

## مقایسه روش‌های مناسب جانه‌ی مقادیر سانسور شده در داده‌های ژئوشیمیایی

سید علی حسینی<sup>۱</sup>، سمانه افتخاری مهابادی<sup>۲</sup>، امید اصغری<sup>۳\*</sup>

۱- دانشجوی دکتری، آزمایشگاه شبیه‌سازی و پردازش داده، دانشکده مهندسی معدن، دانشگاه تهران

۲- استادیار آمار، دانشکده ریاضی، آمار و علوم کامپیوتر، دانشگاه تهران

۳- استادیار، آزمایشگاه شبیه‌سازی و پردازش داده، دانشکده مهندسی معدن، دانشگاه تهران

(دریافت: اسفند ۱۳۹۳ پذیرش: مهر ۱۳۹۴)

### چکیده

در این تحقیق به بررسی روش‌های جانه‌ی مقادیر سانسور شده در مجموعه داده‌های چند متغیره ژئوشیمیایی پرداخته شده است. وجود مقادیر گم شده باعث محدودیت در استفاده از اغلب روش‌های آماری همچون تحلیل مولفه‌های اصلی می‌شود. حذف نمونه‌های شامل داده‌های گم شده باعث اریب شدن نتایج و از دست دادن اطلاعات می‌شود به همین دلیل در نظر گرفتن رویکردی مناسب در مواجهه با داده‌های گم شده یک نیاز اساسی در تحلیل مجموعه داده‌های ناکامل است. در این مقاله، با توجه به ماهیت ترکیبی داده‌های ژئوشیمیایی، چند روش مناسب برای جانه‌ی مقادیر گم شده که در چند سال اخیر ارائه شده‌اند و به سادگی در محیط نرم‌افزار آماری *R* قابل اجرا هستند، معرفی شده‌اند. در نهایت با استفاده از یک مجموعه داده کامل مربوط به منطقه ظرف‌قند، این روش‌ها با یکدیگر مقایسه شده‌اند. نتایج نشان می‌دهند که استفاده از روش‌های چند متغیره برای جانه‌ی و به طور خاص روش *ilr-EM* نسبت به دیگر روش‌ها ارجحیت دارند.

### واژگان کلیدی

داده‌های ژئوشیمیایی، مقادیر سانسور شده، روش‌های جانه‌ی، ماهیت ترکیبی، روش *ilr-EM*

### ارجاع به این مقاله:

حسینی، ع. افتخاری مهابادی، س. اصغری، الف. (۱۳۹۴)، مقایسه روش‌های مناسب جانه‌ی مقادیر سانسور شده در داده‌های ژئوشیمیایی، روش‌های تحلیلی و عددی در مهندسی معدن، ۵(۹)، ۶۳-۷۲.

[http://dx.medra.org/10.17383/S2251-6565\(15\)940916-X](http://dx.medra.org/10.17383/S2251-6565(15)940916-X)

\* عهدہ دار مکاتبات: [o.asghari@ut.ac.ir](mailto:o.asghari@ut.ac.ir)

**1- مقدمه**

جانبه‌ی مقادیر گم شده در مجموعه داده‌ها ترکیبی ارائه شده است که می‌توان از آنها برای کامل کردن مجموعه داده‌های ژئوشیمیایی و تجزیه و تحلیل آنها بهره گرفت. در این مقاله ابتدا چندین روش مناسب برای استفاده در تجزیه و تحلیل داده‌های ژئوشیمیایی مرور می‌شود، از آن جمله می‌توان به روش معرفی شده توسط هرون و همکاران<sup>[4]</sup> که اخیراً در تجزیه و تحلیل داده‌های ژئوشیمیایی نیز به کار گرفته شده است، روش جایگزینی ضربی ساده<sup>[5, 6]</sup> و جایگزینی لاغ نرمال ضربی<sup>[7]</sup> و الگوریتم<sup>[9]</sup>  $EM^{10}$  در فضای تبدیل یافته لگاریتمی<sup>[8, 9]</sup> اشاره نمود. سپس در ادامه با استفاده از یک مجموعه داده کامل ژئوشیمیایی مربوط به منطقه ظفرقند، مجموعه داده‌های جدید که شامل مقادیر سانسور شده با نرخ‌های متفاوت هستند، ایجاد شده و توسط هر کدام از روش‌ها عملیات جانبه‌ی صورت گرفته و در نهایت دقت هر کدام از روش‌ها با بررسی میزان تغییر در آماره‌های مهم توزیع آنها و مقایسه مقدار واقعی و جانبه‌ی شده، مورد ارزیابی قرار خواهد گرفت.

**2- گم شدگی در مجموعه داده‌های ژئوشیمیایی**

مشکلی که اغلب متخصصان تجزیه و تحلیل داده‌های حاصل از نمونه‌گیری تجربی و اندازه‌گیری عناصر با آن روبرو هستند، این است که غلظت برخی از عناصر در نمونه به صورت "کمتر از"<sup>[11]</sup> و یا به ندرت "بیشتر از"<sup>[12]</sup> یک مقدار مشخص گزارش شده است. در مطالعات ژئوشیمیایی عموماً مجموعه داده‌های در دسترس شامل سلول‌هایی است که به صورت کمتر از حد تشخیص دستگاه اندازه‌گیری<sup>[13]</sup> گزارش شده است. چگونگی رسیدگی و برخورد با مقادیر گم شده برای تجزیه و تحلیل داده‌ها از نگرانی اصلی متخصصان به شمار می‌رود<sup>[10]</sup>. مقادیر کمتر از حد تشخیص دستگاه و یا سانسور در واقع مقادیر گم شده هستند، چرا که اطلاعات کاملی در مورد آنها وجود ندارد. با توجه به این که سانسور شدن و یا عدم ثبت مقدار واقعی عنصر، وابسته به غلظت خود عنصر در نمونه است به این‌گونه گم شدگی مکانیزم گم شدگی غیر تصادفی<sup>[14]</sup>  $NMAR$  اطلاق می‌گردد<sup>[1]</sup>. گم شدگی در این مکانیزم غیر قابل چشم پوشی است و تجزیه و تحلیل داده‌ها به روش CCA باعث اریب شدن نتایج می‌شود به ویژه اگر درصد گم شدگی بیشتر از 5٪ باشد.

داده‌های ژئوشیمیایی به صورت گستردگی در مطالعات اکتشافی و بررسی‌های محیط زیستی مورد استفاده قرار می‌گیرد. در عمل پس از جمع‌آوری این‌گونه داده‌ها، از تحلیل‌های آماری تک متغیره، چند متغیره و روش‌های آمار فضایی به منظور توصیف و تحلیل داده‌های ژئوشیمیایی استفاده می‌شود. اغلب در مجموعه داده‌های ژئوشیمیایی وجود داده‌های گم شده<sup>[1]</sup> موضوعی اجتناب ناپذیر است یکی از دلایل رخداد گم شدگی در این‌گونه داده‌ها مقادیر کوچک‌تر از حد تشخیص دستگاه اندازه‌گیری<sup>[2]</sup> (داده‌های سانسور شده) است.

وجود مقادیر گم شده در یک مجموعه داده باعث محدودیت‌هایی در محاسبات و استفاده از روش‌های آماری همچون تحلیل مولفه‌های اصلی می‌شود چرا که استفاده این روش‌ها نیازمند وجود یک مجموعه داده کامل است. یک راه حل مرسوم برای این مشکل، حذف نمونه‌های دارای گم شدگی از مجموعه داده‌ها است که به آن تحلیل داده‌های کامل<sup>[3]</sup> (CCA)<sup>[3]</sup> گفته می‌شود. استفاده از این روش باعث انحراف در نتایج و از دست دادن اطلاعات می‌شود. می‌توان به جای حذف نمونه‌های دارای گم شدگی، سلول‌های دارای گم شدگی در مجموعه داده‌ها را با یک مقدار مناسب جانبه<sup>[4]</sup> نمود<sup>[1]</sup>.

از طرفی، داده‌های ژئوشیمیایی به صورت سهمی از یک مقدار ثابت به عنوان مثال درصد و یا ppm گزارش می‌شوند. این مقادیر، منحصر به فرد، مثبت و مجموع آنها ثابت است (در صورتی که تمام عناصر موجود در نمونه اندازه‌گیری شوند). بنابراین داده‌های ژئوشیمیایی ماهیت ترکیبی<sup>[5]</sup> دارند<sup>[2]</sup>. این ویژگی داده‌های ژئوشیمیایی موجب محدودیت در استفاده از روش‌های آماری چون آنالیز فاکتوری و آنالیز مولفه اصلی می‌شود. برای تحلیل داده‌های ترکیبی روش‌های مختلفی ارائه شده است که بر پایه هندسه آجیسون<sup>[6]</sup> هستند<sup>[3]</sup>. بنابراین انتظار می‌رود یک روش جانبه‌ی مناسب، ماهیت ترکیبی داده‌های ژئوشیمیایی را نیز در نظر بگیرد.

تقریباً در تمام مقالات و کتاب‌ها به این نکته اشاره شده است که، روش‌های جانبه‌ی داده‌های گم شده نباید ساختار کلی داده‌ها را دچار تغییرات شدید کند. چندین روش برای

دستگاه برای تعیین  $\delta_{ij}$  در رابطه (1) از تخمینی از میانگین هندسی مقادیر کمتر از حد تشخیص دستگاه با فرض توزیع لاغ نرمال برای متغیر زام ترکیب استفاده می‌شود. رابطه مقدار جایگزین بر اساس این روش به صورت زیر است.

$$\hat{x}_{ij} = \begin{cases} \delta_{ij} - \exp(\hat{\mu}_j - \hat{\sigma}_j \hat{\lambda}_{ij}) & \text{if } x_{ij} < DL_{ij} \\ \left(1 - \frac{\sum_{k|x_{ik} < DL_{ij}} \delta_{ik}}{c_i}\right) x_{ij} & \text{if } x_{ij} \geq DL_{ij} \end{cases} \quad (2)$$

که در آن:

$$\hat{\lambda}_{ij} = \frac{\phi\left(\frac{(\ln DL_{ij} - \hat{\mu}_j)}{\hat{\sigma}_j}\right)}{\Phi\left(\frac{(\ln DL_{ij} - \hat{\mu}_j)}{\hat{\sigma}_j}\right)}$$

در رابطه فوق (.φ) و (.Φ) به ترتیبتابع چگالی وتابع تجمعی توزیع نرمال استاندارد،  $\hat{\mu}_j$  و  $\hat{\sigma}_j$  پارامترهای تخمین خورده از الگوریتم مبتتنی بر احتمال<sup>17</sup> برای داده‌های سانسور شده در زامین عنصر هستند.

### 3-3-k- نزدیک‌ترین همسایگی (KNN)<sup>18</sup>

روش جانهی KNN عنوان یکی از روش‌های موفق برای داده‌های چند متغیره استاندارد شناخته شده است[13]. در این روش از یک اندازه فاصله برای پیدا کردن  $k$  مشاهده مشابه برای ترکیب‌های دارای گم‌شدنی استفاده می‌شود و مقادیر گم‌شده را توسط اطلاعات متغیرهای مشاهده شده در نزدیک‌ترین همسایگی آن مقدار جانهی می‌کند. با توجه به ماهیت ترکیبی داده‌ها در این روش از اندازه فاصله  $M_i \subset \{1, \dots, D\}$  آجیsson استفاده می‌شود. فرض کنید  $\{x_{ij}, \dots, x_{iD}\}$  مجموعه‌ای از اندیس‌ها برای نشان دادن سلول‌های دارای گم‌شدنی در سطر ثام ماتریس مشاهدات  $X$ ، نشان می‌دهد. برای این  $O_i = \{1, \dots, D\} \setminus M_i$  اندیس مقادیر مشاهده شده را در  $X_i$  نشان می‌دهد. برای جانهی سلول گم‌شده  $X_{ij}$ ، برای هر  $j \in M_i$ ، تمام اجزاء باقی مانده ترکیب که در عنصر زام و اندیس‌های  $O_i$  مشاهده شده‌اند در نظر گرفته می‌شود و  $k$  همسایه  $x_{ik}, \dots, x_{i1}$  برای نمونه  $X_i$  با استفاده از فاصله آجیsson محاسبه می‌شود. ابتدا این مقادیر همسایه انتخاب شده توسط فاکتور زیر تعديل می‌شوند.

### 3- روش‌های جانهی مقادیر سانسور شده

#### 3-1-3- جایگزینی ضربی ساده<sup>15</sup> (SMR)

یکی از روش‌های تک متغیرهای که اغلب برای جانهی مقادیر کمتر از حد تشخیص دستگاه در مطالعات زیست محیطی استفاده می‌شود روش جایگزینی ضربی ساده است [11] که شامل جایگزین کردن 50٪ و یا 70٪ حد تشخیص دستگاه اندازه‌گیری به جای مقادیر گم‌شده است. پالارا و همکاران [7] با بهره گرفتن از این ایده بر روی داده‌های ترکیبی، نشان داده‌اند که یک ضریب تعديل در جایگزینی مقدار سانسور شده برای حفظ مبانی مربوط به داده‌های ترکیبی مانند قید مجموع ثابت نیاز است. فرض کنید مجموعه داده‌ها (ماتریس  $(X_{n \times D})$  شامل  $n$  نمونه ( $i = 1, \dots, n$ ) و  $D$  عنصر اندازه‌گیری شده ( $j = 1, \dots, D$ )، که اندازه‌گیری مربوط به عنصر زام در نمونه  $i$  را نشان می‌دهد و به طور کامل مشاهده نشده است از رابطه (1) به دست می‌آید:

$$x_{ij} = \begin{cases} \delta_{ij} & \text{if } x_{ij} < DL_{ij} \\ \left(1 - \frac{\sum_{k|x_{ik} < DL_{ij}} \delta_{ik}}{c_i}\right) x_{ij} & \text{if } x_{ij} \geq DL_{ij} \end{cases} \quad (1)$$

که  $c_i = \sum_{j=1}^D x_{ij}$  برابر مجموع عناصر برای نمونه  $i$  ام (به عنوان مثال 1,000,000 ppm گزارش شده باشد) و  $\delta_{ij}$  درصدی از  $DL_{ij}$  است. توجه داشته باشید  $\delta$  و  $DL$  به این خاطر اندیس  $ij$  گرفته اند و این امکان را فراهم می‌کنند که غلظت عناصر در نمونه‌ها از روش‌های مختلفی با  $DL$  متفاوت، اندازه‌گیری شده باشند. فرانانز و همکاران [12] نشان داده‌اند که اگر  $\delta_{ij}$  برابر با 65٪ حد تشخیص دستگاه در نظر گرفته شود، ساختار کوواریانس داده‌ها تغییر نخواهد یافت البته اگر میزان گم‌شدنی کمتر از 10٪ باشد.

#### 3-2- جایگزینی لاغ نرمال ضربی<sup>16</sup> (MLNR)

این روش از ویژگی‌های آماری توزیع داده‌ها با فرض یک مدل احتمالی خاص استفاده می‌کند. معمولاً از توزیع لاغ نرمال برای مدل کردن توزیع داده‌های با چولگی به راست مانند غلظت شیمیابی استفاده می‌شود[11]. در این روش به جای در نظر گرفتن درصدی از حد تشخیص

در  $k$  امین تکرار،  $\hat{\sigma}_{ij}^2$  برآورد واریانس شرطی متغیر  $y_j$  در  $k$  امین تکرار،  $\psi$  بردار مربوط به مقدار حد تشخیص دستگاه بعد از تبدیل  $alr$ ,  $\phi$  و  $\Phi$  به ترتیبتابع چگالی وتابع توزیع تجمعی نرمال استاندارد هستند. بعد از همگرایی الگوریتم مقادیر به دست آمده در آخرین تکرار با استفاده از رابطه (7) تبدیل معکوس می‌شوند:

$$\begin{cases} \hat{x}_{ij} = \frac{\exp(\hat{y}_{ij})}{\sum_{j \neq D} \exp(\hat{y}_{ij}) + 1} \\ \hat{x}_{iD} = \frac{1}{\sum_{j \neq D} \exp(\hat{y}_{ij}) + 1} \end{cases} \quad (7)$$

بعد از اجرای تبدیل معکوس  $alr$  مجموعه داده کامل شده در دست خواهد بود.

### 5-3- الگوریتم DA نسبت لگاریتمی جمعی (alr-DA)

این روش شامل الگوریتم داده افزایی (DA) بر اساس مدل  $ALN$  برویده شده<sup>25</sup> نرمال است. این الگوریتم یک حلقه تکراری بر پایه شبیه‌سازی است که مقادیری را از توزیع پسین<sup>26</sup> مقادیر سانسور شده و پارامترهای مدل  $ALN$  بر اساس معادلات تولید می‌کند. مدلسازی  $ALN$  رگرسیون در فضای تبدیل یافته  $alr$  خواهد بود. بنابراین در  $k$  امین تکرار مقدار سانسور شده در فضای  $alr$  به صورت تصادفی از توزیع زیر جانه‌ی می‌شود:

$$\hat{y}_j^{(k)} \sim TruncNormal(y_i' \hat{\beta}_j^{(k)}, \hat{\sigma}_j^{2(k)}) \quad (8)$$

این طرح شباهت زیادی به الگوریتم بر پایه  $EM$  دارد (بخش 4-3). تفاوت اصلی آنها این است که در روش تخمین پارامترها توسط شبیه‌سازی صورت  $alr-DA$  می‌گیرد. بعد از همگرایی الگوریتم با استفاده از رابطه (7) مقادیر تبدیل معکوس و به واحد اصلی برگردانده می‌شوند.

### 6-3- الگوریتم EM نسبت لگاریتمی ایزومتریک (ilr- EM)<sup>27</sup>

در آخر روش ارائه شده توسط فرناندز و همکاران [9] معرفی می‌شود. در این روش هدف کاهش تأثیر مقادیر دور افتاده بر روی مقادیر جانه‌ی شده است. این روش بر پایه یک مدل نرمال چند متغیره بر روی داده‌های تبدیل یافته  $ilr$  هستند. تبدیل  $ilr$  برای یک ترکیب  $D$ -جزئی  $X_i = [x_{i1}, \dots, x_{iD}]$ ,  $i = 1, \dots, n$  می‌شود:

$$f_{ii_l} = \frac{\text{median}_{o \in O_i} x_{io}}{\text{median}_{o \in O_i} x_{i_l o}} \quad \text{for } l = 1, \dots, k \quad (3)$$

با استفاده از این فاکتور تعدیل به عنوان وزن برای مشاهدات،  $k$  همسایه قابل مقایسه می‌شوند. در نهایت مقدار جانه‌ی برای سلول  $X_{ij}$  برابر است با:

$$x_{ij}^* = \text{median}\{f_{ii_1} x_{i_1 j}, \dots, f_{ii_k} x_{i_k j}\} \quad (4)$$

با توجه به این که در ضریب تعدیل از میانه استفاده شده است این روش در صورت وجود مقادیر دور افتاده<sup>19</sup> مقاوم<sup>20</sup> است.

### 4-3- الگوریتم EM نسبت لگاریتمی جمعی (alr-EM)<sup>21</sup>

در این روش ابتدا از تبدیل نسبت لگاریتمی جمعی به عنوان بخشی از الگوریتم استفاده می‌شود. این تبدیل برای ترکیب  $D$ -جزئی نام  $X_i = [x_{i1}, \dots, x_{iD}]$  با در نظر گرفتن یک جزء  $k$  ام به عنوان مخرج به صورت زیر تعریف می‌شود:

$$y_{ij} = \ln\left(\frac{x_{ij}}{x_{ik}}\right), \quad j \in \{1, \dots, D\} \quad (5)$$

از این تبدیل برای مدل کردن آماری بر اساس مدل جمعی لوژستیک نرمال<sup>23</sup> ( $ALN$ ), که یک مدل چند متغیره با بردار میانگین و ماتریس کوواریانس داده‌های تبدیل یافته توسط رابطه (5) است، استفاده می‌شود. در روش  $alr-EM$  طی یک فرآیند تکراری (الگوریتم  $EM$ ) مقادیر گم شده و پارامترهای مدل  $ALN$  متعاقباً به روز می‌شوند تا الگوریتم به همگرایی برسد. برای شروع الگوریتم پارامترهای مدل باید مقدار دهی اولیه شوند که می‌توان آنها را از مقادیر مشاهده شده محاسبه کرد. در  $k$  امین تکرار یک مقدار گم شده  $y_{ij}$  توسط مقدار مورد انتظار تخمین شرطی  $\hat{y}_{ij}$  با توجه به تخمین‌های  $\hat{\mu}^{(k)}$  و  $\hat{\Sigma}^{(k)}$  جانه‌ی می‌شود:

$$\hat{y}_{ij} = y_i' \hat{\beta}_j^{(k)} = \hat{\sigma}_j^{(k)} \frac{\phi\left(\frac{\psi_{ij} - y_i' \hat{\beta}_j^{(k)}}{\hat{\sigma}_j^{(k)}}\right)}{\Phi\left(\frac{\psi_{ij} - y_i' \hat{\beta}_j^{(k)}}{\hat{\sigma}_j^{(k)}}\right)} \quad (6)$$

که در آن  $y_i'$  بردار ترانهاده مقادیر مشاهده شده (بعد از تبدیل  $alr$ ),  $\hat{\beta}_j^{(k)}$  بردار ضرایب رگرسیونی برآورد شده

نوار و لکانوپلتوتونیک ارومیه- دختر واقع است. مطالعات ژئوشیمیابی کل سنگ نشان می‌دهد که سنگ‌های اصلی تشکیل دهنده منطقه اکثراً از نوع بازالت- آندزیت هستند. مشاهدات صحرایی و میکروسکوپی نشان می‌دهند که در این منطقه تناوبی از ماجماتیسم بی مдал اسیدی و بازی وجود دارد [14].

### 5-1- داده‌های اصلی

مجموعه داده اصلی شامل نتایج آنالیز 250 نمونه جمع‌آوری شده از قسمت‌های مختلف ناحیه بوده است که به روش MASS-ICP صورت گرفته است. این مجموعه داده کامل بوده و دارای مقادیر سانسور شده نیست. نمونه‌ها برای 43 عنصر آنالیز شده است و در جدول 1 عناصر و حد تشخیص دستگاه اندازه‌گیری برای هر کدام آورده شده است.

### 5-2- ایجاد داده‌های سانسور به طور مصنوعی در داده‌های اصلی

به منظور ارزیابی روش‌های اشاره شده برای جانهی داده‌های سانسور شده، بر اساس مکانیزم زیر در مجموعه داده‌های اصلی داده‌های سانسور شده ایجاد شده است:

$$CM = \begin{cases} \text{Censored} & \text{if } x_{ij} \leq n \times DL_{x_j} \\ \text{Observed} & \text{otherwise} \end{cases} \quad (11)$$

در حالت اول با استفاده از رابطه (12) چهار مجموعه داده تولید شده است که فقط عنصر بریلیوم دارای داده‌های سانسور است. به ازای مقادیر  $\{2.5, 4.25, 5, 5.5\}$  در رابطه (12) در داده‌های عنصر بریلیوم، 7/5٪، 15٪، 25٪ و 40٪ داده سانسور به طور مصنوعی ایجاد شده است. در حالت دوم 13 عنصر از 43 عنصر بطور تصادفی انتخاب شده است و با استفاده از رابطه (12) درصدهای متفاوتی از سانسور در آنها ایجاد شده است و یک مجموعه داده‌ای جدید ساخته شده است. 13 عنصری که از 43 عنصر دارای داده سانسور هستند و همچنین درصد سانسور در هر کدام در جدول 2 آورده شده است. در این حالت میانگین درصد سانسور در کل مجموعه داده برابر 15٪ است.

$$ilr(x_{ij}) = z_{ij} = \sqrt{\frac{D-j}{D-j+1}} \ln \frac{\prod_{l=j+1}^D x_{il}}{x_{ij}}, \quad \text{for } j=1, \dots, D-1 \quad (9)$$

و تبدیل معکوس  $ilr$  به صورت زیر است:

$$\begin{cases} x_{i1} = \exp\left(-\sqrt{\frac{D-1}{D}} z_{i1}\right) \\ x_{ij} = \exp\left(\sum_{l=1}^{j-1} \frac{1}{\sqrt{(D-l+1)(D-l)}} z_{il} - \sqrt{\frac{D-j}{D-j+1}} z_{ij}\right) \\ x_{iD} = \exp\left(\sum_{l=1}^{D-1} \frac{1}{\sqrt{(D-l+1)(D-l)}} z_{il}\right) \end{cases} \quad (10)$$

این روش بر اساس یک مدل رگرسیون مقاوم سانسور شده بنا گشته است. مراحل اجرای این الگوریتم را می‌توان به صورت زیر خلاصه کرد.

1. انتخاب یک عنصر ( $x_i$ ) که دارای مقادیر گم شده است و اجرای تبدیل  $ilr$  آن بر اساس رابطه (9).
2. استفاده از رابطه (6) بر اساس رگرسیون مقاوم سانسور شده  $z_1, z_2, \dots, z_{D-1}$  به منظور تخمین سلول‌های گم شده و تکرار طرح  $EM$  تا زمان رسیدن به همگرایی.
3. اجرای تبدیل معکوس با استفاده از رابطه (10).
4. تکرار مراحل 3- برای تمام متغیرهای دارای گم شدگی.

### 4- نحوه اجرای روش‌های جانهی

تمام روش‌های ذکر شده توسط نرم‌افزار  $R$  که یک نرم‌افزار محبوب در محاسبات آماری است به سادگی قابل اجرا است. توابع مورد نیاز برای اجرای روش‌ها در دو بسته  $robcomsitions$  و  $zcompositions$  در این <sup>28</sup> نرم‌افزار در دسترس و قابل استفاده است. توابع مربوط به اجرای روش‌ها به ترتیب آورده شده در بالا عبارت‌اند از:  $impKNNa(.)$   $multLN(.)$   $multRepl(.)$   $impRZilr(.)$   $JrDA(.)$   $JrEM(.)$

### 5- مطالعه موردی (ناحیه ظفر قند)

در این تحقیق به منظور اجرای روش‌های مختلف جانهی برای داده ژئوشیمیابی از داده‌های ناحیه ظفر قند که در جنوب شرقی اردستان در 110 کیلومتری شمال شرق اصفهان است، استفاده شده است. از نظر موقعیت زمین‌شناسی بخشی از زون ایران مرکزی محسوب شده و در

جدول 1: لیست عناصر آنالیز شده و حد تشخیص دستگاه (در واحد ppm) متناظر هر عنصر

عنصر	حد تشخیص	عنصر	حد تشخیص	عنصر	حد تشخیص	عنصر	حد تشخیص
0/5	<i>Th</i>	1	<i>Nb</i>	0/2	<i>Bi</i>	0/001	<i>Au</i>
10	<i>Ti</i>	1	<i>Ni</i>	0/1	<i>Cd</i>	100	<i>Al</i>
0/2	<i>Tl</i>	10	<i>P</i>	1	<i>Ce</i>	100	<i>Ca</i>
0/5	<i>U</i>	1	<i>Pb</i>	1	<i>Co</i>	100	<i>Fe</i>
2	<i>V</i>	1	<i>Rb</i>	1	<i>Cr</i>	100	<i>K</i>
0/5	<i>W</i>	50	<i>S</i>	0/5	<i>Cs</i>	100	<i>Mg</i>
0/5	<i>Y</i>	0/5	<i>Sb</i>	1	<i>Cu</i>	100	<i>Na</i>
0/2	<i>Yb</i>	0/5	<i>Sc</i>	1	<i>La</i>	0/1	<i>Ag</i>
1	<i>Zn</i>	0/5	<i>Sn</i>	1	<i>Li</i>	0/5	<i>As</i>
5	<i>Zr</i>	2	<i>Sr</i>	5	<i>Mn</i>	2	<i>Ba</i>
		0/1	<i>Te</i>	0/5	<i>Mo</i>	0/2	<i>Be</i>

که در آن  $X_{imp}$  ماتریس داده‌های اصلی،  $X_{true}$  ماتریس داده‌های جانبه‌ی شده هستند.

در ابتدا مجموعه داده‌های ساخته شده در دو حالت که نحوه تولید آنها در بالا توضیح داده شد مورد استفاده قرار گرفت و با استفاده از روش‌های گفته شده داده‌های سانسور شده در آنها جانبه‌ی شدن. سپس میزان *NRMSE* برای هر روش در مجموعه داده‌های مختلف محاسبه شد که در جدول 3 آورده شده است.

روش *KNN* در مقایسه با سایر روش‌ها نتایج خوبی ارائه نکرده است که احتمالاً مربوط به مکانیزم گم‌شدگی است و این روش برای جانبه‌ی زمانی که مکانیزم گم‌شدگی *NMAR* است مناسب نیست. روش‌های تک متغیره جایگزینی ضربی ساده و نرمال تقریباً نتایج مشابهی ارائه کرده‌اند.

در بین روش‌های چند متغیره روش *ilr-EM* دقیق‌تری دارد. به منظور مقایسه تغییرات ایجاد شده در آماره‌های توزیع، در شکل 1 میزان خطای برآورد میانگین هندسی و انحراف استاندارد در روش‌های مختلف جانبه‌ی برای عنصر بریلیوم محاسبه شده و نمایش داده شده است. زمانی که درصد داده‌های سانسور شده پایین است روش‌های جانشینی ضربی لاغ نرمال، *alr-EM* و *ilr-EM* نتایج مشابهی ارائه کرده‌اند اما زمانی که درصد داده‌های سانسور بالا باشد روش *ilr-EM* بهترین دقیق‌تر است. با افزایش میزان داده‌های سانسور شده از دقیق‌تر روش‌ها نیز کاسته می‌شود اما میزان کاهش دقیق در روش‌های مختلف

جدول 2: عناصر دارای گم‌شدگی در حالت دوم و درصد سانسور هر کدام

عنصر	درصد سانسور
44/6	<i>Ag</i>
14/34	<i>Au</i>
16/73	<i>Bi</i>
34/66	<i>Co</i>
30/67	<i>Cu</i>
15/14	<i>Li</i>
27/5	<i>Mo</i>
1/6	<i>Na</i>
0/4	<i>Nb</i>
6/7	<i>Te</i>
39/44	<i>U</i>
61/75	<i>W</i>
16/73	<i>Zn</i>

## 6- مقایسه روش‌های مختلف جانبه‌ی

برای ارزیابی روش‌های مختلف اشاره شده برای جانبه‌ی داده‌های سانسور شده معیاری برای مقایسه مقدار واقعی و مقدار جانبه‌ی شده به عنوان ریشه دوم نرمال شده میانگین مربع خطای  $NRMSE^{29}$  در نظر گرفته شده است که به صورت زیر تعریف می‌شود:

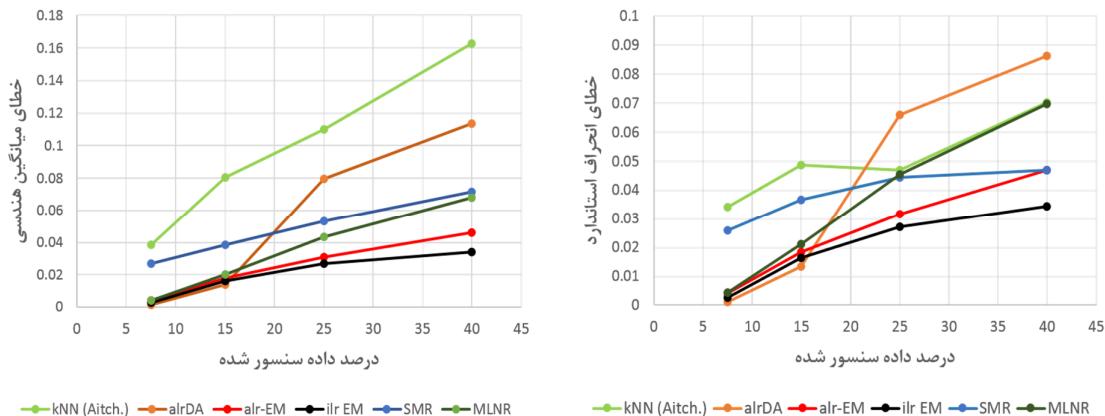
$$NRMSE = \left( \frac{\text{mean}((X_{true} - X_{imp})^2)}{\text{var}(X_{true})} \right)^{1/2} \quad (12)$$

$Q-Q plot$ ,<sup>3</sup> برای بریلیوم داده‌های اصلی و بعد از جانه‌ی مقادیر سانسور شده (در حالتی که ۲۵٪ داده سانسور شده است) به منظور مشاهده بهتر این تغییرات نشان داده شده است. قابل توجه است که از نظر زمان انجام محاسبات رایانه‌ای روش‌های تک متغیره زمان کمتری نیاز دارند.

متفاوت است (شکل ۱). در شکل‌های ۲ هیستوگرام عنصر بریلیوم با داده‌های اصلی و بعد از جانه‌ی مقادیر سانسور شده توسط روش‌های مختلف به منظور مقایسه آورده شده است. همانطور که مشاهده می‌شود روش‌های تک متغیره با جانه‌ی یک مقدار ثابت، باعث تغییر زیادی در هیستوگرام به ویژه زمانی که نرخ سانسور بالا است، می‌شود. در شکل

جدول ۳: NRSME برای روش‌های مختلف جانه‌ی در دو حالت در نظر سانسور

حالت دوم	NRSME				روش جانه‌ی	
	حالات اول					
	%40	%25	%15	7/5%		
50/516	2/738	3/510	5/034	6/053	<i>kNN (Aitch.)</i>	
1/040	2/180	2/482	1/780	1/303	<i>lrDA</i>	
0/530	0/950	1/072	1/047	1/054	<i>alr-EM</i>	
0/504.	0/873	1/035	0/993	1/083	<i>ilr EM</i>	
0/616	1/320	1/410	1/714	2/525	<i>SMR</i>	
0/550	1/210	1/231	1/218	1/082	<i>MLNR</i>	

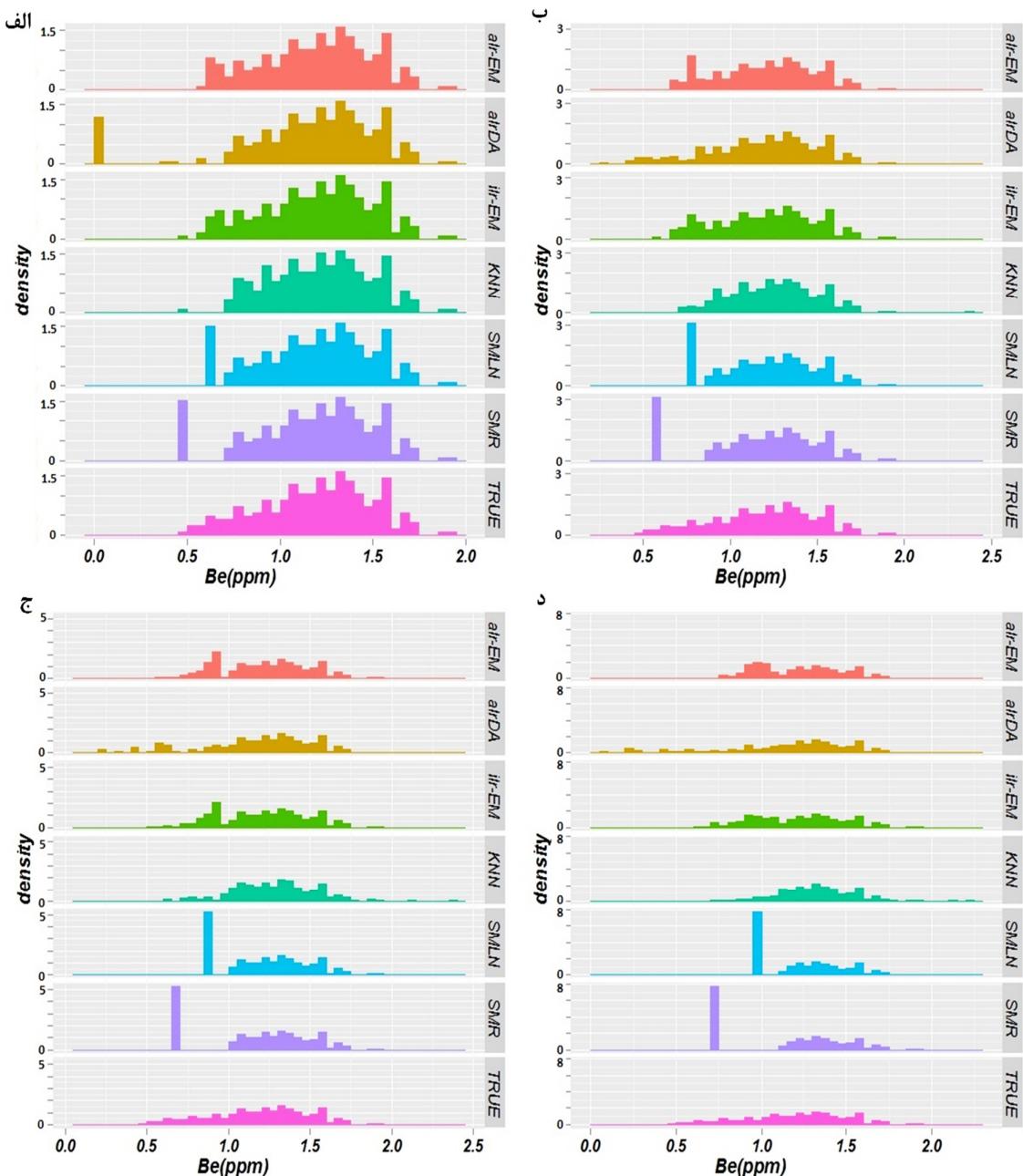


شکل ۱: خطای تخمین میانگین هندسی (چپ) و انحراف استاندارد (راست) برای روش‌های مختلف جانه‌ی در حالت اول

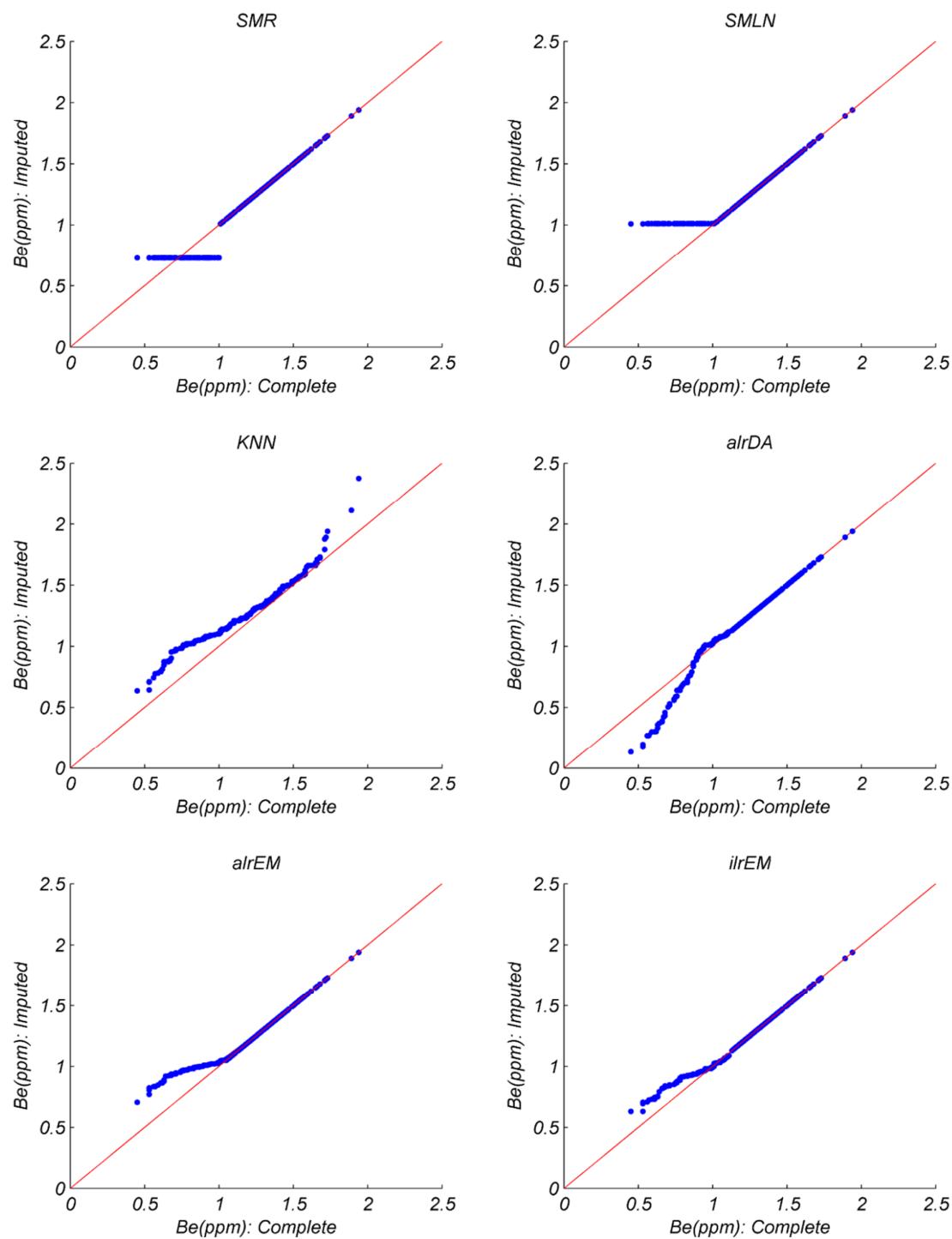
ژئوشیمیابی استفاده کرد. روش‌های تک متغیره مانند جایگزینی ضربی لاغ نرمال زمانی که نسبت داده‌های سانسور پایین باشد (کمتر از ۱۰٪) مناسب است و از نظر محاسبات کامپیوتری با توجه به حجم زیاد داده‌ها زمان کمتری نیاز دارد. بطور کلی روش‌های چند متغیره مدل پایه نتایج بهتری ارائه می‌کنند. استفاده از روش‌های چند متغیره با توجه به این که از روابط بین متغیرها را در نظر می‌گیرند، برای حفظ ساختار کواریانس بین داده‌ها مناسب‌تر است. در بین روش‌های چند متغیره روش *ilr-EM* نتایج بهتری ارائه می‌دهد.

## 7- نتیجه‌گیری

اغلب مجموعه داده‌های ژئوشیمیابی شامل داده‌های سانسور یا گم شده هستند که در برخی مواقع درصد گم شدگی باعث محدودیت استفاده از روش‌های آماری و اریب شدن نتایج تجزیه و تحلیل آنها می‌شود. از طرفی داده‌های ژئوشیمیابی ماهیت ترکیبی دارند که در روش‌های جانه‌ی باید مورد توجه باشد. روش‌های مختلفی برای جانه‌ی داده‌های گم شده برای داده‌هایی که ماهیت ترکیبی دارند ارائه شده است که می‌توان از آنها در داده‌های



شکل 2: هیستوگرام بریلیوم با داده‌های اصلی و بعد از جانبهٔ مقادیر سانسور شده با روش‌های مختلف. (الف) 7/5٪ درصد داده سانسور شده، (ب) 40٪ (ج) 25٪ (د) 15٪ داده سانسور شده.



شکل ۳: بریلیوم داده‌های اصلی و بعد از جانه‌ی مقادیر سانسور شده با روش‌های مختلف در حالتی که ۲۵٪ داده سانسور شده است.

- R. B. (2001). Missing value estimation methods for DNA microarrays. *Bioinformatics*, 17(6), 520-525.
- [14] Sadeghian, M., & Ghaffary, M. (2011). The Petrogenesis of Zafarghand Granitoid Pluton (Se of Ardestan). *Petrology*, (6), 47-70. (In Persian)

**7- مراجع**

- 
- 1- Missing data
  - 2- Below Detection limit
  - 3- Complete Case Analysis
  - 4- Impute
  - 5- Compositional
  - 6- Aitchison
  - 7- Simple multiplicative Replacement
  - 8- Lognormal multiplicative replacement
  - 9- Expectation Maximization
  - 10- Data augmentation
  - 11- Less than
  - 12- Greater than
  - 13- Detection limit
  - 14- Not missing at random
  - 15- Simple multiplicative replacement
  - 16- Multiplicative lognormal replacement
  - 17- likelihood-based
  - 18- k-nearest neighbor
  - 19- Outlier
  - 20- Robust
  - 21- additive logratio Expectation Maximization
  - 22- additive logratio
  - 23- additive logistic normal
  - 24- additive logratio Data augmentation
  - 25- Truncated
  - 26- Posterior distribution
  - 27- isometric logratio Expectation Maximization
  - 28- package
  - 29- Normalized Root Mean Squared Error

- [1] Rubin, D. B., & Little, R. J. (2002). *Statistical analysis with missing data*. Hoboken, NJ: J Wiley & Sons.
- [2] Buccianti, A., & Grunsky, E. (2014). Compositional data analysis in geochemistry: Are we sure to see what really occurs during natural processes?. *Journal of Geochemical Exploration*, 141, 1-5.
- [3] Aitchison, J. (1986). *The Statistical Analysis of Compositional Data*. New York: Chapman and Hall. 416p.
- [4] Hron, K., Templ, M., & Filzmoser, P. (2010). Imputation of missing values for compositional data using classical and robust methods. *Computational Statistics & Data Analysis*, 54(12), 3095-3107.
- [5] de Caritat, P., & Grunsky, E. C. (2013). Defining element associations and inferring geological processes from total element concentrations in Australian catchment outlet sediments: multivariate analysis of continental-scale geochemical data. *Applied Geochemistry*, 33, 104-126.
- [6] Carranza, E. J. M. (2011). Analysis and mapping of geochemical anomalies using logratio-transformed stream sediment data with censored values. *Journal of Geochemical Exploration*, 110(2), 167-185.
- [7] Palarea-Albaladejo, J., & Martín-Fernández, J. A. (2013). Values below detection limit in compositional chemical data. *Analytica chimica acta*, 764, 32-43.
- [8] Palarea-Albaladejo, J., & Martín-Fernández, J. A. (2008). A modified EM alr-algorithm for replacing rounded zeros in compositional data sets. *Computers & Geosciences*, 34(8), 902-917.
- [9] Martín-Fernández, J. A., Hron, K., Templ, M., Filzmoser, P., & Palarea-Albaladejo, J. (2012). Model-based replacement of rounded zeros in compositional data: classical and robust approaches. *Computational Statistics & Data Analysis*, 56(9), 2688-2704.
- [10] Palarea-Albaladejo, J., Martín-Fernández, J. A., & Buccianti, A. (2014). Compositional methods for estimating elemental concentrations below the limit of detection in practice using R. *Journal of Geochemical Exploration*, 141, 71-77.
- [11] Helsel, D. R. (2011). *Statistics for censored environmental data using Minitab and R* (Vol. 77). John Wiley & Sons.
- [12] Martín-Fernández, J. A., Barceló-Vidal, C., & Pawlowsky-Glahn, V. (2003). Dealing with zeros and missing values in compositional data sets using nonparametric imputation. *Mathematical Geology*, 35(3), 253-278.
- [13] Troyanskaya, O., Cantor, M., Sherlock, G., Brown, P., Hastie, T., Tibshirani, R., ... & Altman,

## The Comparison of Appropriate Methods in Imputation of The Censored Values in the Geochemical Datasets

S. A. Hosseini<sup>1</sup>, S. Eftekhari Mahabadi<sup>2</sup>, O. Asghari<sup>3\*</sup>

1- PhD student, Simulation and Data Processing Laboratory, Dept. of Mining, University of Tehran, Iran

2- Assistant Professor, Dept. of Mathematics, Statistics and Computer Science, University of Tehran., Iran

3- Assistant Professor, Simulation and Data Processing Laboratory, Dept. of Mining, University of Tehran, Iran

\* Corresponding Author: [o.asghari@ut.ac.ir](mailto:o.asghari@ut.ac.ir)

(Received: March 2015, Accepted: October 2015)

### Abstract

This study deals with the imputation methods of censored values in the multivariable geochemical data. Presence of the missing values causes limitation in the use of most of statistical methods, e.g. principle component analysis. Excluding the samples which include missing values bias the results and leads to the loss of information. Due to this, consideration of an appropriate approach to deal with missing values is necessary in the analysis of incomplete datasets. In this paper considering the nature of geochemical data, various approaches for imputing the missing values, which have been suggested in the recent years and are easy to be used in the R statistic software, are introduced. Finally, using the complete dataset of the Zafarghand region, these methods are compared with each other. Results show that the application of the multivariable methods in the imputation and particularly the ilr-EM method is preferable to the other methods.

### Keywords

Geochemical Data, Censored Values, Imputation Methods, Compositional Nature, ilr-Em

### Cite This Paper:

Hosseini, S. A., Eftekhari Mahabadi, S., Asghari, O. (2015). "The Comparison of Appropriate Methods in Imputation of the Censored Values in the Geochemical Datasets." Journal of Analytical and Numerical Methods in Mining Engineering 5(9): 63-72. [http://dx.medra.org/10.17383/S2251-6565\(15\)940916-X](http://dx.medra.org/10.17383/S2251-6565(15)940916-X)