



## Analysis of Censored Data Using Dirichlet Process Mixture Model with Generalized Inverse Weibull Distribution as Kernel

Haji Joudaki, B.<sup>1</sup> , Hashemi, R.<sup>2</sup> , Khazaei, S.<sup>2</sup> 

<sup>1</sup>Department of Statistics, Lorestan University, Khorramabad, Iran.

<sup>2</sup>Department of Statistics, Razi University, Kermansha, Iran.

**Corresponding author:** B. Haji Joudaki, Bahram600631@gmail.com

Received: 23/4/2023 Revised: 23/7/2023 Accepted and Published Online: 31/7/2023.

### Introduction

Evaluating complex data distribution, such as multimodal cases, requires more complex statistical models. The complexity of a statistical model is related to the number of model parameters. To achieve a more complex model and consequently achieve better flexibility in statistical inference, one can use an infinite-dimensional family of probability models. Using the Bayesian approach for the infinite-dimensional parameter need a suitable prior distribution. Typically, the prior distribution of such parameters are stochastic processes. Such priors are called nonparametric Bayes priors. The most important nonparametric Bayes prior is the Dirichlet process, first introduced by Ferguson (1973). Due to the discreteness of the Dirichlet process, using a mixed model of Dirichlet processes (DPMM) (Antoniak, 1974) is preferable. The ability to cover multimodal data distribution and perform data clustering are two advantages of DPMMs. In this article, we profit from these two advantages of DPMMs. Kernel selection is an essential issue in working with DPMMs. According to under study data, a flexible distribution should be selected as the kernel of DPMM. For lifetime data, distributions such as Weibull, lognormal, or other lifetime distributions can be chosen. The generalized inverse Weibel distribution is a flexible distribution introduced by de Gusmão, et al. (2011). In this article, the generalized inverse Weibel distribution is considered the kernel of DPMM.

### Material and Methods

In fitting a DPMM, prolongation of the execution time of Markov chain

Monte Carlo simulation algorithms is challenging. Dirichletprocess package in R software is a resolution. The mentioned package has a high ability to fit DPMMs and other non-parametric Bayes models. After obtaining posterior samples, survival density and hazard rate functions can be easily estimated using this package. Finally, by analyzing several real and simulated data sets, the performance of the proposed model is evaluated.

### Results and Discussion

We designed a simulation study to evaluate the performance of the proposed model under different prior distributions for the accuracy parameter of the Dirichlet process based on sample sizes of 100 and 1000. Bayesian and interval estimations of the survival, density, and hazard rate functions of a mixed model of two generalized inverse Weibull distributions reported. The results show that the proposed model has a high potential in estimating the mentioned functions. The proposed model was also used to analyze several real multimodal data sets. The results show that the proposed model performed better compared to other methods. The proposed model is also applied for real data clustering and simulation.

### Conclusion

A mixed model of Dirichlet processes with a generalized inverse Weibull kernel is used to analyze right-censored survival data. The performance of the proposed model is evaluated based on Several simulated and real data sets. Achieved results in this paper show that the proposed model has good potential for estimating density, survival, and hazard rate functions in survival data. Another advantage of the proposed model is its high potential for data clustering.

**Keywords:** Mixture model, Dirichlet process, Generalized inverse Weibull Distribution, Right censoring, Markov chain monte carlo.

**Mathematics Subject Classification (2010):** 62F10, 62N01.



©The Author(s). The Publisher is Iranian Statistical Society.

This is an open access article distributed under the terms and conditions of [\(CC BY-NC 4.0\)](https://creativecommons.org/licenses/by-nc/4.0/)



مجله علوم آماری، پاییز و زمستان ۱۴۰۲

جلد ۱۷، شماره ۲، ص ۲۷۵ -- ۲۹۸

DOI: 10.52547/jss.17.2.03

مقاله پژوهشی

## تحلیل داده‌های سانسور شده با استفاده از مدل آمیخته فرایند دیریکله با هسته وارون وایبل تعمیم یافته

بهرام حاجی جودکی<sup>۱</sup>، سیدرضا هاشمی<sup>۲</sup>، سلیمان خزایی<sup>۲</sup>

گروه آمار، دانشکده علوم پایه، دانشگاه لرستان

گروه آمار، دانشکده علوم، دانشگاه رازی

نویسنده مسئول: بهرام حاجی جودکی، bahram600631@gmail.com

تاریخ دریافت: ۱۴۰۱/۰۰/۰۰ تاریخ بازنگری: ۱۴۰۲/۰۰/۰۰ تاریخ پذیرش و انتشار: ۱۴۰۲/۰۰/۰۰

**چکیده:** در این مقاله یک مدل آمیخته فرایند دیریکله جدید با هسته وارون وایبل تعمیم یافته پیشنهاد شده است. پس از تعیین توزیع پیشین پارامترها در مدل پیشنهادی، برای نمونه‌گیری از توزیع پسین توام پارامترها از روش‌های مونت کارلوی زنجیره مارکف استفاده شده است. عملکرد مدل پیشنهادی با تحلیل چندین مجموعه داده واقعی و شبیه‌سازی شده مورد بررسی قرار گرفته است. در مجموعه داده‌های واقعی برخی از داده‌ها سانسور شده از راست هستند. همچنین در این مقاله پتانسیل مدل پیشنهادی برای خوشه‌بندی کردن داده‌ها بکار گرفته شده است. نتایج بدست آمده نشان دهنده عملکرد مطلوب مدل پیشنهادی است. **واژه‌های کلیدی:** مدل آمیخته، فرایند دیریکله، وارون وایبل تعمیم یافته، سانسور از راست، مونت کارلوی زنجیره مارکف.

کد موضوع بندی ریاضی (۲۰۱۰): 62N01, 62F10

## ۱ مقدمه

یکی از مهم‌ترین دغدغه‌های یک آماردان، انتخاب مدلی برای انجام استنباط آماری است. مسئله انتخاب مدل از دیدگاه نظری و عملی یک مسئله مهم اما دشوار است. اکثر آماردانان بر این باور هستند که عمده مشکلاتی که در استنباط آماری هست را می‌توان مشکلات مربوط به مدل‌سازی آماری در نظر گرفت. در واقع نحوه ترجمه از



©نویسندگان). ناشر انجمن آمار ایران است.

این مقاله با دسترسی آزاد تحت شرایط و ضوابط (CC BY-NC 4.0) توزیع شده است.

مسئله موضوعی به مدل آماری اغلب حیاتی‌ترین بخش یک تحلیل آماری است. در علم آمار برای اینکه بتوان به راحتی به مجموعه‌ای از استنباط‌های آماری دست یافت، معمولاً یک خانواده از مدل‌های احتمال انتخاب می‌شود که این خانواده از مدل‌ها با مجموعه‌ای از پارامترهای که بُعد متناهی دارند اندیس‌گذاری می‌شوند. حال اگر برای پارامترهای مدل، توزیع‌های پیشین‌های انتخاب شود آنگاه می‌توان توزیع پسین پارامترها را محاسبه کرد. در ادامه هر نوع استنباطی در مورد پارامترها بر مبنای توزیع پسین پارامترها خواهد بود که از این رویکرد به عنوان رویکرد بیز پارامتری یاد می‌شود. مدل‌های احتمال ساده، معمولاً توانایی پوشش دادن توزیع داده‌های پیچیده همانند داده‌های چند مُدی را ندارند، بنابراین اگرچه کار کردن با این مدل‌ها راحت‌تر است ولی ممکن است نتایج بدست آمده بر اساس این مدل‌ها با نتایج واقعی اختلاف اساسی داشته باشند. برای پوشش دادن توزیع داده‌های پیچیده، به ناچار باید از مدل‌های آماری پیچیده‌تر استفاده کرد. در کل پیچیدگی یک مدل آماری با تعداد پارامترهای مدل مرتبط است. برای دستیابی به یک مدل پیچیده‌تر و در نتیجه دستیابی به نتایج بهتر، معقول است راهکاری برای انجام استنباط آماری برگزیده شود که در چارچوب این راهکار مجاز به داشتن خانواده‌ای بینهایت بُعدی از مدل‌های احتمال باشیم. در این صورت یک کلاس غنی‌تر و بزرگ‌تر از مدل‌های احتمال در دسترس خواهد بود. برای اینکه بتوان رویکرد بیزی را برای چنین مدل احتمالی پیاده‌سازی کرد باید برای پارامتر بینهایت بُعدی مدل مذکور توزیع پیشینی انتخاب شود. با توجه به نامتناهی بودن بُعد چنین پارامتری نمی‌توان برای آن هر توزیع پیشینی را در نظر گرفت و نوعاً توزیع پیشین چنین پارامترهای از نوع فرایندهای تصادفی هستند که چنین پیشین‌های را اصطلاحاً پیشین‌های بیز ناپارامتری می‌نامند. چنانچه برای انجام استنباط آماری از چنین رویکردی استفاده شود در این صورت گوئیم از رویکرد بیزی ناپارامتری استفاده شده است. مهم‌ترین پیشین بیز ناپارامتری فرایند دیریکله<sup>۱</sup> (DP) است که اولین بار توسط فرگوسن (۱۹۷۳) معرفی و برخی از ویژگی‌های مهم آن بیان و اثبات شد. سوسارلا و ون رابزین (۱۹۷۶)، فرگوسن و فادیا (۱۹۷۹)، لاهیری و پارک (۱۹۹۱) و کائو و همکاران (۱۹۹۲) برای مدل‌بندی تابع توزیع بقا از توزیع پیشین DP استفاده کردند. علیرغم جذابیت‌های توزیع پیشین DP، اما این فرایند یک فرایند گسسته است بنابراین تقریب زدن توابع پیوسته‌ای مانند چگالی، بقا و نرخ‌خطر با آن خالی از اشکال نخواهد بود. با پیش‌توزیع پیشین DP با برخی توزیع‌های پیوسته می‌توان مشکل مذکور را رفع کرد که این منجر به معرفی مدل‌های آمیخته‌ای از فرایندهای دیریکله<sup>۲</sup> (DPMM) می‌شود (آنتونیاک، ۱۹۷۴). توانایی در پوشش دادن توزیع داده‌های چندمُدی و انجام خوشه‌بندی داده‌ها از مهم‌ترین مزایای DPMM‌ها است. در این مقاله این دو مزیت DPMM‌ها مورد استفاده قرار می‌گیرد. برای دستیابی به اطلاعات بیشتر در مورد پیشین‌های بیز ناپارامتری به فرگوسن (۱۹۷۳)، فرگوسن (۱۹۷۴) و قوش و رامامورتی (۲۰۰۳) رجوع شود. انتخاب هسته موضوع مهمی در کار کردن با DPMM‌ها است. در واقع با توجه به نوع داده‌های تحت مطالعه توزیعی انعطاف‌پذیر به عنوان هسته DPMM انتخاب می‌شود. برای داده‌های حقیقی مقدار، معمولاً هسته توزیع نرمال یا نرمال‌چوله انتخاب می‌شود. برای داده‌های طول‌عمر، توزیع‌هایی مانند وایبل، گاما، لگ‌نرمال، بور ۱۲ و

<sup>1</sup>Dirichlet Process<sup>2</sup>Dirichlet Process Mixture Model

یا دیگر توزیع‌های طول عمر انتخاب‌های مناسب‌تری برای هسته هستند. برای داده‌های با مقادیر در بازه (۱, ۰)، معمولاً هسته توزیع بتا است. هسن (۲۰۰۶)، کوتاس (a۲۰۰۶)، کوتاس (b۲۰۰۶)، چنگ و یوان (۲۰۱۳)، بهلوری حجار و خزائی (۲۰۱۸) و حاجی جودکی و همکاران (۲۰۲۲) DPMM با هسته‌های به ترتیب گاما، وایبل، بتا، لگ‌نرمال، بور۱۲، دو پارامتری و بور۱۲ سه پارامتری را برای تحلیل داده‌های بقای سانسور شده از راست بکار بردند. اگرچه DPMM با هسته‌های مذکور مدل‌های انعطاف‌پذیری را ارائه می‌دهند اما هیچ کدام از این مدل‌ها لزوماً همیشه نتایج خوبی را فراهم نمی‌کنند. با انتخاب دیگر توزیع‌های انعطاف‌پذیر به عنوان هسته DPMM می‌توان مدل‌های آمیخته جدیدی معرفی کرد که این مدل‌ها نیز ممکن است در برخی مجموعه داده‌ها در مقایسه با دیگر DPMM ها نتایج بهتری ارائه دهند. یکی از توزیع‌های انعطاف‌پذیری که اخیراً در تحلیل داده‌های بقا و قابلیت اعتماد فراوان بکار می‌رود توزیع وارون وایبل تعمیم‌یافته<sup>۱</sup> (GIW) است که توسط دی گاس مائو و همکاران (۲۰۱۱) معرفی شد. این توزیع همانند توزیع لگ‌نرمال دارای تابع نرخ خطر تک مُدی است. برای دستیابی به اطلاعات بیشتر در مورد توزیع GIW به دی گاس مائو و همکاران (۲۰۱۱) و کائو و همکاران (۲۰۱۸) رجوع شود. در این مقاله هدف این است که یک مدل DPM جدید با هسته GIW معرفی شود و برخی از پتانسیل‌های این مدل جدید خصوصاً در انجام خوشه‌بندی و مدل‌بندی داده‌های بقای چند مُدی نشان داده شود.

سانسور پدیده‌ای رایج در تحلیل داده‌های بقا است که خود دارای انواع مختلفی مانند سانسور از راست، سانسور از چپ، سانسور فاصله‌ای و غیره است. سانسور متداول در داده‌های بقا، سانسور از راست است که خود دارای دو نوع سانسور زمان (نوع ۱) و سانسور شکست (نوع ۲) است. سانسور مورد علاقه در این مقاله سانسور زمان است. یک منبع مفید برای مطالعه انواع سانسور لاولس (۲۰۱۱) است.

در بخش ۲ ابتدا برخی تعاریف را ارائه داده و آنگاه در مورد توزیع GIW صحبت می‌شود. در بخش ۳، DPMM با هسته GIW معرفی می‌شود و فرایند اجرای نمونه‌گیری از توزیع پسین توام پارامترها شرح داده می‌شود. در بخش‌های ۴ و ۵ مدل پیشنهادی برای تحلیل مجموعه داده‌های شیشه‌سازی شده و واقعی بکار برده می‌شود و در نهایت در بخش پایانی نتایج بدست آمده ارائه شده است.

## ۲ تعاریف

**تعریف ۱.** فرض کنید سه‌گانه  $(A, B, P)$  یک فضای احتمال،  $P$  یک اندازه احتمال روی فضای  $A$  و  $\nu$  نیز یک مقدار عددی مثبت است. اگر برای هر افزار اندازه‌پذیر متناهی دلخواه  $B_1, \dots, B_k$  از  $A$  رابطه  $(P(B_1), \dots, P(B_k)) \sim Dir(\nu P_*(B_1), \dots, \nu P_*(B_k))$  برقرار باشد، آنگاه  $P$  یک فرایند دیریکله با پارامترهای  $\nu$  و  $P_*$  است و با نماد  $P \sim DP(\nu, P_*)$  نشان داده می‌شود که در آن نماد  $Dir$  برای نشان دادن توزیع دیریکله بکار رفته است.

<sup>1</sup>Generalized Inverse Weibull

معمولاً برای مدل‌بندی توزیع داده‌های (بقای) چند مودی از مدل‌های آمیخته استفاده می‌شود. یکی از مدل‌های آمیخته‌ای که برای این منظور می‌تواند مورد استفاده قرار گیرد مدل DPM است. این مدل بوسیله آنتونیاک (۱۹۷۴) معرفی شد. در حالت کلی یک مدل DPM به صورت

$$f(t : P) = \int_{\Theta} k(t|\theta) dP(\theta)$$

تعریف می‌شود، که در آن  $k$  هسته مدل آمیخته،  $\theta$  بردار پارامتری و  $\Theta$  فضای پارامتری،  $P$  تابع توزیع پیشین بردار پارامتری  $\theta$  است که نقش وزن‌های مدل آمیخته را ایفا می‌کند. در این چارچوب این اندازه احتمال تصادفی است و فرض می‌شود  $P \sim DP(\nu, P_0)$ . در بخش قبلی اشاره شد که در این مقاله هدف این است که توزیع GIW به عنوان هسته DPMM بکار برده شود. توابع چگالی، بقا و نرخ خطر<sup>۱</sup> توزیع GIW به ترتیب به صورت

$$k(t|\alpha, \beta, \lambda) = \beta \lambda \alpha^\beta x^{-\beta-1} \exp(-\lambda(\frac{\alpha}{t})^\beta), \quad t > 0, \alpha > 0, \beta > 0, \lambda > 0$$

$$S(t|\alpha, \beta, \lambda) = 1 - \exp(-\lambda(\frac{\alpha}{t})^\beta), \quad t > 0$$

$$h(t|\alpha, \beta, \lambda) = \frac{\beta \lambda \alpha^\beta x^{-\beta-1} \exp(-\lambda(\frac{\alpha}{t})^\beta)}{1 - \exp(-\lambda(\frac{\alpha}{t})^\beta)}, \quad t > 0$$

هستند، که در آن  $\alpha$  پارامتر مقیاس و  $\beta$  و  $\lambda$  نیز پارامترهای شکل توزیع هستند. توزیع GIW با پارامترهای مذکور با نماد  $GIW(\alpha, \beta, \lambda)$  نشان داده می‌شود. اگر  $\lambda = 1$  باشد آنگاه توزیع GIW به توزیع وارون وایبل تبدیل می‌شود. در شکل ۱ نمودار توابع چگالی و نرخ خطر توزیع GIW به ازای مقادیر مختلف پارامترهای  $\alpha$ ،  $\beta$  و  $\lambda$  رسم شده است. یکی از ویژگی‌های مهم توزیع‌های طول عمر که معمولاً در تحلیل داده‌های بقا مورد توجه است شکل تابع خطر است. تابع خطر توزیع GIW تک مودی است که این وضعیت در نمودار ب شکل ۱ به وضوح مشاهده می‌شود.

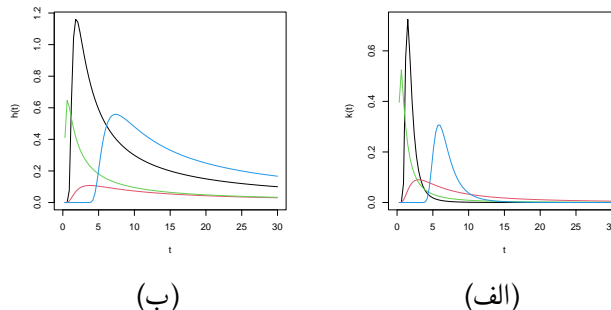
### ۳ معرفی مدل

مدل آمیخته فرایند دریکله (DPMM) با هسته GIW به صورت

$$f(t : P) = \int_{\alpha} \int_{\beta} \int_{\lambda} k(t|\alpha, \beta, \lambda) dP(\alpha, \beta, \lambda)$$

تعریف می‌شود، که در آن  $k$  تابع چگالی توزیع GIW و  $P$  نیز دارای توزیع پیشین  $DP(\nu, P_0)$  است. در این مدل پارامتر  $\nu$  را کمیتی تصادفی در نظر گرفته و برای آن مشابه اسکوبار و وست (۱۹۹۵) توزیع پیشین گامای  $G(a_\nu, b_\nu)$

<sup>۱</sup>Hazard Rate



شکل ۱. الف- نمودار تابع چگالی، ب- نمودار تابع نرخ خطر توزیع GIW به ازای مقادیر مختلف پارامترهای  $\alpha$ ،  $\beta$  و  $\lambda$ .

در نظر گرفته شده است. فرض کنید  $p$  تابع چگالی متناظر با تابع توزیع  $P$  است. برای پارامترهای  $\alpha$ ،  $\beta$  و  $\lambda$  به ترتیب توزیع‌های پیشینی  $\alpha|\delta \sim Exp\{\frac{1}{\delta}\}$ ،  $\beta|\phi \sim U(\circ, \phi)$  و  $\lambda|\gamma \sim Exp\{\frac{1}{\gamma}\}$  در نظر گرفته می‌شود که در آن نمادهای  $Exp$  و  $U$  به ترتیب برای بیان توزیع‌های نمایی و یکنواخت بکار رفته‌اند. با فرض استقلال این پیشین‌ها تابع چگالی پیشین توام پارامترها به صورت  $p_*(\alpha, \beta, \lambda|\delta, \phi, \gamma) = \pi(\alpha|\delta)\pi(\beta|\phi)\pi(\lambda|\gamma)$  خواهد بود. برای ابرپارامترهای  $\delta$ ،  $\phi$  و  $\gamma$  به ترتیب توزیع‌های پیشین مزدوجی به صورت  $\delta \sim IG(a_\delta, b_\delta)$ ،  $\phi \sim Par(a_\phi, b_\phi)$  و  $\gamma \sim IG(a_\gamma, b_\gamma)$  انتخاب شده است که در آن نمادهای  $IG$  و  $Par$  به ترتیب برای نشان دادن توزیع‌های گامای وارون و پاراتو بکار می‌روند. اگر  $a_\phi = a_\gamma = a_\delta = 2$ ، آنگاه  $Var(\delta) = Var(\gamma) = Var(\phi) = \infty$  و نتیجه توزیع‌های پیشین مذکور ناآگاهی‌بخش خواهند بود.

سانسور مورد نظر در این تحقیق سانسور از راست است. تحت این نوع سانسور مجموعه داده‌ها به صورت  $(r_i, t_i^*); i = 1, \dots, n$  هستند.  $t_i^* = \min(t_i, C_i)$  زمان بقای (طول عمر) مشاهده شده فرد  $i$ ام است به طوری که  $t_i$  و  $C_i$  به ترتیب زمان بقای واقعی و زمان سانسور فرد  $i$ ام هستند. همچنین  $r_i$  وضعیت سانسور فرد  $i$ ام را مشخص می‌کند. اگر  $r_i = 1$  باشد طول عمر فرد  $i$ ام مشاهده شده در زمان  $t_i$  و اگر  $r_i = 0$  باشد طول عمر فرد  $i$ ام سانسور شده از راست در زمان  $C_i$  است. نمایش سلسله مراتبی DPMM با هسته GIW را می‌توان به صورت

$$t_i^*|\alpha_i, \beta_i, \lambda_i \stackrel{i.i.d.}{\sim} \{k(t_i^*|\alpha_i, \beta_i, \lambda_i)\}^{r_i} \{S(t_i^*|\alpha_i, \beta_i, \lambda_i)\}^{1-r_i},$$

$$(\alpha_i, \beta_i, \lambda_i)|P \stackrel{i.i.d.}{\sim} P, \quad i = 1, \dots, n,$$

$$P|\delta, \phi, \gamma \sim DP(\nu, P_\circ(\cdot|\delta, \phi, \gamma)), \quad p_*(\alpha, \beta, \lambda|\delta, \phi, \gamma) = \pi(\alpha|\delta)\pi(\beta|\phi)\pi(\lambda|\gamma),$$

$$\alpha|\delta \sim Exp\{\frac{1}{\delta}\}, \quad \beta|\phi \sim U(\circ, \phi), \quad \lambda|\gamma \sim Exp\{\frac{1}{\gamma}\}, \quad \delta \sim IG(a_\delta, b_\delta),$$

$$\phi \sim Par(a_\phi, b_\phi), \quad \gamma \sim IG(a_\gamma, b_\gamma), \quad \nu \sim G(a_\nu, b_\nu)$$

نوشت. با توجه به ارائه سلسله مراتبی از مدل پیشنهادی واضح است که  $(\alpha_i, \beta_i, \lambda_i)$  پارامتر مربوط به داده  $i$ ام و  $\{(\alpha_i, \beta_i, \lambda_i); i = 1, \dots, n\}$  نیز مجموعه پارامترهای مربوط به  $n$  داده است. اگر  $n^*$  تعداد بردارهای متمایز در مجموعه اخیر فرض شود آنگاه مجموعه بردارهای متمایز به صورت  $\{(\alpha_j^*, \beta_j^*, \lambda_j^*); j = 1, \dots, n^*\}$  خواهد بود. در واقع  $n^*$  تعداد خوشه‌ها و  $(\alpha_j^*, \beta_j^*, \lambda_j^*)$  پارامتر مربوط به خوشه  $j$ ام است. با توجه به رابطه بین مقادیر  $(\alpha_i, \beta_i, \lambda_i)$  و  $(\alpha_j^*, \beta_j^*, \lambda_j^*)$  وضعیت خوشه‌بندی  $n$  داده مشخص می‌شود. فرض کنید بردار  $L = (L_1, \dots, L_n)$  وضعیت خوشه‌بندی  $n$  داده را نشان دهد. معمولاً از این بردار به عنوان بردار پیکربندی<sup>۱</sup> یاد می‌شود. داده  $i$ ام در خوشه  $j$ ام قرار می‌گیرد، یعنی  $L_i = j$  است اگر و تنها اگر رابطه  $(\alpha_i, \beta_i, \lambda_i) = (\alpha_j^*, \beta_j^*, \lambda_j^*)$  برقرار باشد. برای ادامه کار بردارهای  $\alpha, \beta, \lambda, r$  و  $t^*$  به ترتیب به صورت  $\alpha = (\alpha_1, \dots, \alpha_n)$ ،  $\beta = (\beta_1, \dots, \beta_n)$ ،  $\lambda = (\lambda_1, \dots, \lambda_n)$ ،  $r = (r_1, \dots, r_n)$ ،  $t^* = (t_1^*, \dots, t_n^*)$  تعریف می‌شوند. همچنین  $\alpha_{-i}$ ،  $\beta_{-i}$  و  $\lambda_{-i}$  بیانگر بردارهای  $\alpha$ ،  $\beta$  و  $\lambda$  بدون  $i$  امین مولفه آنها هستند. با توجه به آنتونیاک (۱۹۷۴) توزیع پسین توام پارامترها به صورت

$$\pi [P, \alpha, \beta, \lambda, \phi, \delta, \gamma, \nu | t^*, r] \propto \pi [P | \alpha, \beta, \lambda, \phi, \delta, \gamma, \nu] \pi [\alpha, \beta, \lambda, \phi, \delta, \gamma, \nu | t^*, r] \quad (۱)$$

است، که در آن

$$P(\cdot) | \alpha, \beta, \lambda, \phi, \delta, \gamma, \nu \sim DP\left(\nu + n, \frac{\nu}{\nu + n} P_0(\cdot | \phi, \delta, \gamma) + \frac{1}{\nu + n} \sum_{i=1}^n I_{(\alpha_i, \beta_i, \lambda_i)}(\cdot)\right)$$

$$\pi [\alpha, \beta, \lambda, \phi, \delta, \gamma, \nu | t^*, r] \propto \pi [\alpha, \beta, \lambda | \phi, \delta, \gamma, \nu, t^*, r] \pi [\phi] \pi [\delta] \pi [\gamma] \pi [\nu]$$

$$\pi [\alpha, \beta, \lambda | \phi, \delta, \gamma, \nu, t^*, r] \propto \prod_{i=1}^n \pi [\alpha_i, \beta_i, \lambda_i | \alpha_{-i}, \beta_{-i}, \lambda_{-i}, \nu, \phi, \delta, \gamma]$$

$$\times \{k(t_i | \alpha_i, \beta_i, \lambda_i)\}^{r_i} \{S(C_i | \alpha_i, \beta_i, \lambda_i)\}^{1-r_i}$$

توزیع شرطی  $\pi [\alpha_i, \beta_i, \lambda_i | \alpha_{-i}, \beta_{-i}, \lambda_{-i}, \nu, \phi, \delta, \gamma]$  نیز در بلکول و مککوین (۱۹۷۳) داده شده است. با توجه به رابطه (۱) فرایند تولید نمونه از توزیع پسین توام پارامترها دو مرحله‌ای است. در گام اول با استفاده از الگوریتم نمونه‌گیری گیبز از توزیع  $\pi [\alpha, \beta, \lambda, \phi, \delta, \gamma, \nu | t^*, r]$  نمونه‌های پسینی پارامترها را بدست آورده و آنگاه در گام دوم با توجه به اینکه توزیع شرطی  $P(\cdot) | (\alpha, \beta, \lambda, \phi, \delta, \gamma, \nu)$  یک DP است با بکار بردن ارائه ستورامان (۱۹۹۴) از DP از آن نمونه تولید می‌شود. ارائه مذکور به فرم  $P(\cdot) | \alpha, \beta, \lambda, \phi, \delta, \gamma, \nu \approx \sum_{j=1}^J \omega_j I_{(\alpha_j, \beta_j, \lambda_j)}(\cdot)$  است که در آن  $J, \omega_j; j = 1, \dots, J$  دنباله مقادیر وزن‌های مدل است که برای محاسبه این مقادیر از روش چوب

<sup>۱</sup>Configuration vector



شکنی<sup>۲</sup> استفاده می‌شود. مقادیر  $j = 1, \dots, J$ ;  $(\alpha_j, \beta_j, \lambda_j)$  نیز از توزیع آمیخته

$$\frac{\nu}{\nu+n} P_{\circ}(\cdot | \phi, \delta, \gamma) + \frac{1}{\nu+n} \sum_{i=1}^n I_{(\alpha_i, \beta_i, \lambda_i)}(\cdot)$$

تولید می‌شوند. برای تعیین  $J$ ، (نقطه برش) می‌توان از روش ارائه شده در کوتاس (۲۰۰۶) استفاده کرد یا مشابه چنگ و یوان (۲۰۱۳) مقدار ثابتی مانند ۴۰۰۰ برای آن در نظر گرفت. یکی از مشکلاتی که در برازش DPMM ها با آن روبرو هستیم طولانی بودن زمان اجرای الگوریتم‌های مونت کارلوی زنجیره مارکف<sup>۱</sup> (MCMC) مربوط به این مدل‌ها است. در سال‌های گذشته چندین بسته مانند Nimble، Dpweibull، DPackage و غیره برای این منظور طراحی شده است. بسته مفیدی که بوسیله راس و مارکویک (۲۰۱۸) برای برازش DPMM ها در نرم‌افزار R معرفی شد بسته Dirichletprocess است. اگر چه در این بسته، به صورت پیش فرض توزیع وایبل به عنوان هسته DPMM انتخاب شده است، اما می‌توان هر توزیع دلخواهی را نیز به عنوان هسته DPMM انتخاب کرد. همچنین در بسته مذکور می‌توان با داده‌های سانسور شده نیز کار کرد. در نیل (۲۰۰۰) برای اجرای الگوریتم‌های MCMC در مدل‌های DPM تعداد ۸ الگوریتم داده شده است که الگوریتم شماره ۸ مهم‌تر و پُرکاربردتر از بقیه الگوریتم‌ها است. در بسته Dirichletprocess می‌توان دو الگوریتم ۴ و ۸ مقاله نیل (۲۰۰۰) را بکار برد که در این مقاله از الگوریتم ۸ استفاده شده است. الگوریتم ۸ نیل (۲۰۰۰) بر اساس نمایش کیسه پولیا<sup>۲</sup> از DP است. برای دستیابی به اطلاعات بیشتر در مورد این ارائه از DP به بلکول و مک‌کوین (۱۹۷۳) رجوع شود. در الگوریتم ۸ نیل (۲۰۰۰) برورسانی مولفه‌های بردار پیکربندی و پارامترهای خوشه‌ها تا زمانی ادامه می‌یابد که همگرایی زنجیره‌های مارکف حاصل شود. در فرایند اجرای الگوریتم نمونه‌گیری گیبز ۴ گام زیر تکرار می‌شود.

#### الگوریتم ۱. نمونه‌گیری گیبز :

- گام ۱- برورسانی مولفه‌های بردار پیکربندی.
- گام ۲- برورسانی پارامترهای خوشه‌ها بر مبنای داده‌هایی که در هر خوشه داریم.
- گام ۳- برورسانی ابرپارامترهای  $\delta$ ،  $\phi$  و  $\gamma$ .
- گام ۴- برورسانی پارامتر  $\nu$ .

برای دستیابی به جزئیات اجرای دو گام اول الگوریتم نمونه‌گیری گیبز به الگوریتم ۸ مقاله نیل (۲۰۰۰) رجوع شود. در فرایند اجرای گام ۳ الگوریتم گیبز باید از توزیع‌های شرطی ابر پارامترهای  $\delta$ ،  $\phi$  و  $\gamma$  نمونه تولید کرد. این

<sup>2</sup>Stick-breaking

<sup>1</sup>Markov Chain Monte Carlo

<sup>2</sup> Poly urn

توزیع‌های شرطی بعد از برخی محاسبات و ساده‌سازی‌ها به صورت

$$\delta|\alpha^*, n^* \sim IG(\text{shape} = a_\delta + n^*, \text{scale} = b_\delta + \sum_{j=1}^{n^*} \alpha_j^*),$$

$$\phi|\beta^*, n^* \sim Par(a_\phi + n^*, \max(b_\phi, \max_{1 \leq j \leq n^*} \beta_j^*)),$$

$$\gamma|\lambda^*, n^* \sim IG(\text{shape} = a_\gamma + n^*, \text{scale} = b_\gamma + \sum_{j=1}^{n^*} \lambda_j^*).$$

هستند، که در آن  $\alpha^* = \{\alpha_j^*; j = 1, \dots, n^*\}$  و  $\beta^* = \{\beta_j^*; j = 1, \dots, n^*\}$  به ترتیب به صورت  $\lambda^* = \{\lambda_j^*; j = 1, \dots, n^*\}$  و  $\omega_j^{(b)}; j = 1, \dots, n^{*(b)}$  استفاده می‌شود. در تکرار  $b$ ام الگوریتم گیبز بر مبنای مقادیر تولید شده  $\{\alpha_j^{*(b)}, \beta_j^{*(b)}, \lambda_j^{*(b)}; j = 1, \dots, n^{*(b)}\}$  و مقادیر وزن‌ها  $\omega_j^{(b)}; j = 1, \dots, n^{*(b)}$ ، توابع چگالی، بقا و نرخ خطر در نقطه  $t$  به صورت

$$f(t) \approx \sum_{j=1}^{n^{*(b)}} \omega_j^{(b)} k(t|\theta_j^{*(b)}), \quad S(t) \approx \sum_{j=1}^{n^{*(b)}} \omega_j^{(b)} S(t|\theta_j^{*(b)}), \quad h(t) \approx \frac{f(t)}{S(t)}.$$

تقریب زده می‌شوند، که در آن  $\theta_j^{*(b)} = (\alpha_j^{*(b)}, \beta_j^{*(b)}, \lambda_j^{*(b)})$  و  $n^{*(b)}$  نیز تعداد خوشه‌ها در تکرار  $b$ ام الگوریتم گیبز است.

## ۴ مطالعه شبیه‌سازی

برای ارزیابی عملکرد مدل پیشنهادی یک مطالعه شبیه‌سازی طراحی می‌شود. برای این منظور از مدل

$$f(t) = \omega_1 f_{GIW}(t|\alpha_1, \beta_1, \lambda_1) + \omega_2 f_{GIW}(t|\alpha_2, \beta_2, \lambda_2), \quad (2)$$

که آمیخته‌ای از دو توزیع GIW است، نمونه‌های در اندازه‌های  $100, 1000$  تولید می‌شود. برای پارامترهای مدل آمیخته (۲) دو مجموعه به صورت

$$I: \alpha_1 = 3, \beta_1 = 2, \lambda_1 = 6, \alpha_2 = 15, \beta_2 = 8, \lambda_2 = 4$$

$$II: \alpha_1 = 1, \beta_1 = 4, \lambda_1 = 8, \alpha_2 = 5, \beta_2 = 8, \lambda_2 = 10$$

نظر گرفته شده است. در ادامه برای ساده‌گی، از مدل (۲) با دو مجموعه  $I$  و  $II$  از پارامترهای این مدل به ترتیب به عنوان مدل‌های  $I$  و  $II$  یاد می‌شود. برای نشان دادن عدم حساسیت استنباط نسبت به انتخاب پارامترهای توزیع پیشین پارامتر  $\nu$  یعنی  $(a_\nu, b_\nu)$  دو مجموعه از مقادیر  $(2, 2)$ ،  $(3, 0.5)$ ،  $(2, 2)$  در نظر گرفته شده‌اند. همچنین با روش آزمایش و خطا این نتیجه حاصل شده است که با انتخاب مقادیر  $b_\delta = 10$ ،  $b_\phi = 10000$  و  $b_\gamma = 100$  نتایج خوبی بدست می‌آید. در فرایند اجرای الگوریتم گیبز، تعداد کل تکرارهای الگوریتم گیبز، تعداد تکرارهای دوره سوخت و فاصله تاخیر به ترتیب برابر با مقادیر  $N = 30000$ ،  $B = 10000$  و  $20$  تعیین شده است و بنابراین حجم نمونه نهایی در استنباط  $1000$  است. برای بررسی همگرایی الگوریتم‌های نمونه‌گیری گیبز از روش‌های گوک، رفتی-لوئیس و هیدلبرگ-ولج استفاده شده است (گلمن و همکاران، ۲۰۱۳). برای دستیابی به نتایج روش‌های مذکور بسته CODA در نرم‌افزار R مورد استفاده قرار گرفته است. نتایج بدست آمده همگرایی الگوریتم‌های نمونه‌گیری گیبز را تایید می‌کنند. در مجموعه داده‌های شبیه‌سازی شده برای اینکه نتایج برآورد توابع چگالی، بقا و نرخ خطر با مقادیر واقعی این توابع مقایسه شود از شاخص‌های اقلیدسی، میانگین قدر مطلق خطا<sup>۱</sup> و هلینگر<sup>۲</sup> استفاده شده است. در جدول‌های ۱ و ۲ سه شاخص مذکور به ترتیب با نمادهای  $d_E$ ،  $d_{MAE}$  و  $d_H$  نشان داده شده‌اند. مقادیر شاخص‌های انحراف در برآورد منحنی‌های واقعی (چگالی، بقا و نرخ خطر) مدل‌های  $I$  و  $II$  به ترتیب در جدول‌های ۱ و ۲ داده شده است. مقادیر انحراف‌های مذکور تحت دو اندازه نمونه  $1000$  و  $100$  تحت دو توزیع پیشین پارامتر  $\nu$  ارائه شده‌اند. با توجه به مقادیر داده شده در جدول ۱ مشاهده می‌شود که برای هر یک از دو اندازه نمونه  $100$  و  $1000$  مقادیر شاخص‌های انحراف تحت دو توزیع پیشین پارامتر  $\nu$  تقریباً یکسان هستند. همین نتیجه با تحلیل مقادیر داده شده در جدول ۲ نیز حاصل می‌شود. در کل با افزایش اندازه نمونه مقادیر شاخص‌ها کاهش یافته است. نتایج نموداری برآورد منحنی‌های واقعی (چگالی، بقا و نرخ خطر) مدل  $I$  تحت پیشین‌های  $G(2, 2)$  و  $G(3, 0.5)$  از پارامتر  $\nu$  به ترتیب در شکل‌های ۲ و ۳ نشان داده شده است. به طور مشابه نتایج برای مدل  $II$  در شکل‌های ۴ و ۵ ارائه شده‌اند. در هر یک از شکل‌های مذکور نتایج مربوط به دو اندازه نمونه  $100$  و  $1000$  به ترتیب در نمودارهای سمت راست و چپ نشان داده شده‌اند. در این نمودارها میانگین‌های پسینی<sup>۳</sup> و کران‌های فاصله اعتبار ۹۵٪ بیزی<sup>۴</sup> منحنی‌های واقعی در نقاط مختلف به ترتیب با رنگ‌های قرمز و مشکی بریده شده نشان داده شده‌اند. همچنین در این نمودارها منحنی‌های واقعی هر دو مدل  $I$  و  $II$  با رنگ آبی نمایش داده می‌شود. با مقایسه نمودارهای ارائه شده در دو سمت چپ و راست شکل‌های ۲، ۳، ۴ و ۵ با یکدیگر مشاهده می‌شود که با افزایش اندازه نمونه مقادیر برآوردها (میانگین‌های پسینی) به مقادیر واقعی منحنی‌های واقعی (چگالی، بقا و نرخ خطر) مدل‌های  $I$  و  $II$  نزدیک‌تر شده‌اند و همچنین طول فواصل اعتبار ۹۵٪ بیزی منحنی‌های واقعی نیز کاهش یافته است. در ادامه هدف مقایسه نتایج نموداری هر یک از دو مدل  $I$  و  $II$  نسبت به دو توزیع پیشین انتخاب شده برای پارامتر  $\nu$  است. نتایج نموداری مدل  $I$  تحت دو توزیع پیشین  $G(2, 2)$  و  $G(3, 0.5)$  پارامتر  $\nu$

<sup>1</sup>Mean Absolute Error

<sup>2</sup>Hellinger

<sup>3</sup>Posterior Mean

<sup>4</sup>Bayesian Credible Interval

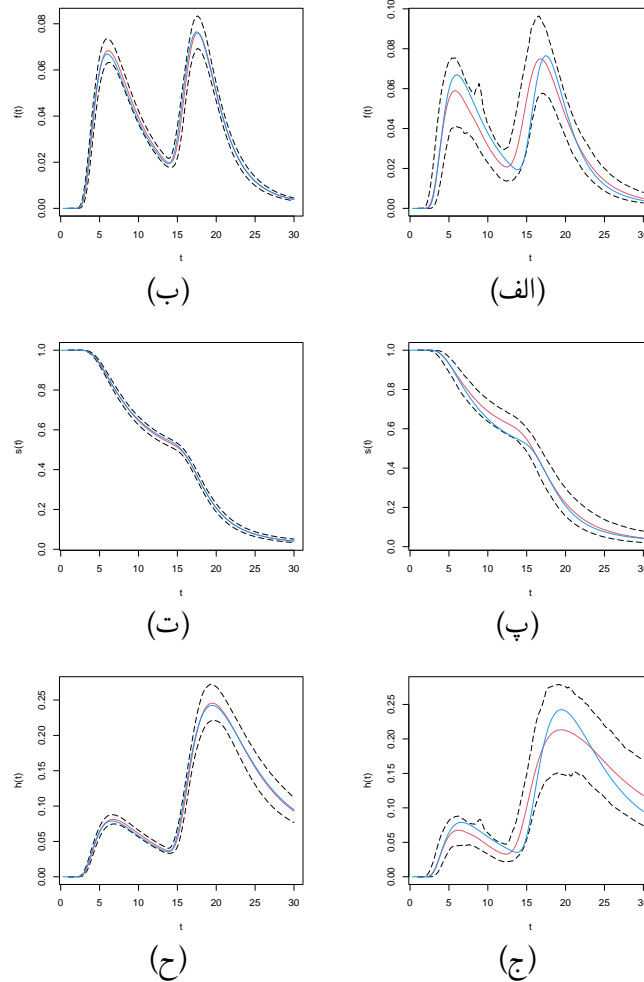
جدول ۱. نتایج محاسبه شاخص‌های انحراف در مدل  $I$  تحت دو اندازه نمونه ۱۰۰ و ۱۰۰۰ و دو توزیع پیشین پارامتر  $\nu$ .

$n = 1000$			$n = 100$			تابع	توزیع پیشین
$d_H$	$d_{MAE}$	$d_E$	$d_H$	$d_{MAE}$	$d_E$		
۰/۰۳۹۸	۰/۰۰۱۰	۰/۰۲۱۵	۰/۱۷۹۷	(۰/۰۰۴۹)	۰/۱۰۰۲	چگالی	$G(2, 2)$
۰/۰۲۳۲	۰/۰۰۲۲	۰/۰۴۹۹	۰/۱۷۴۴	۰/۰۱۶۴	۰/۳۲۶۹	بقا	
۰/۰۴۹۹	۰/۰۰۱۹	۰/۰۳۷۱	۰/۲۶۷۷	۰/۰۱۳۸	۰/۲۳۹۴	نرخ خطر	
۰/۰۳۹۹	۰/۰۰۱۱	۰/۰۲۱۸	۰/۱۷۸۸	۰/۰۰۴۹	۰/۱۰۰۲	چگالی	$G(3, 0.05)$
۰/۰۲۳۹	۰/۰۰۲۳	۰/۰۵۱۱	۰/۱۸۰۶	۰/۰۱۷۱	۰/۳۳۶۷	بقا	
۰/۰۵۲۳	۰/۰۰۲۴	۰/۰۴۱۴	۰/۲۶۵۱	۰/۰۱۳۷	۰/۲۳۸۱	نرخ خطر	

جدول ۲. نتایج محاسبه شاخص‌های انحراف در مدل  $II$  تحت دو اندازه نمونه ۱۰۰ و ۱۰۰۰ و دو توزیع پیشین پارامتر  $\nu$ .

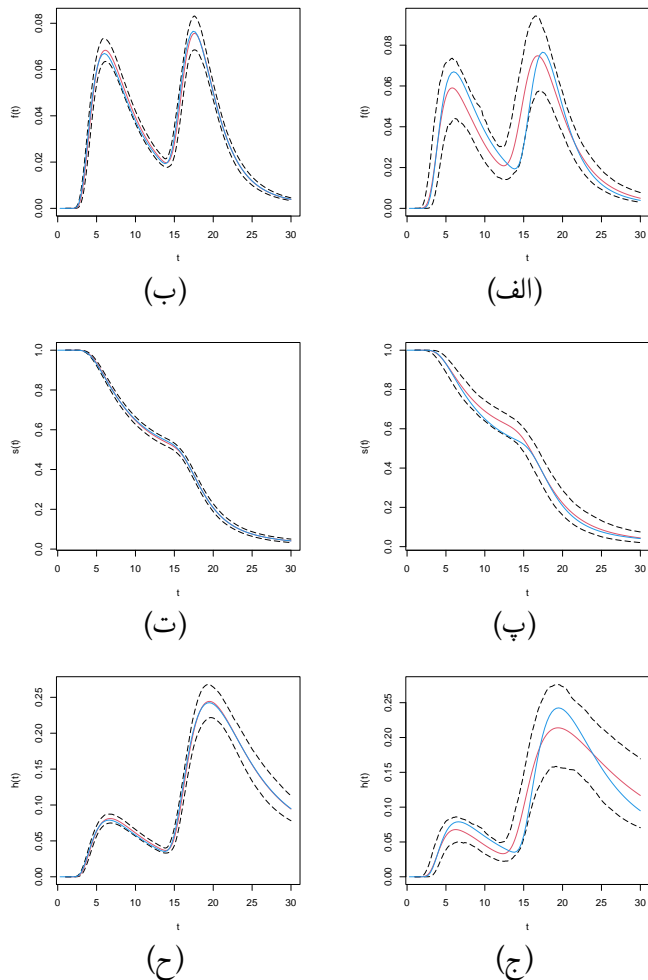
$n = 1000$			$n = 100$			تابع	توزیع پیشین
$d_H$	$d_{MAE}$	$d_E$	$d_H$	$d_{MAE}$	$d_E$		
۰/۰۵۳۶	۰/۰۰۱۷	۰/۰۴۶۵	۰/۱۱۹۴	۰/۰۰۴۵	۰/۱۲۷۲	چگالی	$G(2, 2)$
۰/۰۴۴۲	۰/۰۰۳۲	۰/۰۶۵۳	۰/۰۹۹۵	۰/۰۰۶۴	۰/۱۱۹۴	بقا	
۰/۰۸۶۰	۰/۰۰۷۳	۰/۱۳۲۴	۰/۱۵۵۸	۰/۰۱۲۹	۰/۲۲۹۹	نرخ خطر	
۰/۰۵۳۹	۰/۰۰۱۸	۰/۰۴۴۷	۰/۱۲۶۶	۰/۰۰۴۶	۰/۱۳۶۲	چگالی	$G(3, 0.05)$
۰/۰۴۳۴	۰/۰۰۳۰	۰/۰۶۳۳	۰/۰۹۷۱	۰/۰۰۶۲	۰/۱۱۷۰	بقا	
۰/۱۰۲۸	۰/۰۱۰۲	۰/۱۷۸۲	۰/۱۶۴۴	۰/۰۱۳۵	۰/۲۴۰۲	نرخ خطر	

به ترتیب در شکل‌های ۲ و ۳ داده شده‌اند. با مقایسه نمودارهای در دو شکل اخیر مشاهده می‌شود که استنباط‌های پسینی تحت دو توزیع پیشین پارامتر  $\nu$  تقریباً یکسان هستند. همین نتیجه در مقایسه شکل‌های ۴ و ۵ مربوط به مدل  $II$  نیز برقرار است و بنابراین استنباط‌های پسینی نسبت به توزیع پیشین انتخاب شده برای پارامتر  $\nu$  نوعی عدم حساسیت را نشان می‌دهند. در جدول ۳ استنباط‌های پسینی تعداد خوشه‌های ( $n^*$ ) مدل  $I$  تحت دو توزیع پیشین  $\nu \sim G(2, 2)$  و  $\nu \sim G(3, 0.05)$  ارائه شده است. این استنباط‌ها شامل مقادیر احتمال‌های تجربی، میانه، میانگین، انحراف معیار و فواصل اعتبار ۹۵٪ پسینی  $n^*$  است که تحت اندازه نمونه‌های ۱۰۰، ۳۰۰، ۵۰۰ و ۱۰۰۰ بدست آمده‌اند. نتایج موجود در جدول اخیر بعد از کنار گذاشتن ۱۰۰۰۰ تکرار اولیه در الگوریتم گیبز حاصل شده‌اند. در جدول ۳ با توجه به مقادیر احتمال تجربی  $n^*$  که تحت توزیع پیشین  $\nu \sim G(2, 2)$  بدست آمده‌اند می‌توان



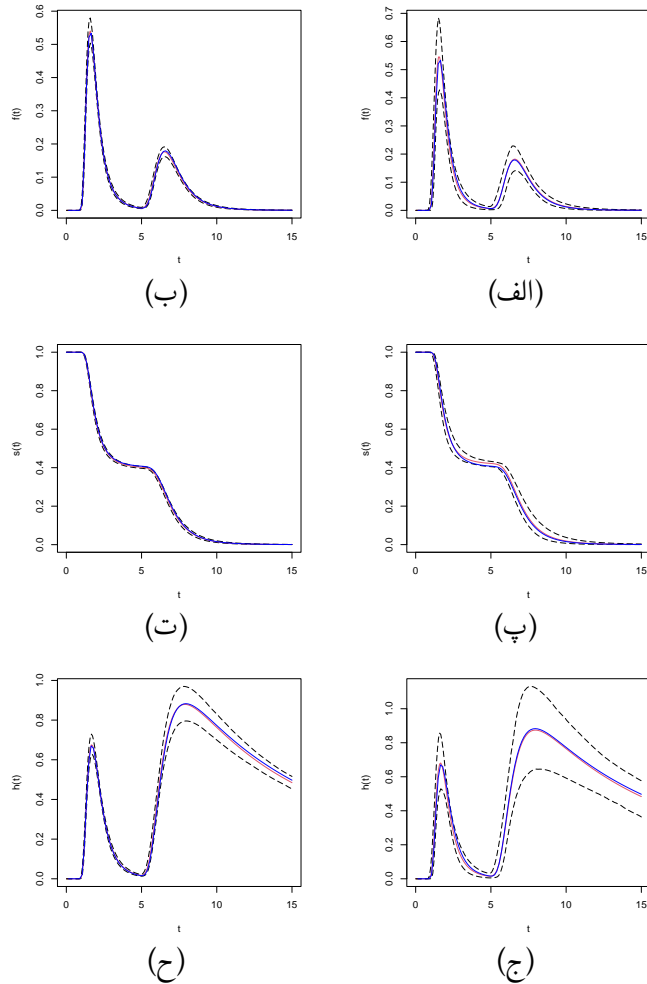
شکل ۲. منحنی برآوردهای نقطه‌ای ( قرمز رنگ) و فاصله‌ای (تیره رنگ) از توابع چگالی ( سطر اول برای الف-  $n = 100$  و ب-  $n = 1000$ ), بقا ( سطر دوم پ-  $n = 100$  و ت-  $n = 1000$ ) و نرخ خطر ( سطر سوم ج-  $n = 100$  و ح-  $n = 1000$ ) مدل  $I$  در نقاط مختلف و تحت پیشین  $G(3, 0.05)$ . منحنی‌های واقعی با رنگ آبی نشان داده شده است.

مشاهده کرد که بیشترین مقدار احتمال تجربی  $n^* = 2$  است. البته با افزایش اندازه نمونه مقدار احتمال مذکور کاهش و مقدار بقیه احتمال‌ها افزایش یافته است. اما در کل بیشترین احتمال تجربی مربوط به  $n^* = 2$  است و بنابراین برآورد ماکسیمم درست‌نمایی  $n^*$  برابر با ۲ یعنی تعداد خوشه‌های واقعی داده‌ها است. نتایج مربوط به محاسبه احتمال تجربی مقادیر مختلف  $n^*$  که تحت توزیع پیشین  $\nu \sim G(3, 0.05)$  بدست آمده‌اند مشابه نتایج تحت توزیع



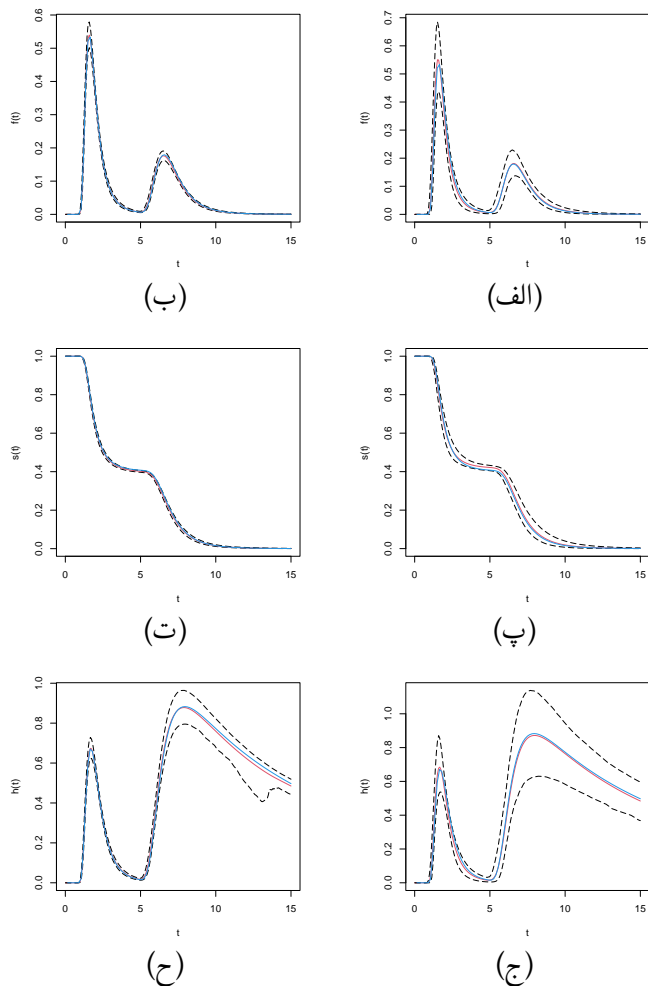
شکل ۳. منحنی برآوردهای نقطه‌ای (قرمز رنگ) و فاصله‌ای (تیره رنگ) از توابع چگالی (منحنی‌های سطر اول برای الف- $n=100$  و ب- $n=1000$ )، بقا (سطر دوم پ- $n=100$  و ت- $n=1000$ ) و نرخ خطر (سطر سوم ج- $n=100$  و ح- $n=1000$ ) مدل  $I$  در نقاط مختلف و تحت پیشین  $G(2, 2)$ . منحنی‌های واقعی با رنگ آبی نشان داده شده است.

پیشین  $\nu \sim G(2, 2)$  هستند اما در حالت کلی مقدار احتمال  $P(n^* = 2)$  تحت توزیع پیشین  $\nu \sim G(2, 2)$  در مقایسه با توزیع پیشین  $\nu \sim G(3, 0.5)$  مقدار بزرگتری است ولی مقادیر بقیه احتمال‌های تجربی تحت توزیع پیشین  $\nu \sim G(3, 0.5)$  بزرگتر هستند. در کل تحت هر دو توزیع پیشین با افزایش اندازه نمونه مقادیر میانگین‌ها و انحراف‌های نمونه پسینی  $n^*$  افزایش یافته است اما مقادیر میانه‌های پسینی و فواصل اعتبار ۹۵٪ پسینی



شکل ۴. منحنی برآوردهای نقطه‌ای (قرمز رنگ) و فاصله‌ای (تیره رنگ) از توابع چگالی (سطر اول برای الف-  $n = 100$  و ب-  $n = 1000$ )، بقا (سطر دوم پ-  $n = 100$  و ت-  $n = 1000$ ) و نرخ خطر (سطر سوم ج-  $n = 100$  و ح-  $n = 1000$ ) در نقاط مختلف و تحت پیشین  $G(3, 0.05)$ . منحنی‌های واقعی با رنگ آبی نشان داده شده است.

$n^*$  تقریباً ثابت هستند. در ادامه این بخش با استفاده از مدل پیشنهادی نتایج خوشه‌بندی داده‌های تولید شده از مدل  $I$  ارائه می‌شود. این نتایج تحت دو توزیع پیشین پارامتر  $\nu$  و تحت اندازه نمونه‌های گوناگون در جدول ۴ داده شده است. نتایج مذکور بر مبنای خوشه‌بندی داده‌ها در آخرین تکرار الگوریتم نمونه‌گیری گیبز بدست آمده‌اند. تعداد خوشه‌های واقعی مدل  $I$  برابر ۲ است و بنابراین در فرایند شبیه‌سازی از این مدل هر داده تولید شده یا متعلق به



شکل ۵. منحنی برآوردهای نقطه‌ای (قرمز رنگ) و فاصله‌ای (تیره رنگ) از توابع چگالی (سطر اول برای الف-  $n = 100$  و ب-  $n = 1000$ )، بقا (سطر دوم پ-  $n = 100$  و ت-  $n = 1000$ ) و نرخ خطر (سطر سوم ج-  $n = 100$  و ح-  $n = 1000$ ) در نقاط مختلف و تحت پیشین  $G(2, 2)$ . منحنی‌های واقعی با رنگ آبی نشان داده شده است.

خوشه ۱ و یا خوشه ۲ است. مدل پیشنهادی نیز هر داده را یا در خوشه واقعی خود قرار می‌دهد و یا در خوشه‌ای غیر از خوشه واقعی داده قرار می‌دهد. نتایج متناظر با توزیع پیشین  $\nu \sim G(3, 0.5)$  داخل پیرانتز نشان داده شده است. در ادامه نتایج متناظر با  $n = 100$  و توزیع پیشین  $\nu \sim G(2, 2)$  تحلیل می‌شوند. تحلیل نتایج به ازای توزیع پیشین  $\nu \sim G(3, 0.5)$  و بقیه اندازه نمونه‌ها مشابه است. در نمونه  $n = 100$  تایی تولید شده از مدل  $I$



جدول ۰۳. استنباط‌های پسینی  $n^*$  تحت دو توزیع پیشین  $\nu$  و اندازه نمونه‌های مختلف در مدل  $I$ .

$P(n^* = 2)$	$P(n^* = 3)$	$P(n^* = 4)$	$P(n^* = 5)$	میانگین	انحراف معیار	فاصله اعتبار ۹۵٪	$n$	پیشین $\nu$
۰٫۸۶۱	۰٫۳۳۸	۰٫۰۰۱	۰٫۰۰۰	۲	۲٫۰۴۰	۰٫۱۹۸	[۲, ۳]	۱۰۰
۰٫۸۸۱	۰٫۱۱۲	۰٫۰۰۶	۰٫۰۰۱	۲	۲/۱۲۵	۰٫۳۵۱	[۲, ۳]	۳۰۰
۰٫۸۰۳	۰٫۰۹۲	۰٫۰۰۵	۰٫۰۰۰	۲	۲/۱۰۲	۰٫۳۲۲	[۲, ۳]	۵۰۰
۰٫۵۵۵	۰٫۴۰۹	۰٫۳۳۴	۰٫۰۰۲	۲	۲/۴۸۲	۰٫۵۷۳	[۲, ۴]	۱۰۰۰
۰٫۸۴۸	۰٫۰۵۱	۰٫۰۰۱	۰٫۰۰۰	۲	۲/۰۵۳	۰٫۲۲۹	[۲, ۳]	۱۰۰
۰٫۸۲۷	۰٫۱۵۹	۰٫۰۱۳	۰٫۰۰۱	۲	۲/۱۸۹	۰٫۴۳۳	[۲, ۳]	۳۰۰
۰٫۸۲۹	۰٫۱۵۵	۰٫۰۱۵	۰٫۰۰۱	۲	۲/۱۸۸	۰٫۴۳۴	[۲, ۳]	۵۰۰
۰٫۴۲۶	۰٫۵۳۶	۰٫۳۳۶	۰٫۰۰۲	۳	۲/۶۱۴	۰٫۵۶۶	[۲, ۴]	۱۰۰۰

تعداد داده‌های خوشه‌های اول و دوم به ترتیب ۵۵ و ۴۵ است که مدل پیشنهادی برای ۵۵ داده خوشه اول در ۴۵ مورد خوشه واقعی داده‌ها را به درستی خوشه ۱ پیش‌بینی کرده است و در ۱۰ مورد به اشتباه خوشه داده‌ها را ۲ پیش‌بینی کرده است. اما مدل پیشنهادی برای ۴۵ داده مربوط به خوشه دوم در ۳۸ مورد خوشه واقعی داده‌ها را به درستی پیش‌بینی کرده است و در ۷ مورد به اشتباه خوشه‌بندی را انجام است. مدل پیشنهادی برای ۱۰۰ داده در ۸۳ مورد به درستی خوشه‌بندی را انجام داده است. بنابراین نرخ درست پیش‌بینی کردن خوشه داده‌ها برابر با ۸۳٪ است. مقدار مذکور تحت توزیع پیشین  $\nu \sim G(3, 0.05)$  برابر با ۸۲٪ است. به همین ترتیب بقیه اعداد و ارقام جدول ۵ می‌تواند تحلیل شود.

جدول ۰۴. نتایج خوشه‌بندی داده‌های تولید شده از مدل  $I$  تحت دو توزیع پیشین  $\nu$  و اندازه نمونه‌های مختلف.

$P(n^* = 2)$	$P(n^* = 3)$	$P(n^* = 4)$	$P(n^* = 5)$	میانگین	انحراف معیار	فاصله اعتبار ۹۵٪	$n$	پیشین $\nu$
۰٫۸۶۱	۰٫۳۳۸	۰٫۰۰۱	۰٫۰۰۰	۲	۲٫۰۴۰	۰٫۱۹۸	[۲, ۳]	۱۰۰
۰٫۸۸۱	۰٫۱۱۲	۰٫۰۰۶	۰٫۰۰۱	۲	۲/۱۲۵	۰٫۳۵۱	[۲, ۳]	۳۰۰
۰٫۸۰۳	۰٫۰۹۲	۰٫۰۰۵	۰٫۰۰۰	۲	۲/۱۰۲	۰٫۳۲۲	[۲, ۳]	۵۰۰
۰٫۵۵۵	۰٫۴۰۹	۰٫۳۳۴	۰٫۰۰۲	۲	۲/۴۸۲	۰٫۵۷۳	[۲, ۴]	۱۰۰۰
۰٫۸۴۸	۰٫۰۵۱	۰٫۰۰۱	۰٫۰۰۰	۲	۲/۰۵۳	۰٫۲۲۹	[۲, ۳]	۱۰۰
۰٫۸۲۷	۰٫۱۵۹	۰٫۰۱۳	۰٫۰۰۱	۲	۲/۱۸۹	۰٫۴۳۳	[۲, ۳]	۳۰۰
۰٫۸۲۹	۰٫۱۵۵	۰٫۰۱۵	۰٫۰۰۱	۲	۲/۱۸۸	۰٫۴۳۴	[۲, ۳]	۵۰۰
۰٫۴۲۶	۰٫۵۳۶	۰٫۳۳۶	۰٫۰۰۲	۳	۲/۶۱۴	۰٫۵۶۶	[۲, ۴]	۱۰۰۰

## ۵ تحلیل مجموعه داده‌های واقعی

در این بخش دو مجموعه داده واقعی تحلیل می‌شود. اولین مجموعه داده تحت عنوان acidity از مجموعه داده‌های بسته نرم‌افزاری BNPdensity در نرم‌افزار R است که مربوط به اندازه‌گیری میزان اسید در نمونه‌های از آب دریاچه‌ها است. در ادامه از این مجموعه داده به عنوان مجموعه داده اسید یاد می‌شود. دومین مجموعه داده تحت عنوان mouse از مجموعه داده‌های نرم‌افزار SPSS است که مربوط به زمان‌های بقای یک مجموعه از موش‌ها است

جدول ۵. نتایج خوشه‌بندی داده‌های تولید شده از مدل  $I$  تحت دو توزیع پیشین  $\nu$  و اندازه نمونه‌های مختلف. مقادیر داخل پرانتز متناظر با پیشین  $G(3, 0.5)$  است.

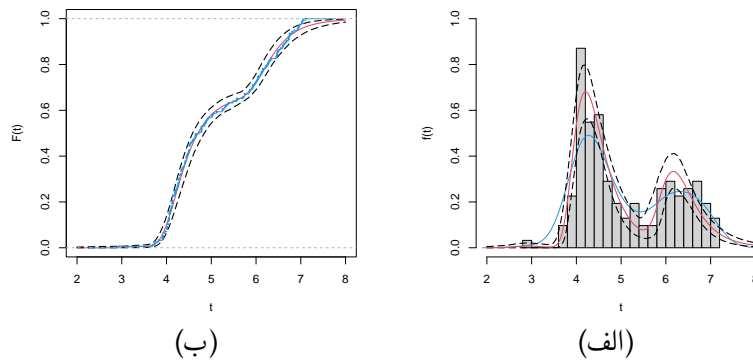
$n$	خوشه واقعی	۱	۲	درصد پیش‌بینی درست خوشه داده‌ها
۱۰۰	۱	۴۵(۴۳)	۱۰(۱۲)	۰.۸۳(۰.۸۲)
	۲	۷(۶)	۳۸(۳۹)	
۳۰۰	۱	۱۵۴(۱۵۰)	۱۵(۱۹)	۰.۸۶(۰.۸۷)
	۲	۲۷(۱۹)	۱۰۴(۱۱۲)	
۵۰۰	۱	۲۴۷(۲۳۳)	۳۵(۴۵)	۰.۸۱(۰.۸۲)
	۲	۵۹(۴۰)	۱۵۹(۱۷۸)	
۱۰۰۰	۱	۵۴۲(۵۲۴)	۷۸(۹۶)	۰.۸۴(۰.۸۳)
	۲	۸۶(۷۹)	۲۹۴(۳۰۱)	

که برخی از آنها سانسور شده از راست هستند. از این مجموعه داده به عنوان مجموعه داده بقای موش‌ها یاد می‌شود.

## ۵.۱ تحلیل مجموعه داده‌های اسید

داده‌های این مثال مربوط به اندازه‌گیری میزان اسید در نمونه‌های از آب ۱۵۵ دریاچه در ایالت ویسکانسین آمریکا است. در این مثال در فرایند اجرای الگوریتم نمونه‌گیری گیبز، تعداد کل تکرارها، تکرارهای دوره سوخت و فاصله تاخیر به ترتیب  $N = 50000$ ،  $B = 10000$  و  $20$  در نظر گرفته شده‌اند. بنابراین حجم نمونه نهایی که در استنباط استفاده می‌شود برابر با ۲۰۰۰ است. در شکل ۶ برآوردهای بیزی (میانگین پسینی) و کران‌های فواصل اعتبار ۹۵٪ پسینی از توابع چگالی و توزیع تجمعی برای مجموعه داده‌های اسید به ترتیب با رنگ‌های قرمز و تیره بریده شده نشان داده شده‌اند. در نمودار سمت راست شکل ۶، بافت نگار داده‌ها به همراه استنباط‌های پسینی برای تابع چگالی رسم شده است. همچنین در این نمودار منحنی آبی رنگ نتایج برآورد چگالی داده‌ها با استفاده از دستور density در نرم‌افزار R را نشان می‌دهد. بافت نگار مجموعه داده‌های اسید دو قله‌ای (مدی) است و بنابراین توزیع‌های پارامتری معمول که عمدتاً تک قله‌ای هستند نمی‌توانند توزیع چنین داده‌های را پوشش دهند، اما شکل مذکور به وضوح نشان می‌دهد که DPMM با هسته GIW توزیع این مجموعه داده را به خوبی پوشش داده‌است.

در نمودار سمت چپ شکل ۶ نتایج برآورد تابع توزیع تجمعی داده‌ها ارائه شده است. در این نمودار منحنی آبی رنگ نتایج برآورد تابع توزیع داده‌ها با روش تابع توزیع تجربی است. نتایج ارائه شده در این نمودار نشان می‌دهد که DPMM با هسته GIW و روش تابع توزیع تجربی تقریباً نتایج یکسانی ارائه داده‌اند. همچنین در این نمودار مشاهده می‌شود کران‌های فواصل اعتبار ۹۵٪ پسینی برآورد تابع توزیع تجمعی در نقاط مختلف خیلی نزدیک بهم هستند. مدل DPM با هسته GIW مجموعه داده‌های اسید را با احتمال‌های  $P(n^* = 2) = 0.0006$ ،  $P(n^* = 3) = 0.9973$  و  $P(n^* = 4) = 0.0021$  به ترتیب در دو، سه و چهار خوشه قرار داده‌است. با توجه

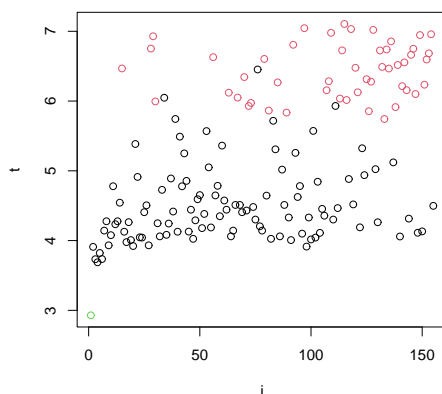


شکل ۶. منحنی برآوردهای نقطه‌ای ( قرمز رنگ) و فاصله‌ای (تیره رنگ) از تابع الف- چگالی در نقاط مختلف به همراه بافت‌نگار مجموعه داده‌های اسید و ب-توزیع تجمعی به همراه نتایج روش تابع توزیع تجربی (رنگ آبی) برای مجموعه داده‌های اسید در نقاط مختلف با مدل پیشنهادی و پیشین  $G(2, 2)$ .

به مقادیر مذکور محتمل‌ترین مقدار (برآورد ماکسیمم درست‌نمایی) برای تعداد خوشه‌ها ۳ است یعنی مدل پیشنهادی داده‌ها را در سه خوشه قرار می‌دهد. در شکل ۷ نتایج خوشه‌بندی داده‌ها با استفاده از مدل پیشنهادی نشان داده شده است. در این نمودار شماره داده‌ها با حرف لاتین  $i$  روی محور  $x$  ها نشان داده می‌شود. با توجه به نتایج ارائه شده در این نمودار مشاهده می‌شود که داده‌ها در سه خوشه قرار گرفته‌اند. داده‌های خوشه‌های اول، دوم و سوم به ترتیب با رنگ‌های مشکی، قرمز و سبز نشان داده شده‌اند. تعداد داده‌های سه خوشه مذکور به ترتیب به صورت ۱۰۵، ۴۹ و ۱ هستند که البته با توجه به نمودار خوشه‌بندی داده‌ها مشاهده می‌شود که در خوشه سوم تنها داده شماره ۱ هست.

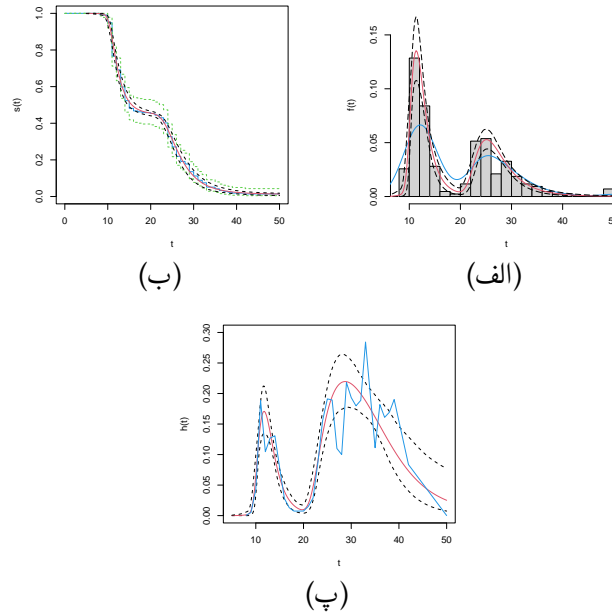
## ۵.۲ تحلیل مجموعه داده‌های بقای موش‌ها

در این زیربخش هدف این است که یکی از مجموعه داده‌های نرم‌افزار SPSS تحت عنوان مجموعه داده‌های بقای موش‌ها تحلیل شود. در این مجموعه داده اطلاعات متغیرهایی مانند زمان بقا، وضعیت سانسور شدن داده‌ها و غیره برای ۲۱۴ موش داده شده است که تنها ۳ تا از داده‌های بقا، سانسور شده از راست هستند. در این مثال در فرایند اجرای الگوریتم نمونه‌گیری گیبز، تعداد کل تکرارها، تکرارهای دوره سوخت و فاصله تاخیر به ترتیب  $N = 20000$ ،  $B = 10000$  و ۲۰ انتخاب شده‌اند. در شکل ۸ برآوردهای بیزی (میانگین پسینی) و همچنین فواصل اعتبار ۹۵٪ بیزی از توابع چگالی، بقا و نرخ خطر داده‌ها در نقاط مختلف زمانی به ترتیب با رنگ‌های قرمز و مشکی بریده شده نشان داده می‌شود. بافت‌نگار زمان بقای موش‌ها همراه با استنباط‌های پسینی برای توابع چگالی در سمت راست سطر بالای شکل ۸ نشان داده شده‌اند. همچنین در این نمودار منحنی آبی رنگ نتایج برآورد چگالی داده‌ها با استفاده از دستور density را نشان می‌دهد. این نمودارها به وضوح نشان می‌دهند که مدل پیشنهادی توزیع داده‌ها را به



شکل ۷. خوشه‌بندی مجموعه داده‌های اسید با استفاده از مدل پیشنهادی و پیشین  $G(2, 2)$ .

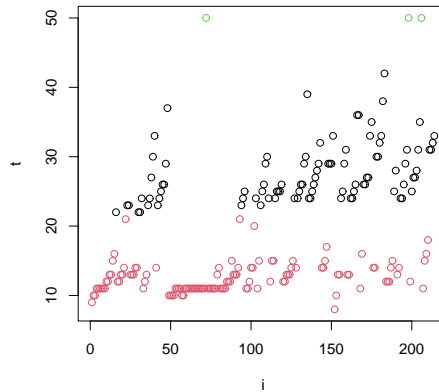
خوبی پوشش داده است. استفاده از روش‌های ناپارامتری مانند کاپلان-مهیر در تحلیل داده‌های بقا متداول است. در صورتی که حجم زیادی از داده‌های بقا سانسور نشده باشند مانند وضعیت موجود در این مثال، روش کاپلان-مهیر نتایج خوبی در برآورد تابع بقا ارائه می‌دهد (لاولس، ۲۰۱۱). در نمودار سمت چپ سطر بالایی شکل ۸ علاوه بر استنباط‌های پسینی تابع بقا، مقادیر برآورد کاپلان-مهیر به همراه کران‌های فواصل اطمینان ۹۵٪ ناپارامتری تابع بقا ارائه شده‌اند. با توجه به این نمودار مشاهده می‌شود که در کل نتایج روش‌های پیشنهادی و کاپلان-مهیر در برآورد تابع بقا خیلی نزدیک بهم هستند. اگرچه طول فواصل اعتبار ۹۵٪ بیزی تابع بقا در مقایسه با طول فاصله اطمینان ۹۵٪ ناپارامتری خیلی کوتاه‌تر است. در سال‌های اخیر در نرم‌افزار R برای برآورد تابع نرخ خطر بسته‌های نرم‌افزاری مختلفی توسعه داده شده است که یکی از این بسته‌ها bshazard است. در این بسته نرم‌افزاری برای دستیابی به یک برآورد ناپارامتری از تابع نرخ خطر در چارچوب مدل‌های آمیخته خطی تعمیم‌یافته از B-اسپلاین‌ها استفاده می‌شود (ربورا و همکاران، ۲۰۱۴). از آنجایی که بسته مذکور معمولاً نتایج خوبی را در برآورد تابع نرخ خطر ارائه می‌کند. بنابراین از این بسته نیز برای برآورد نرخ خطر داده‌ها استفاده شده است. استنباط‌های پسینی مربوط به برآورد تابع نرخ خطر که با DPMM با هسته  $GIW$  بدست آمده‌اند، در آخرین نمودار شکل ۸ ارائه شده‌اند. همچنین در این نمودار نتایج بسته نرم‌افزاری bshazard با رنگ آبی نشان داده شده‌اند. مقایسه نتایج این دو روش نشان می‌دهند که DPMM با هسته  $GIW$  نتایج بهتری فراهم کرده‌است. همچنین در این نمودار مشاهده می‌شود که طول فاصله اعتبار ۹۵٪ بیزی تابع نرخ خطر با افزایش زمان به سرعت افزایش می‌یابد. در شکل ۹ نتایج خوشه‌بندی داده‌های بقای موش‌ها با روش پیشنهادی آورده شده است. این نتایج بر اساس وضعیت خوشه‌بندی داده‌ها در آخرین تکرار الگوریتم گیزی بدست آمده‌اند. همانطور که ملاحظه می‌شود داده‌ها در سه خوشه قرار گرفته‌اند.



شکل ۸. منحنی برآوردهای نقطه‌ای (قرمز رنگ) و فاصله‌ای (تیره رنگ) از تابع الف- چگالی به همراه بافت‌نگار داده‌ها، ب- بقا به همراه برآوردهای کاپلان-مهیر (آبی رنگ) و برآوردهای فاصله‌ای ناپارامتری (سبز رنگ) از تابع بقا و پ- نرخ خطر به همراه نتایج بسته bshazard (آبی رنگ) در برآورد نرخ خطر برای مجموعه داده‌های بقای موش‌ها با مدل پیشنهادی و پیشین  $G(2, 2)$ .

## ۶ بحث و نتیجه‌گیری

در این مقاله برای تحلیل داده‌های بقای سانسور شده از راست، مدل آمیخته‌ای از فرایندهای دیریکله با هسته وارون وایبل تعمیم‌یافته بکار برده شده است. با تحلیل چندین مجموعه داده شبیه‌سازی شده و واقعی، عملکرد مدل پیشنهادی مورد بررسی قرار گرفت. نتایج داده شده در این مقاله نشان می‌دهد مدل پیشنهادی از پتانسیل خوبی برای برآورد توابع چگالی، بقا و نرخ خطر داده‌های بقا برخوردار است. از دیگر مزیت‌های مدل پیشنهادی می‌توان به پتانسیل بالای آن برای خوشه‌بندی کردن داده‌ها اشاره کرد. عدم توانایی مدل پیشنهادی در بررسی تاثیر متغیرهای کمکی روی توزیع طول عمر از جمله نقایص این مدل است. بنابراین در پژوهش‌های آینده هدف نویسندگان این مقاله توسعه دادن مدل پیشنهادی به منظور برطرف کردن نقیصه مذکور است.



شکل ۹. نتایج خوشه‌بندی داده‌های بقای موش‌ها با استفاده از مدل پیشنهادی و پیشین  $G(2, 2)$ .

## تقدیر و تشکر

نویسنده‌های مقاله از راهنمایی‌ها و پیشنهادهای مفید و موثر داوران، سردبیر و ویراستار محترم مجله علوم آماری که سبب ارتقای کیفی مقاله شد، کمال تشکر و قدردانی را دارند.

## مراجع

- Antoniak, C.E. (1974), Mixtures of Dirichlet Processes with Applications to Non-parametric Problems, *The Annals of Statistics*, **2**, 1152-1174.
- Blackwell, D. and MacQueen, J.B. (1973), Ferguson Distributions via Polya Urn Schemes, *The Annals of Statistics*, **1**(2), 353-355.
- Bohlouri Hajjar, S. and Khazaei, S. (2018), Bayesian Nonparametric Survival Analysis Using Mixture of Burr XII Distributions, *Communications in Statistics-Simulation and Computation*, **47**(9), 2724-2738.
- Cheng, N. and Yuan, T. (2013), Nonparametric Bayesian Lifetime Data Analysis Using Dirichlet Process Lognormal Mixture Model, *Naval Research Logistics*, **60**(3), 208-221.

- De Gusmão, F. R., Ortega, E. M. and Cordeiro, G. M. (2011), The Generalized Inverse Weibull Distribution, *Statistical Papers*, **52**, 591-619.
- Escobar, M.D. and West, M. (1995), Bayesian Density Estimation and Inference Using Mixtures, *Journal of the American Statistical Association*, **90**(430), 577-588.
- Ferguson, T.S. (1973), Bayesian Analysis of Some Nonparametric Problems, *The Annals of Statistics*, **1**(2), 209–230.
- Ferguson, T.S. (1974), Prior Distributions on Spaces of Probability Measures, *The Annals of Statistics*, **2**(4), 615-629.
- Ferguson, T.S. and Phadia, E. G. (1979), Bayesian Nonparametric Estimation Based on Censored Data, *The Annals of Statistics*, 163-186.
- Haji Joudaki, B., Hashemi, R. and Khazaei, S. (2022), Survival Analysis Using Dirichlet Process Mixture Model with Three-Parameter Burr XII Distribution as Kernel, *Communications in Statistics-Simulation and Computation*, 1–19. <https://doi.org/10.1080/03610918.2022.2076868>
- Hanson, T. (2006), Modeling Censored Lifetime Data Using a Mixture of Gammas Baseline, *Bayesian Analysis*, **1**(3), 575–594.
- Kaur, K., Mahajan, K.K. and Arora, S. (2018), Bayesian and Semi-Bayesian Estimation of the Parameters of Generalized Inverse Weibull Distribution, *Journal of Modern Applied Statistical Methods*, **17**(1), 1–32.
- Kottas, A. (2006), Nonparametric Bayesian Survival Analysis Using Mixtures of Weibull Distributions, *Statistical Planning and Inference*, **136**(3), 578–596.
- Kottas, A. (2006), Dirichlet Process Mixtures of Beta Distributions, with Applications to Density and Intensity Estimation, *In Workshop on Learning with Non-parametric Bayesian Methods, 23rd International Conference on Machine Learning (ICML)*, **47**.

- Gelman, A., Carlin, J.B., Stern, H.S., Dunson, D.B., Vehtari, A. and Rubin, D.B. (2013), *Bayesian Data Analysis*, Chapman and Hall/CRC.
- Ghosh, J.K. and Ramamoorthi, R.V. (2003), *Bayesian Nonparametrics*, Springer Series in Statistics, Springer-Verlag, New York.
- Kuo, L., Smith, A.F., MacEachern, S. and West, M. (1992), *Bayesian Computations in Survival Models via the Gibbs Sampler*, Springer Netherlands, 11–24.
- Neal, R.M. (2000), Markov Chain Sampling Methods for Dirichlet Process Mixture Models, *Journal of Computational and Graphical Statistics*, **9**(2), 249–265.
- Lahiri, P. and Park, D.H. (2000), Nonparametric Bayes and Empirical Bayes Estimators of Mean Residual Life at Age  $t$ , *Journal of Statistical Planning and Inference*, **29**(1-2), 125–136.
- Lawless, J. F. (2011), *Statistical Models and Methods for Lifetime Data*, John Wiley and Sons, New York.
- Rebora, P., Salim, A. and Reilly, M. (2014), Bshazard: A Flexible Tool for Nonparametric Smoothing of the Hazard Function, *The R Journal*, **6**(2), 114–122.
- Ross, G.J. and Markwick, D. (2018), Dirichletprocess: An R Package for Fitting Complex Bayesian Nonparametric Models.
- Sethuraman, J. (1994), A Constructive Definition of Dirichlet Priors, *Statistica Sinica*, 639–650.
- Susarla, V. and Van Ryzin, J. (1976), Nonparametric Bayesian Estimation of Survival Curves from Incomplete Observations, *Journal of the American Statistical Association*, **71**, 897–902.