر نگرفتن داده گمشده یا استفاده از روش‌های نامناسب برخورد با این داده‌ها، به‌عنوان مثال جایگذاری داده‌های گمشده با میانگین می‌تواند به مشکلاتی همچون آماسیدگی^۱ در نقطه‌ای مانند میانگین منجر شود. بسیاری از محققین روش‌هایی را برای مقابله با مشکلات داده گمشده ارائه کرده‌اند. اکثر این روش‌ها وابسته به نوع سازوکار گمشدگی داده‌ها هستند و برحسب آن می‌توان روش تحلیل مناسبی را انتخاب نمود. در این راستا، رویین^۲ (۱۹۷۶) سیستم طبقه‌بندی را برای مشاهدات داده گمشده معرفی نمود که به‌طور گسترده مورد استفاده قرار می‌گیرد.

به‌طور معمول سه مکانیسم گمشدگی شامل گمشدگی تصادفی^۳، گمشدگی کاملاً تصادفی^۴ و گمشدگی غیرتصادفی^۵ وجود دارد. در این راستا، بردار پاسخ برای فرد i ام به‌صورت $y_i = (y_{i1}, \dots, y_{ij})$ است. برای هر متغیر پاسخ برداری از متغیرهای نشانگر

1. Inflation
2. Rubin
3. Missing at Random (MAR)
4. Missing Completely at Random (MCAR)
5. Missing not at Random (MNAR)

به صورت $m_i = (m_{i1}, \dots, m_{ij})$ ، که در آن $m_{ij} = 0$ اگر y_{ij} مشاهده شود و $m_{ij} = 1$ وقتی y_{ij} مشاهده نشود وجود دارد. بردار مشاهدات برای هر فرد به صورت y_i^{obs} (مقادیر مشاهده شده) و y_i^{mis} (مقادیر گمشده) تفکیک می‌شود.

در گمشدگی تصادفی داده‌ها به طور تصادفی گمشده هستند. بدین مفهوم که احتمال گمشدن داده‌ها در یک متغیر Y به برخی از متغیرهای اندازه‌گیری شده دیگر وابسته، اما به مقدار گمشده Y وابسته نیست. به عبارت دیگر بین داده‌های گمشده در Y و مقادیر Y پس از جدا کردن سایر متغیرها رابطه‌ای وجود ندارد.

امروزه بسیاری از محققین روش‌هایی را برای مقابله با مشکل داده گمشده ارائه کرده‌اند که اغلب آنها وابسته به نوع سازوکار گمشدگی داده‌ها هستند. از جمله می‌توان به روش‌های ویلکس^۱ (۱۹۳۲) و لرد^۲ (۱۹۵۵) اشاره کرد که بر توزیع نرمال تمرکز کرده‌اند، درحالی‌که در زمینه داده‌های بقا، اسلاچر و جکسون^۳ (۱۹۸۹) یک مدل لوگ خطی برای تجزیه و تحلیل داده‌های بقا سانسور شده همراه با عوامل قطعی در نظر گرفتند. هیرینگ و ابراهیم^۴ (۲۰۰۱) بر مدل خطرات متناسب کاکس^۵ تمرکز کرده‌اند. در این دو مطالعه الگوی کلی داده‌های گمشده مطرح و پارامترها با استفاده از الگوریتم EM^۶ برآورد شده‌اند. لیتل و آن^۷ (۲۰۰۴) یک روش مبتنی بر احتمال قوی MAR را برای داده‌های چند چند متغیره با مقادیر گمشده بر اساس رگرسیون روی اسپلاین‌های نمره‌های تمایل ارائه دادند. تسیاتیس^۸ (۲۰۰۷) برای داده‌های گمشده و سانسور شده یک روش با مدل‌های نیمه پارامتری معرفی نمود. چن و همکاران^۹ (۲۰۱۴) یک مدل بندی برای داده‌های طولی و بقا بر روی داده‌های گمشده و سانسور شده از چپ معرفی نمودند که بر پایه مدل

1. Wilks
2. Lord
3. Schluchter & Jackson
4. Herring & Ibrahim
5. Cox Proportional Hazards Model
6. Expectation Maximization
7. Little & An
8. Tsiatis
9. Chen et al

کاکس، مدل آمیخته خطی و مدل‌های بیزی بود. خو و همکاران^۱ (۲۰۱۷) از معادلات برآورد ساز مورد انتظار برای تولید یک نظریه کلی استفاده نمودند. شایان ذکر است که رویکرد آن‌ها در شرایطی که متغیرهای کمکی و برآمدها پیوسته یا گسسته هستند قابل استفاده است، اما برای برآورد واریانس از روش بوت استرپ^۲ بایستی استفاده نمود. این روش نیمه پارامتری، یک طرح با داده‌های مقطعی و متغیرهای مستقل از زمان را در نظر می‌گیرد. در نهایت اونگولو و همکاران^۳ (۲۰۱۹) به برآورد پارامترهای مدل برای داده‌های داده‌های بقا بر اساس داده‌های گمشده به وسیله اعمال یک توزیع آمیخته برای متغیرهای کمکی و روش بهینه‌سازی مقید پرداختند.

یکی از ویژگی‌های لازم برای کفایت مدل آماری، به خصوص زمانی که داده‌ها حاوی مقادیر گمشده‌اند، شناسایی پذیری است. در صورت عدم شناسایی پذیری، می‌توان ادعا نمود که مدل آماری رتبه کامل نبوده و مدل مناسبی برای داده‌ها نیست. یک مدل رتبه کامل مدلی است که پارامتر افزونه ندارد. به‌طور خاص اگر در یک مدل آماری، داده‌ها گمشده وجود داشته باشد، ممکن است پارامتر افزونه حاصل شود و روی شناسایی پذیری اثر بگذارد. کچپول و مورگان^۴ (۱۹۹۷) در قضیه ۴ مقاله خود نشان دادند دادند اگر یک مدل پارامتر افزونه باشد، به‌صورت محلی شناسایی پذیر نیست. مدل‌های آماری که پارامتر افزونه نیستند، به‌صورت رتبه کامل تعریف می‌شوند. اگر یک مدل رتبه کامل باشد، مقادیر ویژه ماتریس اطلاع تجربی ارزیابی شده توسط برآورد ماکسیمم درست‌نمایی^۵ (MLE) همگی منفی و از صفر دور خواهند بود. برای تشخیص رتبه کامل بودن به ماتریس ژاکوبین نیاز است. اگر رتبه نمادین این ماتریس با تعداد پارامترهای مدل (تعداد ردیف‌های ماتریس ژاکوبین^۶) برابر بود، مدل رتبه کامل^۷ است. اما اگر رتبه

1. Xu et al
2. Bootstrap
3. Ungolo et al
4. Catchpole & Morgan
5. Maximum likelihood Estimation
6. Jacobian Matrix
7. Full Rank

رتبه نمادین ماتریس ژاکوبین کمتر از تعداد پارامترهای مدل حاصل شود، مدل پارامتر افزونه^۱ و در نتیجه شناسایی پذیر نیست. پارامتر افزونه به حالتی گفته می‌شود که بتوان مدل آماری مدنظر را با بردار پارامتری کوچک‌تری نیز تعریف نمود. بنابراین مدل دارای تعداد زیادی پارامتر (بیش از تعداد مورد نیاز) است و در نتیجه شناسایی پذیر نیست.

در طرح‌های بازنشستگی اطلاعات متفاوتی برای هر فرد موجود است و وابستگی مرگ و میر به این اطلاعات به خصوص زمانی که ریسک روی زیرمجموعه‌هایی از اعضا متمرکز است، بسیار اهمیت دارد. در مدل‌بندی داده‌ها در طرح‌های بازنشستگی با مواردی مواجه می‌شویم که برخی از اطلاعات در مورد هر فرد موجود نیست. یکی از رویکردهای موجود در برآورد پارامترهای مدل‌های مرگ و میر با حضور مقادیر گمشده استفاده از روش ماکسیمم درستنمایی است. اما در اغلب موارد به دلیل عدم دسترسی به برخی از داده‌ها در برآورد پارامترهای مدل خطای برآورد حاصل شود، که می‌توان آن‌ها را با مدل‌سازی از یک جمعیت بزرگ‌تر کاهش داد. به همین علت معمولاً از طرح‌های بازنشستگی ترکیبی که همگن باقی بمانند، استفاده می‌شود. این امر باعث افزایش اندازه نمونه، کاهش ریسک پارامتر و در نتیجه، کاهش الزام سرمایه می‌گردد.

۳. روش شناسایی پژوهش

۳-۱. مشاهدات و تابع درستنمایی

در این بخش، ساختار تابع درستنمایی توأم برای برآورد بردار پارامترهای مدل در حالت مشاهدات کامل ارائه و سپس برای مواردی که داده گمشده وجود دارد، توسعه داده شده است.

۳-۱-۱. مشاهدات کامل

در تجزیه و تحلیل رگرسیونی، متغیرهای کمکی اغلب قطعی در نظر گرفته می‌شوند، در این بخش آن‌ها را به عنوان متغیر تصادفی مدل‌بندی می‌کنیم. مجموعه متغیرهای کمکی تصادفی یک فرد با بردار تصادفی p بعدی Z نشان می‌دهیم. مقدار تحقق یافته z ، شدت مرگ و میر فرد $\mu_x(Z; t)$ را تعیین می‌کند. تابع توزیع و تابع چگالی یا احتمال Z ، به ترتیب با $F_Z(z; \zeta)$ و $f_Z(z; \zeta)$ نشان داده می‌شود که ζ یک بردار پارامتری ناشناخته است. فرض کنید n تعداد افراد مشاهده شده در یک دوره زمانی محدود است. نمادهای زیر را تعریف می‌کنیم:

- ✓ x_i سن فرد i ام در آغاز دوره مشاهده است،
- ✓ t_i مقدار تحقق یافته متغیر تصادفی T_i ، بیانگر زمان زندگی شده فرد i ام در طول دوره است،
- ✓ d_i شمارنده‌ای است که مقدار تحقق یافته متغیر تصادفی D_i را نشان می‌دهد؛ اگر فرد i ام فوت کرده باشد مقدار آن یک و اگر بر اثر سانسور از مشاهده خارج شده باشد مقدار آن صفر است،
- ✓ Z_i مقدار تحقق یافته متغیر تصادفی Z_i ، که متغیرهای کمکی فرد i ام را مشخص می‌کند.

در طول یک دوره مشاهدات متناهی برای آن دسته از افرادی که فوت می‌کنند، T_i بیانگر طول عمر باقی مانده برای سن x_i است. برای افراد زنده (سانسور شده از راست)، تمام آن چیزی که می‌دانیم این است که طول عمر باقی مانده آن‌ها بیشتر از T_i خواهد بود.

فرض می‌کنیم متغیرهای کمکی Z_1, \dots, Z_n مستقل و هم توزیع $(i.i.d.)$ و روی مقادیر تحقق یافته خود شرطی هستند. به علاوه سنین مشاهده شده x_1, \dots, x_n و طول عمرهای مشاهده شده T_1, \dots, T_n به صورت $i.i.d.$ هستند. در نهایت فرض می‌شود که سانسورها

به صورت ناآگاهی بخش هستند. یعنی می‌توان طول عمرها را به عنوان متغیر تصادفی در نظر و توزیع آن را مستقل از زمان سانسور، برآورد نمود، (مک‌دونالد^۱ و همکاران، ۲۰۱۸). در ادامه فرض می‌کنیم مشاهدات به صورت کامل در دسترس هستند، به طوری که هر بردار متغیر کمکی Z_i به طور کامل مشاهده شده است. اگر سهم فرد i ام در تابع درستنمایی را با L_i نمایش دهیم، بنابراین تابع درستنمایی کل به صورت $L = \prod L_i$ خواهد بود.

برای افرادی که فوت می‌شوند ($d_i = 1$) سهم تابع درستنمایی به وسیله ضرب چگالی شرطی T_i روی x_i و Z_i مربوطه در تابع چگالی یا احتمال Z_i به دست می‌آید که به صورت زیر محاسبه می‌گردد:

$$L_i(\tau, \zeta) = f_{x_i}(t_i | z_i; \tau) f_Z(z_i; \zeta), \quad (1)$$

که در آن (τ, ζ) بردار پارامترهای نامعلوم هستند و در رابطه فوق چگالی توام (T, Z) به عنوان مدل کامل^۲ (برای سن x) تعریف می‌شود، و همچنین به چگالی T شرطی شده روی Z مدل جزئی^۳ گفته می‌شود. عبارت $f_{x_i}(t_i | z_i; \tau)$ به این معنی است که فرد x ساله بعد از سن $x + t$ ، دیگر زنده نیست.

برای افرادی که طول عمرهای مشاهده شده آن‌ها سانسور شده از راست است ($d_i = 0$)، زمان دقیق فوت مشاهده نشده است. لذا سهم تابع درستنمایی برای این افراد به صورت زیر محاسبه می‌شود:

$$L_i(\tau, \zeta) = S_{x_i}(t_i, z_i; \tau) f_Z(z_i; \zeta), \quad (2)$$

که $S_{x_i}(t_i, z_i; \tau)$ تابع بقا و بیانگر این است که فرد x ساله حداقل t سال زنده مانده است و داریم:

$$S_{x_i}(t_i, z_i; \tau) = \exp \left\{ - \int_0^{t_i} \mu_{x_i+s}(z_i; \tau) ds \right\}, \quad (3)$$

1. Macdonald
2. Full Model
3. Partial Model

که μ_{x_i+s} بیانگر شدت مرگ و میر فرد i ام دارای سن $x + s$ است. با تلفیق این دو رابطه و ارتباط شدت مرگ و میر، تابع احتمال و بقا، سهم مشارکت در تابع درستنمایی فرد i ام به صورت زیر است:

$$L_i(\tau, \zeta) = \exp\left\{-\int_0^{\tau} \mu_{x_i+s}(z_i; \tau) ds\right\} \mu_{x_i+\tau}(z_i; \tau)^{d_i} f_Z(z_i; \zeta). \quad (4)$$

اگر عناصر Z برای تمام افراد مشاهده شود، تابع درستنمایی به یک تابع از τ و یک تابع از ζ تبدیل می‌شود:

$$L(\tau, \zeta) = \underbrace{\prod_i \exp\left\{-\int_0^{\tau} \mu_{x_i+s}(z_i; \tau) ds\right\} \mu_{x_i+\tau}(z_i; \tau)^{d_i}}_{L^T(\tau)} \underbrace{\prod_i f_Z(z_i; \zeta)}_{L^Z(\zeta)}. \quad (5)$$

در این حالت τ می‌تواند به صورت مستقل از ζ برآورد و توزیع متغیرهای کمکی در تحلیل رگرسیونی نادیده گرفته می‌شود. در بخش بعد مشاهده می‌شود که اگر بعضی از مشاهدات Z_i گمشده باشند، این کار امکان‌پذیر نیست و پارامترهای τ و ζ باید توأمآ برآورد شوند.

۳-۱-۲. مشاهدات گمشده

در بررسی بسیاری از مجموعه داده‌ها، با داده‌های گمشده روبرو می‌شویم. گمشدگی داده‌ها می‌تواند به علل مختلفی مانند پاسخ داده نشده یا سانسور شده به وجود آید. حال فرض کنید بعضی از مولفه‌های متغیرهای کمکی Z برای برخی از افراد تحت مطالعه گمشده هستند. لذا بردار p بعدی $Z_i = (Z_{i,p}, \dots, Z_{i,1})$ به دو بردار Z_i^{obs} و Z_i^{mis} تقسیم می‌شود. برای هر فرد تابع چگالی بردار متغیرهای کمکی $Z_i = (Z_i^{obs}, Z_i^{mis})$ به صورت زیر تجزیه می‌شود:

$$f_{Z_i}(z_i) = f_{Z_i^{obs}}(z_i^{obs}; \zeta) f_{Z_i^{mis}|Z_i^{obs}}(z_i^{mis}|z_i^{obs}; \zeta). \quad (6)$$

متغیرهای کمکی دارای گمشدگی تصادفی هستند، لذا عبارت $f_{Z_i^{mis}|Z_i^{obs}}(z_i^{mis}|z_i^{obs}; \zeta)$ را می‌توان نادیده گرفت. مشابه رابطه (۴)، میزان سهم

درست‌نمایی برای هر فرد، بر پایه طول عمر مشاهده شده T ، سن x و مجموعه متغیرهای کمکی مشاهده شده Z^{obs} به صورت زیر است:

$$L_i(\tau, \zeta) = \int_{S_{z_i^{mis}}} \exp \left\{ - \int_0^{\tau} \mu_{x_i+s}(z_i^{obs}, y; \tau) ds \right\} \mu_{x_i+\tau}^{di}(z_i^{obs}, y; \tau) f_Z(z_i^{obs}, y; \zeta) dy, \quad (7)$$

که روی فضای وضعیت $S_{z_i^{mis}}$ از Z_i^{mis} انتگرال گرفته می‌شود. اگر مولفه‌های Z_i^{mis} دارای فضای وضعیت گسسته باشند، انتگرال روی $S_{z_i^{mis}}$ تبدیل به مجموع خواهد شد. در حالت داده‌های گمشده، تابع درست‌نمایی کل را نمی‌توان مشابه رابطه (۵) به دو تابع جداگانه از τ و ζ عامل‌بندی نمود. بنابراین نمی‌توان τ را به طور مستقل از ζ برآورد کرد و بایستی آن‌ها را به صورت توأم برپایه طول عمر مشاهده شده و متغیرهای کمکی برآورد نماییم.

همچنین زمانی که داده‌ها گمشده و پارامترهای τ و ζ هر دو نامعلوم هستند، شناسایی‌پذیری مدل کامل را نمی‌توان تصدیق نمود.

۴. تجزیه و تحلیل داده‌ها

۴-۱. تحلیل و مدل‌بندی بر اساس داده‌های کامل

هدف این بخش مدل‌بندی داده‌ها برای مشاهدات اعضای یک صندوق بازنشستگی است. برای هر عضو صندوق، دو متغیر کمکی وجود دارد. این مجموعه داده کامل نامیده می‌شوند.

۴-۱-۱. مجموعه داده‌های کامل

مجموعه داده‌های کامل دارای مشخصه‌های زیر است:

✓ ۱۸۷۴۱ مستمری در حال پرداخت است،

- ✓ مجموع سال‌های در معرض ریسک افراد برابر $4/172601$ سال است،
 - ✓ تعداد افراد فوت شده ۴۹۵۶ نفر است،
 - ✓ دوره مشاهده از ۱۰ نوامبر ۱۹۹۲ تا ۳۱ دسامبر ۲۰۰۹ است،
 - ✓ مقدار مزایای سالانه برای هر فرد مشاهده شده را با B نشان می‌دهیم که یک متغیر کمکی با دو سطح مزایای بالا و مزایای پایین است،
 - ✓ مشخصات جغرافیایی جمعیتی برای هر فرد بر پایه سیستم طبقه‌بندی موزاییک^۱، را با C نمایش می‌دهیم و شامل سه سطح ۰ و ۱ و ۲ است. سطح ۰ نماینده محروم‌ترین مناطق و سطح ۲ مناطقی با کم‌ترین محرومیت است.
- دو ایتم آخر، بردار متغیرهای کمکی Z_i را تشکیل می‌دهند. در این مجموعه داده مزایا بر حسب پوند است و می‌توان با آن به‌عنوان متغیر پیوسته برخورد نمود، اما به چند دلیل آن را به صورت متغیر گسسته لحاظ می‌نماییم. اول: مکدونالد و همکاران (۲۰۱۸) نشان دادند که استفاده مستقیم آن در مدل مناسب نیست. دوم: اعمال متغیر پیوسته به‌عنوان یک متغیر طبقه‌بندی شده، از لحاظ آمساک در رویکرد مدل‌بندی راحت‌تر است. سوم: مادریگال^۲ و همکاران (۲۰۱۱) مشاهده کردند که مقدار پایین حقوق بازنشستگی می‌تواند یک شاخص گمراه‌کننده برای ثروتمندی فرد بوده که می‌تواند به این دلیل رخ دهد که یک کارمند (یا مستخدم، عضو) برای خدمات طولانی مدت حقوق پایینی داشته یا برای خدمات کوتاه مدت حقوق بالایی داشته است.

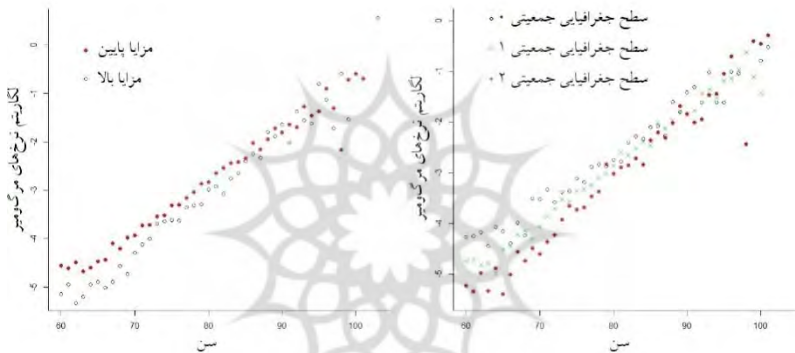
۴-۱-۲. تجزیه و تحلیل مقدماتی داده‌ها

در این بخش ویژگی‌های مجموعه داده کامل را توصیف می‌کنیم. ابتدا نرخ خام مرگ و میر D_x/E_x در هر سال واحد از سن محاسبه و در شکل ۱ نمایش داده شده است. نرخ‌های مرگ و میر به‌طور جداگانه برای هر دو سطح مزایا (نمودار سمت چپ) و سه گروه مشخصات جغرافیایی جمعیتی (نمودار سمت راست) رسم شده است. در این نمودار D_x

1. Mosaic
2. Madrigal

بیانگر تعداد افرادی است که بین سن‌های دقیق x و $x + 1$ فوت شده‌اند و E_x (در معرض ریسک^۱) بیانگر مجموع تمام زمان‌هایی است که تمام افراد بین این دو سن زنده هستند. براساس شکل (۱) سه نتیجه مهم حاصل می‌شود:

- ✓ تقریباً یک رابطه لگ-خطی بین سن و نرخ‌های مرگ و میر وجود دارد.
- ✓ بین مزایای بالاتر و مرگ و میر پایین‌تر ارتباط وجود دارد.
- ✓ بین مشخصات جغرافیایی جمعیتی دارای کم‌ترین محرومیت، با مرگ و میر پایین‌تر ارتباط وجود دارد.



شکل ۱. نمودار خام نرخ مرگ و میر D_x/E_x بر اساس سن

منبع: یافته‌های تحقیق

پژوهشگاه علوم انسانی و مطالعات فرهنگی
رتال جامع علوم انسانی

۳-۱-۴. انتخاب مدل پارامتری

برای مدل‌بندی و به‌منظور ساده نگه داشتن مدل، از اختلاف کوچک بین نرخ مرگ و میر برای گروه‌های مختلف در سنین بالا چشم‌پوشی و از مدل گومپرتز (گومپرتز، (۱۸۲۵)) با جداسازی گروه‌های خاص استفاده می‌کنیم که در آن نیروی مرگ و میر برای فرد i ام از رابطه زیر حاصل می‌شود:

$$\mu(b_i, c_i; \tau) = \exp\{\alpha + \beta x_i + \gamma \mathbb{1}_{[b_i=High]} + \delta_1 \mathbb{1}_{[c_i=1]} + \delta_2 \mathbb{1}_{[c_i=2]}\}, \quad (۸)$$

که $\tau = (\alpha, \beta, \gamma, \delta_1, \delta_2)$ بردار پارامترهای نامعلوم، پارامتر β تأثیر سن، γ تأثیر مزایای بالا و δ_1 و δ_2 تأثیر مشخصات جغرافیایی جمعیتی مختلف است. تابع شدت مرگ و میر ارائه شده در رابطه ۸ متناظر با مدل بقای $S_{x_i}(t_i | b_i, c_i; \tau)$ است که از رابطه زیر محاسبه می‌گردد:

$$\begin{aligned} S_{x_i}(t_i | b_i, c_i; \tau) &= \exp \left\{ - \int_0^{t_i} \mu_{x+s} ds \right\} \\ &= \exp \left\{ - \left(\frac{\exp(\beta t_i) - 1}{\beta} \right) \exp(\alpha + \beta x_i + \gamma \mathbb{1}_{[b_i=High]} \right. \\ &\quad \left. + \delta_1 \mathbb{1}_{[c_i=1]} + \delta_2 \mathbb{1}_{[c_i=2]} \right\}. \end{aligned} \quad (9)$$

تابع چگالی برای طول عمر باقی مانده به صورت زیر حاصل می‌شود:

$$f_{x_i}(t_i | b_i, c_i; \tau) = S_{x_i}(t_i | b_i, c_i; \tau) \mu_{x_i+t_i}(b_i, c_i; \tau). \quad (10)$$

طبق رابطه (۸) پارامترهای γ ، δ_1 و δ_2 بیانگر تفاوت نرخ مرگ و میر بین گروه‌های متفاوت است. مرگ و میر پایه $\exp(\alpha + \beta x)$ در مدل، بیانگر نرخ مرگ و میر برای فرد x ساله که دارای مزایا پایین و در مناطق جغرافیایی جمعیتی با محرومیت بیشتر زندگی می‌کند، است.

۴-۱-۴. نتایج تجربی برای مجموعه داده‌های کامل

در ادامه به بررسی بهترین برازش مدل در رابطه ۸ بر اساس مدل‌های آشیانه‌ای با ترکیب‌های مناسب از پارامترهای γ ، δ_1 و δ_2 می‌پردازیم. در این راستا ۴ مدل مختلف حاصل می‌شود. ساختار این مدل‌ها بر اساس نیروی مرگ و میر در جدول زیر ارائه شده است:

جدول ۱. ساختار ۴ مدل بر اساس نیروی مرگ و میر

ساختار مدل	مدل
$\exp\{\alpha + \beta x_i\};$	M_0
$\exp\{\alpha + \beta x_i + \gamma \mathbb{1}_{[b_i=High]}\};$	M_1
$\exp\{\alpha + \beta x_i + \delta_1 \mathbb{1}_{[c_i=1]} + \delta_2 \mathbb{1}_{[c_i=2]}\};$	M_2
$\exp\{\alpha + \beta x_i + \gamma \mathbb{1}_{[b_i=High]} + \delta_1 \mathbb{1}_{[c_i=1]} + \delta_2 \mathbb{1}_{[c_i=2]}\}.$	M_3

منبع: یافته‌های تحقیق

در جدول ۱، مدل M_0 یک مدل گومپرتز بدون متغیر کمکی، مدل‌های M_1 و M_2 تنها دارای یک متغیر کمکی (به ترتیب مقدار مزایا و مشخصات جغرافیایی جمعیتی) و مدل M_3 هر دو متغیر کمکی را دارد. برای مجموعه داده کامل، تابع درستنمایی به صورت رابطه (۵) است. لذا تنها برآورد τ نیاز است، پس کافی است فقط عامل $L^T(\tau)$ را در نظر گرفته و $L^Z(\zeta)$ را نادیده بگیریم. به عنوان مثال رابطه $L^T(\tau)$ برای مدل M_3 به صورت زیر است:

$$L^T(\tau) = \prod_i \exp \left\{ - \int_0^{t_i} \exp\{\alpha + \beta(x_i + s) + \gamma \mathbb{1}_{[b_i=High]} + \delta_1 \mathbb{1}_{[c_i=1]} + \delta_2 \mathbb{1}_{[c_i=2]}\} ds \right\} \quad (11)$$

$$\times \exp\{\alpha + \beta(x_i + t_i) + \gamma \mathbb{1}_{[b_i=High]} + \delta_1 \mathbb{1}_{[c_i=1]} + \delta_2 \mathbb{1}_{[c_i=2]}\}^{d_i},$$

لذا لگاریتم درستنمایی به صورت زیر حاصل می‌شود:

$$l(\tau) = - \sum_i \exp \left\{ - \left(\frac{\exp(\beta t_i) - 1}{\beta} \right) \exp(\alpha + \beta x_i + \gamma \mathbb{1}_{[b_i=High]} + \delta_1 \mathbb{1}_{[c_i=1]} + \delta_2 \mathbb{1}_{[c_i=2]}) \right\} \quad (12)$$

$$+ \sum_{d_i=1} (\alpha + \beta(x_i + t_i) + \gamma \mathbb{1}_{[b_i=High]} + \delta_1 \mathbb{1}_{[c_i=1]} + \delta_2 \mathbb{1}_{[c_i=2]}).$$

تابع درستنمایی برای سایر مدل‌ها به طریق مشابه حاصل می‌گردد. برای بررسی نیکویی برازش مدل‌ها، معیار اطلاع بیزی^۱ (BIC) را محاسبه و آن‌ها با یکدیگر مقایسه می‌کنیم. این نتایج برای مجموعه داده‌های کامل در جدول ۲ ارائه شده است.

جدول ۲. برآورد پارامترها، لگاریتم درستنمایی، BIC و برآورد خطای استاندارد برآورد پارامترها.

مدل	$\hat{\alpha}$	$\hat{\beta}$	$\hat{\gamma}$	$\hat{\delta}_1$	$\hat{\delta}_2$	$L^T(\tau)$	BIC
M_0	-۱۱/۵۸	۰/۱۱	-	-	-	-۲۰۵۹۲/۶۱	۴۱۲۰۴/۹۰
$\hat{\sigma}(\hat{\theta}) \times 10^3$	۹۹/۱۳۰	۶۷/۱	-	-	-	-	-
M_1	-۱۱/۵۴	۰/۱۱	-۰/۲۸	-	-	-۲۰۵۶۱/۸۸	۴۱۱۵۳/۲۷
$\hat{\sigma}(\hat{\theta}) \times 10^3$	۵۱/۱۳۱	۶۸/۱	۶۳/۱۳	-	-	-	-
M_2	-۱۱/۴۲	۰/۱۱	-	-۰/۲۴	-۰/۴۸	-۲۰۵۲۸/۹۳	۴۱۰۹۷/۲۱
$\hat{\sigma}(\hat{\theta}) \times 10^3$	۹۷/۱۳۲	۶۸/۱	-	۵۷/۳۳	۰۴/۴۳	-	-
M_3	-۱۱/۴۱	۰/۱۱	-۰/۲۱	-۰/۲۲	-۰/۴۳	-۲۰۵۱۲/۶۹	۴۱۰۷۴/۵۷
$\hat{\sigma}(\hat{\theta}) \times 10^3$	۲۵/۱۳۳	۶۹/۱	۸۵/۳۶	۷۵/۳۳	۸۹/۴۳	-	-

منبع: یافته‌های تحقیق

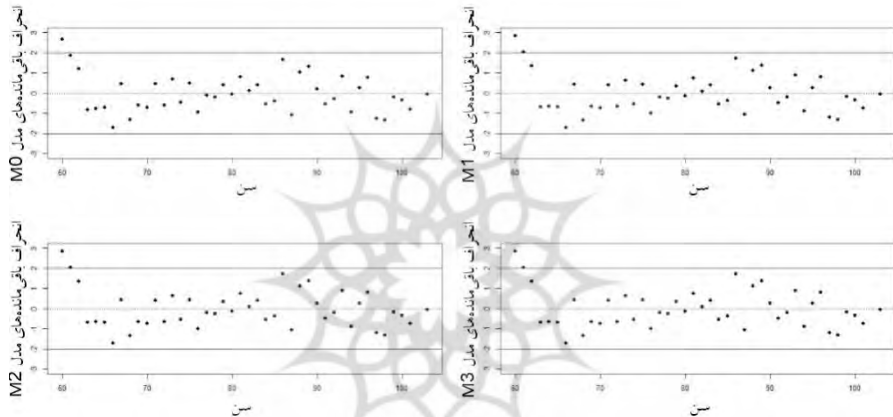
بر اساس نتایج جدول ۲، مدل‌های حاوی متغیرهای کمکی B و C (به صورت جداگانه یا توأم)، برازش مدل را بهبود می‌بخشند. کمترین مقدار BIC برای مدل M_3 است که هر دو متغیر کمکی را دارد. بنابراین M_3 مدل بهتری برای داده‌های موجود است. برای بررسی کیفیت برازش مدل‌ها، باقی‌مانده انحراف پواسون به صورت رابطه زیر محاسبه می‌شود:

$$r_i = \text{sign}(d_i - \Lambda_i) \sqrt{2 \left[d_i \log \frac{d_i}{\Lambda_i} - (d_i - \Lambda_i) \right]}, \quad (13)$$

که d_i مجموع تعداد مرگ و میرها در سن i ام و $\Lambda_i = \sum_j \lambda_{i,j}$ است. پارامتر توزیع پواسون λ برای مدل M_1 و $i = 1, \dots, 44$ و $i = 1, \dots, 18741$ از رابطه زیر محاسبه می‌شود:

$$\lambda_{i,j} = \left(\frac{\exp(\beta \min((60 + i - x_{i,j}), (la_{i,j} - x_{i,j}))) - 1}{\beta} \right) \times \exp(\alpha + \beta x_{i,j} + \gamma \mathbb{1}_{[b_{i,j}=High]}), \quad (14)$$

که $\Lambda_i = \sum_{j=1}^{18741} \lambda_{i,j}$ و $la_{i,j} = x_{i,j} + t_{i,j}$ باقی مانده انحراف پواسون برای همه مدل‌ها در شکل ۲ رسم شده است. مشاهده می‌شود که باقی مانده‌ها در اطراف صفر متمرکز و از مدل خاصی تبعیت نمی‌کنند.



شکل ۲. نمودار باقی مانده انحراف پواسون در برابر سن برای هر ۴ مدل

منبع: یافته‌های تحقیق

پژوهشگاه علوم انسانی و مطالعات فرهنگی
پرتال جامع علوم انسانی

با توجه به مقادیر ویژه ماتریس هسین در برآورد پارامترها که همگی منفی و از صفر دور هستند، لذا تمامی مدل‌ها شناسایی پذیر خواهند بود. همچنین با توجه به مجموعه داده‌های کامل، توزیع توام متغیرهای کمکی B و C به وسیله فراوانی نسبی آن‌ها که در جدول ۳ گزارش شده، برآورد شده است. این مقادیر همان مقادیر بردار پارامتری γ هستند.

جدول ۳. فراوانی نسبی سطح مزایا و مشخصات جغرافیایی جمعیتی برای طرح بازنشستگی

	مشخصات جغرافیایی جمعیتی (C)			مشخصات جغرافیایی جمعیتی (C) سطح مزایا (B)
	۲	۱	۰	
۰/۷۵	۰/۱۳	۰/۴۳	۰/۱۹	پایین
۰/۲۵	۰/۰۹	۰/۱۳	۰/۰۳	بالا
۱	۰/۲۲	۰/۵۶	۰/۲۲	

منبع: یافته‌های تحقیق

۴-۲. تحلیل و مدل‌بندی بر اساس مشاهدات گمشده

در این بخش با استفاده از حذف برخی از اطلاعات افراد، داده گمشده را شبیه‌سازی می‌کنیم. برای این منظور مجموعه داده کامل را به‌طور تصادفی به چند بخش تقسیم و یکی از متغیرهای کمکی یا هر دو آن‌ها را در هر زیر مجموعه داده کامل حذف می‌کنیم.

۴-۲-۱. تولید مجموعه داده حاوی مشاهدات گمشده

برای تولید مجموعه داده‌های حاوی مشاهدات گمشده، مجموعه داده کامل را که شامل n فرد است، به‌صورت تصادفی به دو زیرمجموعه تقسیم و آن‌ها را با P_1 و P_2 نمایش می‌دهیم که هر کدام به ترتیب شامل $n_1 = 9370$ و $n_2 = 9371$ عضو هستند ($n_1 + n_2 = n$). فرض می‌کنیم برای هر فرد در P_1 تنها متغیر کمکی B را مشاهده کردیم (C را حذف کردیم) و برای هر فرد در P_2 تنها متغیر کمکی C را مشاهده کردیم (B را حذف کردیم). این دو زیرمجموعه را به‌عنوان دو طرح بازنشستگی جداگانه در نظر می‌گیریم که علاقه‌مند به مدل‌بندی آن‌ها به‌صورت توأم هستیم.

از آنجائی که مجموعه داده ترکیبی توسط یک منبع تولید شده (مجموعه داده کامل)، لذا هیچ ناهمگنی بین داده‌ها وجود ندارد و می‌توان یک قانون مرگ‌ومیر یکسان برای اعضای P_1 و P_2 در نظر گرفت.

به دلیل این که متغیرهای گمشده رسته‌ای هستند، انتگرال در رابطه ۷ به یک مجموع متناهی روی فضای نمونه مشاهدات گمشده تبدیل می‌شود. برای هر فرد در P_1 که مشخصات جغرافیایی جمعیتی مشاهده نشده است، سهم تابع درستنمایی به صورت زیر است:

$$L_i(\tau, \zeta) \propto \sum_{c \in \{0,1,2\}} S_{x_i}(t_i | b_i, c; \tau) \mu_{x_i+t_i}^{d_i}(b_i, c; \tau) f_{B,C}(b_i, c; \zeta). \quad (15)$$

و برای اعضای P_2 که متغیر مزایا مشاهده نشده است، سهم تابع درستنمایی به فرم زیر است:

$$L_i(\tau, \zeta) \propto \sum_{b \in \{Low, High\}} S_{x_i}(t_i | b, c_i; \tau) \mu_{x_i+t_i}^{d_i}(b, c_i; \tau) f_{B,C}(b, c_i; \zeta). \quad (16)$$

بنابراین تابع درستنمایی به صورت زیر است:

$$L(\tau, \zeta | t, x, b_{obs}, c_{obs}, d) \propto \prod_{\{i \in P_1\}} L_i \prod_{\{i \in P_2\}} L_i. \quad (17)$$

هدف برآورد توأم پارامتر τ از تابع مرگ و میر، پارامتر ζ از تابع احتمال توأم برای سطوح مزایا و مشخصات جغرافیایی جمعیتی است. یک روش منطقی ساختن یک توزیع توأم مناسب برای B و C است.

۴-۲-۲. شناسایی پذیری

حال دو شرط ریاضی که برای شناسایی پذیری بردار پارامتری $\theta = (\tau, \zeta)$ لازم است را ارائه می‌کنیم. فرض می‌کنیم که هیچ کدام از متغیرهای کمکی B و C برای یک فرد مشاهده نشده‌اند. سهم درستنمایی توسط توزیع آمیخته متناهی که اجزای آن تمام نتایج ممکن (B, C) است، ارائه می‌شود. براساس رابطه (۷) تابع چگالی احتمال آمیخته با ۶ مولفه، را می‌توان به صورت زیر نوشت:

$$\begin{aligned}
 f_x(t; \tau, \zeta) = & \zeta_1 f_x(t|b = Low, c = 0; \tau) + \zeta_2 f_x(t|b = High, c = 0; \tau) \\
 & + \zeta_3 f_x(t|b = Low, c = 1; \tau) \\
 & + \zeta_4 f_x(t|b = High, c = 1; \tau) \\
 & + \zeta_5 f_x(t|b = Low, c = 2; \tau) \\
 & + \zeta_6 f_x(t|b = High, c = 2; \tau).
 \end{aligned}
 \tag{18}$$

شروط زیر برای شناسایی پذیری مدل آمیخته متناهی لازم است (مکلاکن و پیل^۱، ۲۰۰۰):

۱. برای سنین مختلف با طول عمرهای متفاوت، مقادیر متفاوت (B,C) باید تابع چگالی احتمال متفاوتی را نتیجه دهند. یعنی برای $h \neq k$ باید داشته باشیم:

$$f_x(t|(b, c) = h; \tau) \neq f_x(t|(b, c) = k; \tau). \tag{19}$$

که این شرط در مورد تابع بقا، به این معنا است که مقادیر متفاوت (B,C) تابع مرگ و میر متفاوتی را بایستی نتیجه دهند.

۲. مقادیر ζ_j برای $j = 1, \dots, 6$ بایستی به صورت $0 < \zeta_j < 1$ باشد.

شایان ذکر است که شرطها لازم هستند، اما کافی نیستند، (تیتراگتون و همکاران^۲، ۱۹۸۵). اگر دو شرط بالا برقرار شد، مدل آمیخته در رابطه ۱۸ شناسایی پذیر است.

زمانی که داده گمشده وجود دارد، شناسایی پذیری مدل را از روشهای دیگری نیز می توان بررسی نمود. هنگام برآورد پارامترهای مدل ممکن است برآوردهایی به دست آیند که دارای واریانسهای بسیار بزرگ و یا بسیار کوچک و نزدیک به صفر هستند و یا ممکن است برآوردها در فضای پارامتری موردنظر، تعریف نشوند، که می توان نتیجه گرفت مدل شناسایی پذیر نیست. همچنین شناسایی پذیری را می توان از طریق مقادیر ویژه ماتریس هسین^۳ (حاصل شده در مرحله برآوردیابی) نیز بررسی نمود. در صورتی که مقادیر ویژه^۴، همگی منفی و از صفر دور به دست آیند، نتیجه می شود مدل آماری شناسایی پذیر است.

1. McLachlan & Peel
2. Titterington et al
3. Hessian Matrix
4. Eigenvalues

۴-۲-۳. الگوریتم برازش

زمانی که داده‌ها گمشده هستند، به حداکثر رساندن تابع درستنمایی ممکن است به آسانی صورت نگیرد. مسئله کلی محاسباتی این است که چگونه محدودیت‌ها روی بردار پارامتری ζ را با توجه به توزیع توام B و C ، کنترل نماییم (عناصر ζ بین صفر و یک محدود شده‌اند و مجموع آن‌ها برابر یک است). به همین دلایل هر بهینه‌سازی مقید معمولی ممکن است شامل عدم پایداری و استواری شود.

الگوریتم EM که ابتدا توسط دمپستر و همکاران^۱ (۱۹۹۷) ارائه شد، برآورد را برای برخی مسائل آسان می‌نمود. در ادامه یک اصلاح جدید برای تبدیل مسئله بهینه‌سازی مقید به بهینه‌سازی بدون قید که منجر به همگرایی سریع‌تر نسبت به الگوریتم EM و نیز استواری برآوردها نسبت به مقادیر اولیه متفاوت می‌شود، ارائه می‌شود. برای تبدیل مسئله بهینه‌سازی مقید به یک مسئله بدون قید، از تبدیل نسبت لگاریتمی ایزومتریک^۲ (ILR) که توسط اگزکو و همکاران^۳ (۲۰۰۳) ارائه شد، استفاده می‌شود. این تبدیل حاوی نگاشت حفظ فاصله از k بعدی ساده به فضای \mathbb{R}^{k-1} بعدی است و بیشتر در تجزیه و تحلیل داده‌های ترکیبی برای پارامترسازی مجدد بردار $\zeta = (\zeta_1, \dots, \zeta_k)$ استفاده می‌شود.

$$\pi = ILR(\zeta) = \Psi^T \log(\zeta) \quad (20)$$

که π یک بردار $k-1$ بعدی و Ψ یک ماتریس با ابعاد $(k \times (k-1))$ است که ستون‌های $(\psi_1, \dots, \psi_{k-1})$ پایه‌های یک‌معامد را برای ابرصفحه از \mathbb{R}^k متعامد به بردار k واحدی ۱ نشان را می‌دهند که به صورت زیر تعریف می‌شود:

$$\psi_i = \sqrt{\frac{i}{i+1}} \left[\underbrace{1, \dots, 1}_i, \underbrace{-1, \dots, -1}_{k-i-1} \right] \quad (21)$$

1. Dempster et al
2. Isometric Log-Ratio Transform
3. Egozcue et al

برای به دست آوردن ζ برای هر π موجود، معکوس تبدیل ILR در رابطه ۲۰ به صورت زیر است:

$$\zeta = ILR^{-1}(\pi) = \left[\frac{\exp(\psi_1\pi)}{\sum_{i=1}^k \exp(\psi_i\pi)}, \dots, \frac{\exp(\psi_k\pi)}{\sum_{i=1}^k \exp(\psi_i\pi)} \right]. \quad (22)$$

حال از تبدیل ILR برای تغییر مسئله بهینه‌سازی مقید به یک مسئله بدون قید، استفاده می‌کنیم. لذا الگوریتم بهینه‌سازی پیشنهادی شامل گام‌های زیر می‌شود:

۱. مقادیر اولیه‌ای برای بردارهای پارامتری $\hat{\tau}^{(0)}$ و $\hat{\zeta}^{(0)}$ در نظر بگیرید.

۲. براساس تبدیل ILR مقدار $\hat{\pi}^{(0)}$ را از روی $\hat{\zeta}^{(0)}$ محاسبه نمایید.

$$\hat{\pi}^{(0)} = ILR(\hat{\zeta}^{(0)}). \quad (23)$$

۳. در تابع لگاریتم درست‌نمایی z را با مولفه‌های فردی زیر جایگذاری نمایید.

$$\frac{\exp(\psi_j\pi)}{\sum_{i=1}^k \exp(\psi_i\pi)}, \quad j = 1, \dots, k \quad (24)$$

۴. مقادیر $\hat{\tau}$ و $\hat{\pi}$ را به وسیله

$$(\hat{\tau}, \hat{\pi}) = \arg \max_{(\tau, \pi)} \log L(\tau, ILR^{-1}(\pi)). \quad (25)$$

۵. با استفاده از طرح تکراری نیوتن رافسون با مقادیر اولیه $\hat{\tau}^{(0)}$ و $\hat{\pi}^{(0)}$ محاسبه نمایید.

۶. تبدیل ILR معکوس را برای محاسبه $\hat{\zeta}$ از روی $\hat{\pi}$ استفاده نمایید.

$$\hat{\zeta} = ILR^{-1}(\hat{\pi}). \quad (26)$$

از آنجائی که جواب رابطه درست‌نمایی می‌تواند از انتخاب مقادیر اولیه، تأثیر بگیرد. در پیاده‌سازی، چندین مقدار اولیه در نظر گرفته شد و MLE ای انتخاب شد که بیش‌ترین مقدار را برای تابع لگاریتم درست‌نمایی نتیجه می‌داد.

۴-۲-۴. نتایج تجربی برای مجموعه مشاهدات گمشده

برای تحلیل خصوصیات برآوردگر $(\hat{\tau}, \hat{\zeta})$ ، ابتدا مجموعه داده کامل را به صورت تصادفی به دو بخش مطابق ساختار بخش ۲-۴ تقسیم می‌کنیم. بدین ترتیب دو مجموعه داده مصنوعی P_1 و P_2 ایجاد می‌گردد و سپس آن‌ها را به صورت مصنوعی ترکیب می‌نماییم که برای برآورد $\theta = (\tau, \zeta)$ استفاده می‌شوند. برآوردگر بر پایه P_1 و P_2 با $\hat{\theta}_1$ نمایش داده می‌شود. با توجه به تابع چگالی رابطه ۱۸، لگاریتم درستنمایی برای مجموعه P_1 به صورت زیر است:

$$\begin{aligned}
 l_i(\tau, \zeta) \propto & (\zeta_1 + \zeta_2) \exp \left\{ - \left(\frac{\exp(\beta t_i) - 1}{\beta} \right) \exp(\alpha + \beta x_i \right. \\
 & \left. + \gamma \mathbb{1}_{[b_i=High]} \right\} \\
 & \times \exp \{ \alpha + \beta(x_i + t_i) + \gamma \mathbb{1}_{[b_i=High]} \}^{d_i} \\
 & + (\zeta_3 + \zeta_4) \exp \left\{ - \left(\frac{\exp(\beta t_i) - 1}{\beta} \right) \exp(\alpha + \beta x_i + \gamma \mathbb{1}_{[b_i=High]} \right. \\
 & \left. + \delta_1) \right\} \\
 & \times \exp \{ \alpha + \beta(x_i + t_i) + \gamma \mathbb{1}_{[b_i=High]} + \delta_1 \}^{d_i} \\
 & + (\zeta_5 + \zeta_6) \exp \left\{ - \left(\frac{\exp(\beta t_i) - 1}{\beta} \right) \exp(\alpha + \beta x_i + \gamma \mathbb{1}_{[b_i=High]} \right. \\
 & \left. + \delta_2) \right\} \\
 & \times \exp \{ \alpha + \beta(x_i + t_i) + \gamma \mathbb{1}_{[b_i=High]} + \delta_2 \}^{d_i}.
 \end{aligned} \tag{۲۷}$$

لگاریتم درستنمایی برای مجموعه P_2 به صورت زیر است:

$$\begin{aligned}
 l_i(\tau, \zeta) \propto & (\zeta_1 + \zeta_3 + \zeta_5) \\
 & \times \exp \left\{ - \left(\frac{\exp(\beta t_i) - 1}{\beta} \right) \exp(\alpha + \beta x_i + \delta_1 \mathbb{1}_{[c_i=1]} + \delta_2 \mathbb{1}_{[c_i=2]}) \right\} \\
 & \times \exp \{ \alpha + \beta(x_i + t_i) + \delta_1 \mathbb{1}_{[c_i=1]} + \delta_2 \mathbb{1}_{[c_i=2]} \}^{d_i} \\
 & + \\
 & \times \exp \left\{ - \left(\frac{\exp(\beta t_i) - 1}{\beta} \right) \exp(\alpha + \beta x_i + \gamma + \delta_1 \mathbb{1}_{[c_i=1]} + \delta_2 \mathbb{1}_{[c_i=2]}) \right\}
 \end{aligned} \tag{۲۸}$$

$$\times \exp\{\alpha + \beta(x_i + t_i) + \gamma + \delta_1 \mathbb{1}_{[c_i=1]} + \delta_2 \mathbb{1}_{[c_i=2]}\}^{d_i}.$$

ابتدا تنها یکبار سناریو تقسیم‌بندی تصادفی به P_1 و P_2 را انجام می‌دهیم. برآورد پارامترها در این سناریو را با $(\hat{\theta}_1)$ نمایش می‌دهیم. سناریوی دیگر ۱۰۰۰ بار تقسیم کردن مجموعه داده کامل به صورت تصادفی است و میانگین برآورد پارامترها $(\bar{\theta} = \frac{1}{1000} \sum_k \hat{\theta}_k)$ و همچنین انحراف استاندارد پارامترها $(\sigma_{PAR} = (\sigma_{PAR} \times 10^3))$

نمایش می‌دهیم. کلیه نتایج در جدول ۴ ارائه شده است. $(\sqrt{\frac{1}{999} \sum_k (\hat{\theta}_k - \bar{\theta})^2})$ محاسبه می‌نماییم. برآورد پارامترها برای مجموعه داده کامل را با $(\hat{\theta})$

جدول ۴. برآورد و انحراف استاندارد پارامترها در داده کامل و گمشده با سناریو تقسیم‌بندی مختلف

برآورد پارامتر	$\hat{\alpha}$	$\hat{\beta}$	$\hat{\gamma}$	δ_1	δ_2	ζ_1	ζ_2	ζ_3	ζ_4	ζ_5
$\hat{\theta}$	-۴۱/۱۱	۱۱/۰	-۲۱/۰	-۲۲/۰	-۴۳/۰	۱۹/۰	۰۳/۰	۴۳/۰	۱۳/۰	۱۳/۰
$\hat{\theta}_1$	-۴۱/۱۱	۱۱/۰	-۰۸/۰	-۲۷/۰	-۳۷/۰	۲۱/۰	۰۰/۰	۵۴/۰	۰۲/۰	۰۰/۰
$\bar{\theta}$	-۴۲/۱۱	۱۱/۰	-۱۳/۰	-۲۳/۰	-۳۴/۰	۲۲/۰	۰۰/۰	۵۲/۰	۰۴/۰	۰۱/۰
$\sigma_{PAR} \times 10^3$	۴۳/۱۷	۲۱/۰	۲۴/۱۰۱	۲۲/۳۳	۹۴/۱۱۴	۰۹/۳	۸۷/۰	۷۴/۱۷	۴۹/۱۷	۰۸/۱۷

منبع: یافته‌های تحقیق

برای بررسی دلایل برآورد ضعیف به دست آمده از ترکیب P_1 و P_2 ، به بررسی شناسایی‌پذیری پارامترها می‌پردازیم. برای این منظور، براساس تابع مرگ و میر موجود تابع ϕ را برای مقادیر مختلف (B, C) و با توجه به سطح پایه و سن x به صورت زیر تعریف می‌کنیم:

$$\phi = \exp[\gamma \mathbb{1}_{[b=High]} + \delta_1 \mathbb{1}_{[c=c_1]} + \delta_2 \mathbb{1}_{[c=c_2]}].$$

جدول ۵ مقادیر $\hat{\phi}$ را برای مقادیر مختلف (B, C) برای مجموعه داده کامل، نشان می‌دهد.

جدول ۵. مقادیر برآورد ϕ ($\hat{\Phi}$) برای مجموعه داده کامل

۲	۱	۰	مشخصات جغرافیایی جمعیتی (C)
			سطح مزایا (B)
۰/۶۵	۰/۸۰	۱	پایین
۰/۵۳	۰/۶۵	۰/۸۱	بالا

منبع: یافته‌های تحقیق

در جدول ۵ مشاهده می‌شود که فردی با مشخصات جغرافیایی جمعیتی «صفر» و سطح مزایای بالا، با فردی که سطح مزایا پایین دارد و مشخصات جغرافیایی جمعیتی او «یک» است، مرگ‌ومیر بسیار مشابهی (در سن‌های یکسان) دارند. این موضوع برای افرادی که مزایای بالا و مشخصات جغرافیایی جمعیتی آن‌ها یک است، با افرادی که دارای مزایای پایین و مشخصات جغرافیایی جمعیتی «دو» هستند، نیز برقرار است. بنابراین شرایط شناسایی پذیری، برقرار نمی‌باشد.

همچنین شناسایی‌پذیری پارامترها در مدل مطرح‌شده با داده‌های موجود، می‌تواند از طریق بررسی مقادیر ویژه ماتریس هسین تحلیل شود. مقادیر ویژه برای تمام ۱۰۰۰ مجموعه داده شبیه‌سازی شده محاسبه شده است. مشاهده می‌شود که دو مقدار بزرگ‌تر مقادیر ویژه ماتریس هسین بسیار نزدیک به صفر هستند که نشان دهنده این است که ماتریس اطلاع برآورد شده، رتبه کامل نیست و مدل برای داده‌های موجود، دارای پارامتر افزونه است. علاوه بر این ممکن است یک مسئله در ریسک مدل موجود باشد. ما یک مدل گومپرتز برای داده‌ها استفاده کرده‌ایم که ممکن است نتواند تابع مرگ‌ومیر درستی را برای این تجربه مرگ‌ومیر منعکس کند و یا ممکن است که اختلافات مرگ‌ومیر به دلیل ویژگی‌های مختلف اقتصادی و اجتماعی به‌طور متناسب بر روی لگاریتم خطر تأثیر نگذارد.

۴-۳. دسترسی به اطلاعات بیشتر

مسئله شناسایی پذیری زمانی رخ می‌دهد که دو متغیر کمکی که هیچ‌گاه به صورت مشترک مشاهده نشده‌اند، مانع از استنباط بیشتر در مورد تابع مرگ‌ومیر می‌شود. اکنون برای بررسی تأثیر دسترسی به مجموعه داده بیشتر حالت‌های مختلفی را مدنظر قرار می‌دهیم. مجموعه داده P_3 اندازه نمونه کوچکتری به تعداد n_3 دارد که نماینده طرح بازنشستگی است که هیچ اطلاعاتی از دو متغیر کمکی، گمشده نیستند. برای نشان دادن استواری نتایج، مجموعه داده P_4 با اندازه نمونه بزرگتر به تعداد n_4 در نظر گرفته می‌شود که نماینده طرح بازنشستگی است که در آن B و C هر دو گمشده هستند.

حال چهار طرح بازنشستگی P_1, P_2, P_3 و P_4 داریم که از طریق تقسیم‌بندی مجموعه داده‌های کامل به صورت تصادفی و حذف عوامل مقتضی به دست آورده می‌شوند. در این رویکرد، ۵ درصد به طور کامل مشاهده شده‌اند (P_3)، ۲۰ درصد فقط B مشاهده شده (P_1)، ۲۰ درصد فقط C مشاهده شده (P_2) و در ۵۵ درصد باقی‌مانده هر دو B و C مشاهده نشده‌اند (P_4). از این رو $n_1 = n_2 = 3748$ ، $n_3 = 937$ و $n_4 = 10308$ است. این نحوه تقسیم‌بندی تصادفی ۱۰۰۰ بار تکرار می‌شود. بنابراین سهم تابع درست‌نمایی برای هر فرد در P_3 (بدون مشاهده گمشده) به صورت زیر است:

$$L_i(\tau, \zeta) \propto S_{x_i}(t_i | b_i, c_i; \tau) \mu_{x_i+t_i}^{d_i}(b_i, c_i; \tau) f_{B,C}(b_i, c_i; \zeta). \quad (29)$$

لذا:

$$L_i(\tau, \zeta) \propto (\zeta_1 + \zeta_2 + \zeta_3 + \zeta_4 + \zeta_5 + \zeta_6) \times \exp \left\{ - \left(\frac{\exp(\beta t_i) - 1}{\beta} \right) \exp(\alpha + \beta x_i + \gamma \mathbb{1}_{[b_i=High]} + \delta_1 \mathbb{1}_{[c_i=1]} + \delta_2 \mathbb{1}_{[c_i=2]}) \right\} \times \exp \{ \alpha + \beta(x_i + t_i) + \gamma \mathbb{1}_{[b_i=High]} + \delta_1 \mathbb{1}_{[c_i=1]} + \delta_2 \mathbb{1}_{[c_i=2]} \}^{d_i}. \quad (30)$$

و سهم تابع درست‌نمایی برای افراد در P_4 (هر دو B و C گمشده) به صورت زیر است:

$$L_i(\tau, \zeta) \propto \sum_{b \in \{Low, High\}} \sum_{c \in \{0,1,2\}} S_{x_i}(t_i | b, c; \tau) \mu_{x_i+t_i}^{d_i}(b, c; \tau) f_{B,C}(b, c; \zeta). \quad (31)$$

در نتیجه داریم:

$$\begin{aligned}
 L_i(\tau, \zeta) \propto & \zeta_1 \exp \left\{ - \left(\frac{\exp(\beta t_i) - 1}{\beta} \right) \exp(\alpha + \beta x_i) \right\} \exp \{ \alpha + \beta(x_i + t_i) \}^{d_i} \\
 & + \zeta_2 \exp \left\{ - \left(\frac{\exp(\beta t_i) - 1}{\beta} \right) \exp(\alpha + \beta x_i + \gamma) \right\} \exp \{ \alpha + \beta(x_i + t_i + \gamma) \}^{d_i} \\
 & + \zeta_3 \exp \left\{ - \left(\frac{\exp(\beta t_i) - 1}{\beta} \right) \exp(\alpha + \beta x_i + \delta_1) \right\} \exp \{ \alpha + \beta(x_i + t_i + \delta_1) \}^{d_i} \\
 & + \zeta_4 \exp \left\{ - \left(\frac{\exp(\beta t_i) - 1}{\beta} \right) \exp(\alpha + \beta x_i + \gamma + \delta_1) \right\} \exp \{ \alpha + \beta(x_i + t_i + \gamma + \delta_1) \}^{d_i} \\
 & + \zeta_5 \exp \left\{ - \left(\frac{\exp(\beta t_i) - 1}{\beta} \right) \exp(\alpha + \beta x_i + \delta_2) \right\} \exp \{ \alpha + \beta(x_i + t_i + \delta_2) \}^{d_i} \\
 & + \zeta_6 \exp \left\{ - \left(\frac{\exp(\beta t_i) - 1}{\beta} \right) \exp(\alpha + \beta x_i + \gamma + \delta_2) \right\} \exp \{ \alpha + \beta(x_i + t_i + \gamma + \delta_2) \}^{d_i}.
 \end{aligned} \tag{۳۲}$$

در ادامه، مقادیر برآورد پارامترها و خطاهای استاندارد به دست آمده از مجموعه داده‌های مختلف را با هم مقایسه می‌کنیم. برآوردها بر اساس ترکیب چهار مجموعه داده (با $P_1 - P_4$ مشخص شده)، بر اساس ترکیب P_1, P_2 و P_3 (با $P_1 - P_3$ مشخص شده) یعنی به غیر از مجموعه داده بزرگتر که هر دو B و C مشاهده نشده‌اند؛ و تنها برای P_3 (بدون داده گمشده ولی اندازه نمونه خیلی کوچک) را در نظر می‌گیریم.

برآورد پارامترها در مجموعه داده کامل $(\hat{\theta})$ ، برآورد پارامترها در یک بار تقسیم مجموعه داده کامل به صورت تصادفی $(\hat{\theta}_1)$ ، برآورد خطاهای استاندارد آنها $(\hat{\sigma}(\hat{\theta}_1) \times 10^3)$ ، میانگین برآورد پارامترها به وسیله ۱۰۰۰ بار تقسیم کردن مجموعه داده‌های کامل به صورت تصادفی $(\bar{\theta} = \frac{1}{1000} \sum_k \hat{\theta}_k)$ ، میانگین برآورد خطای استاندارد و انحراف استاندارد پارامترها به ترتیب $(\bar{\sigma}(\hat{\theta}) = \frac{1}{1000} \sum_k \hat{\sigma}(\hat{\theta}_k) \times 1000)$ و $\sigma_{PAR} =$

$$\sqrt{\frac{1}{999} \sum_k (\hat{\theta}_k - \bar{\theta})^2} \times 1000.$$

محاسبه و در جدول ۶ نمایش داده شده است.

این برآوردها را با نتایج حاصل از مجموعه داده کامل (خط اول جدول ۶) مقایسه می‌کنیم. خطاهای استاندارد با استفاده از منفی معکوس ماتریس اطلاع تجربی برآورد شده‌اند. پروسه برآورد در بخش ۵، شامل (τ, ζ) است. از این رو برآورد خطای استاندارد برای (τ, ζ) نیازمند این است که از روش دلنا برای توابع برداری استفاده شود. این روش تنها برای برآورد خطاهای استاندارد ζ استفاده می‌شود. تابع مورد استفاده، همان رابطه ۲۲ است. بنابراین مشتق تابع مد نظر نسبت به π ، یک برداری k بعدی به صورت زیر است:

$$\zeta'_j = \frac{\psi_j \exp(\psi_j \pi) \sum_{i=1}^k \exp(\psi_i \pi) - \exp(\psi_j \pi) \sum_{i=1}^k \psi_i \exp(\psi_i \pi)}{(\sum_{i=1}^k \exp(\psi_i \pi))^2} \quad (33)$$

که $j = 1, \dots, k$ است. بنابراین، واریانس آن به طور تقریبی از رابطه $var(\zeta) = \zeta' \Sigma \zeta$ حاصل می‌شود، که در آن Σ ماتریس واریانس کوواریانسی است که از طریق منفی معکوس ماتریس هسین به دست آمده در برآورد بردار پارامتری π حاصل می‌شود. در نهایت با جذر گرفتن از رابطه فوق، خطای استاندارد مربوط به بردار پارامتری ζ حاصل می‌گردد.

همان‌طور که از مقادیر $\hat{\theta}_k$ و $\hat{\theta}$ مشاهده می‌شود برآورد پارامترها با برآوردهای به دست آمده براساس مجموعه داده کامل، مطابقت دارند. اما، چون بخش‌هایی از مجموعه داده کامل از استنباط خارج می‌شوند، خطای استاندارد برآورده شده پارامترها $\hat{\sigma}(\hat{\theta})$ به دلیل کوچک‌تر بودن اندازه نمونه افزایش می‌یابد.

افزایش خطای استاندارد برآورد شده، برای پارامترهای γ ، δ_1 و δ_2 به این دلیل است که تفاوت‌های مرگ و میر بر پایه مزایا و مشخصات جغرافیایی جمعیتی، وابسته به متغیرهایی است که گاهی اوقات گمشده هستند. برای پارامترهای α و β این امر تأثیر کم‌تری دارد. به این دلیل که متغیر کمکی سن برای تمام افراد مشاهده شده است و خط مبنا برای همه افراد در چهار مجموعه داده طرح‌های بازنشستگی شبیه‌سازی شده، مشترک است. این تفاوت‌ها بیشتر هنگام مقایسه خطای استاندارد برآورد شده $\hat{\gamma}$ ، δ_1 و δ_2 وقتی که از هر چهار طرح بازنشستگی $(P_1 - P_4)$ با برآورد فقط بر اساس طرح‌های بازنشستگی $P_1 - P_3$ استفاده می‌شود، تأیید می‌شوند.

برای بررسی شناسایی‌پذیری مدل کامل با ۴ مجموعه داده $P_1 - P_4$ ، مشابه قبل می‌توان مقادیر ویژه ماتریس هسین حاصل از برآورد هر چهار مجموعه داده، را محاسبه نمود.

مشاهده می‌شود که در تمام ۱۰۰۰ طرح تصادفی داده گمشده، مقادیر ویژه همگی منفی و دور از صفر هستند. همچنین نزدیک‌ترین مقدار ویژه به صفر در تمام شبیه‌سازی‌ها برابر با $-9/93$ است. بنابراین مدل بررسی شده، شناسایی پذیر است.

جدول ۶. برآورد پارامترها و برآورد خطای استاندارد مربوطه در سناریوهای مختلف

مجموعه داده	برآورد	$\hat{\alpha}$	$\hat{\beta}$	$\hat{\gamma}$	$\hat{\delta}_1$	$\hat{\delta}_2$	$\hat{\zeta}_1$	$\hat{\zeta}_2$	$\hat{\zeta}_3$	$\hat{\zeta}_4$	$\hat{\zeta}_5$
کامل	$\hat{\theta}$	-۱۱/۴۱	۰/۱۱	-۰/۲۱	-۰/۲۲	-۰/۴۳	۰/۱۹	۰/۰۳	۰/۴۳	۰/۱۳	۰/۱۳
$P_1 - P_4$	$\hat{\theta}_1$	-۱۱/۴۶	۰/۱۱	-۰/۰۸	-۰/۱۶	-۰/۳۹	۰/۱۹	۰/۰۲	۰/۴۳	۰/۱۴	۰/۱۳
	$\hat{\sigma}(\hat{\theta}_1)$ $\times 10^3$	۱۳۶/۲۶	۱/۷۱	۶۷/۹۶	۶۳/۱۶	۸۴/۷۱	۶/۶۹	۴/۰۸	۹/۴۴	۷/۹۴	۷/۹۰
	$\hat{\theta}$	-۱۱/۴۵	۰/۱۱	-۰/۱۸	-۰/۱۹	-۰/۳۶	۰/۱۹	۰/۰۳	۰/۴۳	۰/۱۳	۰/۱۳
	$\hat{\sigma}(\hat{\theta})$ $\times 10^3$	۱۳۶/۴۲	۱/۷۲	۶۸/۴۱	۶۱/۶۵	۸۲/۱۶	۶/۹۴	۴/۶۸	۹/۴۴	۷/۹۴	۷/۶۶
	σ_{PAR} $\times 10^3$	۲۵/۴۹	۰/۲۹	۵۰/۹۷	۴۸/۵۸	۶۳/۲۹	۶/۱۴	۴/۴۸	۸/۵۶	۷/۳۰	۷/۴۷
$P_1 - P_3$	$\hat{\theta}_1$	-۱۱/۵۲	۰/۱۱	-۰/۱۰	-۰/۱۹	-۰/۴۷	۰/۱۹	۰/۰۲	۰/۴۳	۰/۱۴	۰/۱۳
	$\hat{\sigma}(\hat{\theta}_1)$ $\times 10^3$	۱۹۹/۲۱	۲/۵۲	۷۷/۰۹	۶۹/۰۱	۹۴/۲۹	۶/۶۹	۴/۰۹	۹/۴۴	۷/۹۵	۷/۹۱
	$\hat{\theta}$	-۱۱/۴۳	۰/۱۱	-۰/۲۰	-۰/۲۱	-۰/۴۱	۰/۱۹	۰/۰۳	۰/۴۳	۰/۱۳	۰/۱۳
	$\hat{\sigma}(\hat{\theta})$ $\times 10^3$	۲۰۰/۶۱	۲/۵۴	۷۴/۰۵	۶۶/۴۶	۸۸/۵۵	۶/۹۵	۴/۶۹	۹/۴۵	۷/۹۵	۷/۶۷
	σ_{PAR} $\times 10^3$	۱۵۴/۳۲	۱/۹۷	۵۹/۷۳	۵۴/۹۷	۷۱/۳۹	۶/۴۲	۴/۵۷	۸/۴۵	۷/۳۸	۷/۳۷
P_3	$\hat{\theta}_1$	-۱۱/۵۴	۰/۱۱	-۰/۱۶	-۰/۱۳	-۰/۶۶	۰/۱۹	۰/۰۲	۰/۴۵	۰/۱۴	۰/۱۱
	$\hat{\sigma}(\hat{\theta}_1)$ $\times 10^3$	۶۱۱/۵۲	۷/۷۴	۱۷۳/۳۷	۱۵۰/۹۲	۲۲۷/۵۲	۱۲/۹۰	۴/۳۶	۱۶/۲۵	۱۱/۴۴	۱۰/۳۵
	$\hat{\theta}$	-۱۱/۵۱	۰/۱۱	-۰/۲۱	-۰/۲۲	-۰/۴۴	۰/۱۹	۰/۰۳	۰/۴۳	۰/۱۳	۰/۱۳
	$\hat{\sigma}(\hat{\theta})$ $\times 10^3$	۶۰۴/۹۷	۷/۶۶	۱۶۶/۵۷	۱۵۲/۲۵	۱۹۹/۰۵	۱۲/۸۰	۵/۱۴	۱۶/۱۵	۱۱/۱۰	۱۰/۹۸
	σ_{PAR} $\times 10^3$	۶۳۴/۰۷	۸/۱۳	۱۵۲/۹۵	۱۴۶/۰۲	۱۹۰/۲۱	۱۲/۶۹	۵/۰۱	۱۵/۹۴	۱۰/۸۶	۱۰/۷۰

منبع: یافته‌های تحقیق

۴-۴. بررسی تأثیر مالی در استنباط آمار براساس مشاهدات گمشده

۴-۴-۱. بررسی نرخ مرگ و میر و عامل مستمری

فرض کنید متغیر تصادفی Y_x نشان دهنده ارزش فعلی جریان‌های نقدی در یک صندوق بازنشستگی برای عضو x ساله است که در طول عمر باقی‌مانده خود آن را دریافت می‌کند. اگر این جریان نقدی به میزان یک واحد هر ساله به‌طور پیوسته پرداخت شود، مقدار مورد انتظار برای Y_x ، عامل مستمری نامیده شده که با \bar{a}_x نمایش و به‌صورت زیر محاسبه می‌شود (دیکسون و همکاران^۱، ۲۰۱۳):

$$\bar{a}_x = \mathbb{E}(Y_x) = \int_0^{\infty-x} \exp(-rt) S_x(t) dt. \quad (34)$$

به دلیل گمشدگی داده‌ها، \hat{t} به $\hat{\zeta}$ وابسته است، توزیع نمونه‌گیری $\hat{\theta} = (\hat{t}, \hat{\zeta})$ به صورت مجانبی دارای توزیع نرمال چندمتغیره^۲ با میانگین θ و ماتریس واریانس کوواریانس Σ است، که از طریق منفی معکوس ماتریس اطلاع^۳ فشر^۳ برآورد می‌شود.

برای بررسی نرخ مرگ و میر و عامل مستمری، مقادیر عامل مستمری و درصد تغییر آن در مقایسه با عامل مستمری در مدل M_0 محاسبه و نتایج در جدول ۷ ارائه شده است. برآورد پارامترها با مشاهدات هر چهار طرح $P_1 - P_4$ در نظر گرفته می‌شوند. تجزیه و تحلیل عددی برای دو نرخ بهره ۱ و ۳ درصد و دو فرد به سن‌های ۶۵ و ۷۱ سال انجام شده است. اولین سن، همان سن معمول بازنشستگی و دومین سن به‌طور تقریبی، میانگین سن افراد زنده در مجموعه داده کامل در زمان اتمام دوره مشاهده است.

نتایج جدول نشان می‌دهد که برای سن‌های انتخابی و دو نرخ بهره محاسبه‌شده، بین عامل‌های مستمری به‌دست آمده از مدل مبتنی بر متغیرهای کمکی^۲ اقتصادی-اجتماعی و مقادیر به‌دست آمده از مدل M_0 تفاوت‌های قابل توجهی وجود دارد. این اختلافات برای نرخ بهره پایین‌تر، بیشتر است که بیانگر این است که وقتی نرخ بهره پایین است،

1. Dickson et al
2. Multivariate Normal Distribution
3. Fisher Information Matrix

مرگ و میر عامل خطر مهم‌تری برای مستمری است. برای بررسی تأثیر مالی استفاده از مدل مرگ و میر معرفی شده، می‌توان الزام سرمایه را بررسی نمود که نتیجه آن وابسته به عامل مستمری و در نتیجه مدل مرگ و میر است.

جدول ۷. مقادیر مستمری، درصد تغییرات بر اساس مدل M_0 و نرخ‌های بهره مختلف

$r = 3\%$		$r = 1\%$		تقسیم‌بندی	مجموعه داده	سن
درصد تغییر		درصد تغییر	\bar{a}_x			
-	۷۰/۱۳	-	۷۷/۱۶	M_0	$P_1 - P_4$	۶۵
-۸۱/۵	۹۰/۱۲	-۹۳/۶	۶۱/۱۵	$b = L, c = 0$		
-۲۶/۳	۲۵/۱۳	-۹۳/۳	۱۱/۱۶	$b = H, c = 0$		
-۹۸/۰	۵۶/۱۳	-۲۱/۱	۵۷/۱۶	$b = L, c = 1$		
۵۴/۱	۹۱/۱۳	۸۴/۱	۰۸/۱۷	$b = H, c = 1$		
۷۷/۵	۴۹/۱۴	۰۱/۷	۹۵/۱۷	$b = L, c = 2$		
۶۶/۶	۶۱/۱۴	۱۰/۱۰	۴۷/۱۸	$b = H, c = 2$		
-	۰۲/۱۱	-	۹۸/۱۲	M_0	$P_1 - P_4$	۷۱
-۲۸/۷	۲۲/۱۰	-۳۳/۸	۹۰/۱۱	$b = L, c = 0$		
-۱۵/۴	۵۷/۱۰	-۷۸/۴	۳۶/۱۲	$b = H, c = 0$		
-۳۲/۱	۸۸/۱۰	-۵۴/۱	۷۸/۱۲	$b = L, c = 1$		
۸۳/۱	۲۲/۱۱	۱۰/۲	۲۶/۱۳	$b = H, c = 1$		
۱۶/۷	۸۱/۱۱	۳۴/۸	۰۷/۱۴	$b = L, c = 2$		
۲۹/۸	۹۴/۱۱	۱۱/۱۲	۵۶/۱۴	$b = H, c = 2$		

منبع: یافته‌های تحقیق

۴-۴-۱. ریسک برآورد نادرست الزام سرمایه

هدف این بخش نحوه کاهش الزام سرمایه براساس طرح بازنشستگی موجود، و ترکیب طرح‌های بازنشستگی با دیگر با هدف تخمین پارامترهای نرخ مرگ و میر است.

با استفاده از رویکرد ریچارد^۱ (۲۰۱۶) برای محاسبه SCR^A ، بر اساس مجموعه بزرگی از مقادیر پارامتر جایگزین، ارزیابی‌های مکرری از کل پرتفوی مستمری انجام می‌دهیم. مقدار پرتفوی به‌عنوان تابعی از پارامتر θ با $P(\theta)$ نشان داده می‌شود. الزام سرمایه ناشی از ریسک پارامترها با برآورد نادرست، به شرح زیر محاسبه می‌شود:

$$\left(\frac{\text{چندک } 99/5 \text{ تابع } P(\theta)}{P(\theta) \text{ میانگین}} - 1 \right) \times 100. \quad (35)$$

برای یافتن برآورد صدک و میانگین $P(\theta)$ ، مقدار شبیه‌سازی شده برای θ محاسبه می‌شود. مقادیر شبیه‌سازی شده را با θ' نمایش داده و k امین مقدار شبیه‌سازی شده برای پرتفوی ($P(\theta')$) به‌صورت زیر محاسبه می‌شود:

$$P(\theta'_k) = \sum_{i=1}^n w_i \bar{a}_{x_i}(\theta'_k). \quad (36)$$

که در این رابطه، w_i نشان‌دهنده مبلغ مستمری شخص نام است. بنابراین ریسک برآورد نادرست الزام سرمایه را به‌عنوان مثال برای طرح بازنشستگی P_1 به سه روش محاسبه می‌کنیم:

- ✓ پارامتر تابع مرگ و میر τ فقط با استفاده از داده‌های موجود در P_1 برآورد می‌شود. به‌منظور ساده نگه داشتن کارها و جلوگیری از مسئله شناسایی پذیری به‌دلیل نبود C ، تابع مرگ و میر مانند مدل M_1 است که پارامترهای آن را می‌توان با حداکثر درست‌نمایی تخمین زد.
- ✓ تمام پارامترهای تابع مرگ و میر با استفاده از ترکیب طرح‌های $P_1 - P_3$ ، از طریق رویکرد داده‌های گمشده، برآورد می‌شوند.
- ✓ تمام پارامترهای تابع مرگ و میر با استفاده از ترکیب طرح‌های $P_1 - P_4$ برآورد می‌شوند.

1. Richards
2. The Solvency Capital Requirement

این کار را برای برای طرح بازنشستگی P_3 تکرار می‌کنیم که در این طرح هر دو متغیر کمکی موجود است. الزامات سرمایه برای ریسک برآورد نادرست براساس این دو رویکرد و دو نرخ بهره، در جدول ۸ نشان داده شده است.

جدول ۸ الزام سرمایه برای طرح بازنشستگی P_1 و P_3 بر اساس دو نرخ بهره فرضی و طرح‌های مورد استفاده برای برآورد نرخ مرگ و میر

نمونه‌ها						نرخ بهره و الزام سرمایه
P_1		$P_1 - P_3$		$P_1 - P_4$		طرح بازنشستگی P_1
%۳	%۱	%۳	%۱	%۳	%۱	نرخ بهره
%۲۷/۲	%۸۲/۲	%۶۲/۱	%۹۵/۱	%۰۹/۱	%۳۰/۱	الزام سرمایه
P_3		$P_1 - P_3$		$P_1 - P_4$		طرح بازنشستگی P_3
%۳	%۱	%۳	%۱	%۳	%۱	نرخ بهره
%۶۱/۴	%۷۹/۵	%۵۴/۱	%۸۸/۱	%۰۳/۱	%۲۶/۱	الزام سرمایه

منبع: یافته‌های تحقیق

از نتایج این جدول مشاهده می‌شود که در نظر گرفتن سایر مجموعه داده‌ها برای برآورد نرخ مرگ و میر، می‌تواند الزام سرمایه را کاهش دهد. در نظر گرفتن مشاهدات برای وقتی که در آن‌ها بعضی از متغیرهای کمکی یا تمامی آن‌ها گمشده هستند، اندازه نمونه را افزایش داده و سبب کاهش عدم حتمیت پارامتر می‌شود. این نتایج نشان می‌دهند که در نظر گرفتن نمونه‌هایی با مشاهدات گمشده، می‌تواند تغییرپذیری برآورد نمونه‌گیری را کاهش دهد و به نوبه خود سبب کاهش الزام سرمایه شود. و یکی از راه‌های به دست آوردن چنین مشاهدات اضافی، ترکیب داده‌های تجربی از طرح‌های مختلف بازنشستگی است با این فرض که اعضای آن طرح‌ها از قانون مرگ و میر یکسانی پیروی کنند.

۵. جمع‌بندی و پیشنهادها

در این پژوهش به ترکیب داده‌های حاصل از طرح‌های بازنشتگی مختلف، با فرض اینکه قانون احتمال یکسانی برای طول عمر آتی دارند، اما اطلاعات متغیر کمکی متفاوت است، پرداخته شد. اگر متغیرهای کمکی حتی در یک تجربه مرگ و میر گمشده باشد، می‌توان از همین روش‌ها استفاده کرد.

همچنین به دست آوردن داده‌های کامل برای زیرمجموعه کوچکی از اعضا به ما این امکان را می‌دهد داده‌های حاصل از دو یا چند تجربه را با هم تلفیق و از موانع شناسایی‌پذیری جلوگیری کنیم. این نتیجه از لحاظ عملی اهمیت دارد، زیرا به دست آوردن اطلاعات کامل برای همه اعضای صندوق‌های بازنشتگی مختلف ممکن است دشوار یا گران باشد.

روش‌های مطرح شده می‌تواند برای محاسبه کمیت‌های مالی مورد علاقه براساس عامل‌های مستمری، برای بیم‌سنج‌ها مفید باشد. این روش مجموعه داده‌های مختلف با تجربه مرگ و میر برابر یا مشابه را با هم ترکیب، اندازه نمونه را افزایش و ریسک پارامتر را کاهش می‌دهد، بنابراین منجر به کاهش الزام سرمایه می‌شوند. متغیرهای اقتصادی-اجتماعی از جمله سطح مزایا و مشخصات جغرافیایی جمعیتی در صورت پایین بودن نرخ بهره بیشتر مورد توجه قرار می‌گیرند.

بررسی این‌که آیا می‌توان بدون کنار گذاشتن داده‌های مربوطه در وضعیت پارامتر افزونه، استنباط آماری کرد، مساله مورد بررسی قرار نگرفته است و این مسئله می‌تواند در چارچوب استنباط بیزی بهتر حل شود، زیرا هنگامی که ماتریس اطلاع فیشر معین مثبت اکید نیست، برآوردگر ماکسیمم درست‌نمایی، توزیع نرمال مجانبی ندارد (وانتانابه^۱، ۲۰۱۰). همچنین یکی از مفروضات مدل‌بندی این بود که اختلافات مرگ و میر بین گروه‌های اقتصادی-اجتماعی با افزایش سن تغییر نمی‌کند. لذا تأثیر و تجزیه و تحلیل یک عبارت اثر متقابل می‌تواند در بررسی‌های آتی مدنظر قرار گیرد.

ملاحظات اخلاقی

حامی مالی

این مقاله حامی مالی ندارد.

مشارکت نویسندگان

تمام نویسندگان در آماده سازی این مقاله مشارکت کرده‌اند.

تعارض منافع

بنا به اظهار نویسندگان، در این مقاله هیچ گونه تعارض منافی وجود ندارد.

تعهد کپی رایت

طبق تعهد نویسندگان، حق کپی راست (CC) رعایت شده است.



منابع

ذکایی، محمد و مقصودی، مسطوره. (۱۳۸۹). بازسازی مدل‌های مرگ‌ومیر بر پایه شکنندگی با استفاده از تعمیم توزیع گومپرتز. *فصلنامه صنعت بیمه*، ۲۵(۴): ۵۹-۸۵.

شجاعی‌آذر، زهرا و حسن‌زاده، امین. (۱۳۹۳). کاربرد مدل‌های فاز-نوع در مدل‌بندی مرگ‌ومیر. *پژوهشنامه بیمه*، ۲۹(۱): ۱۰۵-۱۲۶.

کمیجانی، اکبر، کوششی، مجید و نیاکان، لیلی. (۱۳۹۲). برآورد و پیش‌بینی نرخ مرگ‌ومیر در ایران با استفاده از مدل لی-کارتر. *پژوهشنامه بیمه*، ۲۸(۴): ۱-۲۵.

مهدوی، غدیر، دقیقی اصل، علیرضا و لطفی، نیر. (۱۳۹۰). کاربرد یک مدل مرگ‌ومیر با چند عامل ریسک در فسخ قراردادهای بیمه عمر (مورد مطالعه: یک شرکت بیمه). *پژوهشنامه بیمه*، ۲۶(۳): ۱-۲۸.

Catchpole, E. A. & Morgan, B. J. T. (1997). Detecting parameter redundancy. *Biometrika*, 84(1): 187-196.

Chen, Q., May, R. C., Ibrahim, J. G., Chu, H. & Cole, S. R. (2014). Joint modeling of longitudinal and survival data with missing and left-censored time-varying covariates. *Statistics in Medicine*, 33(26): 4560-4576.

Dempster, A. P., Laird, N. M. & Rubin D. B. (1977). Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society*, 39(1): 1-38.

Dickson, D., Hardy, M. & Waters, H. (2013). Actuarial mathematics for life contingent risks. international series on actuarial science. Cambridge University Press.

Egozcue, J. J., Pawlowsky-Glahn, V., Mateu-Figueras, G. & Barceló-Vidal, C. (2003). Isometric logratio transformations for compositional data analysis. *Mathematical Geology*, 35(3): 279-300.

Gompertz, B. (1825). On the nature of the function expressive of the law of human mortality, and on a new mode of determining the

- value of life contingencies. *Philosophical Transactions of the Royal Society of London*, 115: 513–583.
- Herring, A. H. & Ibrahim, J. G. (2001). Likelihood-Based methods for missing covariates in the cox proportional hazards model. *Journal of the American Statistical Association*, 96(453): 292–302.
- Lin, X. S. & Liu, X. (2007). Markov aging process and Phase-Type law of mortality. *North American Actuarial Journal*. 11(4): 92–109.
- Little, R. & An, H. (2004). Robust Likelihood-Based analysis of multivariate data with missing values. *Statistica Sinica*, 14(3): 949–968.
- Lord, F. M. (1955). Estimation of parameters from incomplete data. *Journal of the American Statistical Association*, 50(271): 870–876.
- Madrigal, A. M., Matthews, F. E., Patel, D., Gaches, A. & Baxter, S. (2011). What longevity predictors should be allowed for when valuing pension scheme liabilities? *British Actuarial Journal*, 16(1): 1–38.
- Macdonald, A. S., Richards, S. J. & Currie, I. D. (2018). Modelling mortality with actuarial applications. International Series on Actuarial Science. Cambridge University Press.
- McLachlan, G. & Peel, D. (2000). Finite mixture models. Wiley Series in Probability and Statistics, New York.
- Richards, S. J. (2016). Mis-Estimation risk: Measurement and impact. *British Actuarial Journal*, 21(3): 429–457.
- Rubin, D. B. (1976). Inference and missing data. *Biometrika*, 63(3): 581–592.
- Schluchter, M. D. & Jackson, K. L. (1989). Log-Linear analysis of censored survival data with partially observed covariates. *Journal of the American Statistical Association*, 84(405): 42–52.
- Titterington, D. M., Smith, A. F. M. & Makov, U. E. (1985). Statistical analysis of finite mixture distributions. New York, Wiley.
- Tsiatis, A. (2007). Semiparametric theory and missing data. Springer Science & Business Media.

- Ungolo, F., Christiansen, M. C., Kleinow, T. & MacDonald, A. S. (2019). Survival analysis of pension scheme mortality when data are missing. *Scandinavian Actuarial Journal*, 2019 (6): 523–547.
- Watanabe, S. (2010). Asymptotic equivalence of bayes cross validation and widely applicable information criterion in singular learning theory. *Journal of Machine Learning Research*, 11(116): 3571–3594.
- Wilks, S. S. (1932). Moments and distributions of estimates of population parameters from fragmentary samples. *The Annals of Mathematical Statistics*, 3(3): 163–195.
- Xu, Y., Kim, J. K. & Li, Y. (2017). Semiparametric estimation for measurement error models with validation data. *Canadian Journal of Statistics*, 45(2): 185–201.
- Yashin, A. (2001). Mortality models incorporating theoretical concepts of ageing. *In Forecasting Mortality in Developed Countries*, 261–280.

