

## مقابله با مخاطرات ناشی از غلظت آلاینده PM<sub>2.5</sub> با به کارگیری روش‌های رگرسیون و شباهت مکانی - زمانی و تخمین مقادیر گم‌شده در سری زمانی آنها (مطالعه موردی: شهر تهران)

مرجان فرجی

دانشجوی دکتری مهندسی نقشه‌برداری گرایش سنجش از دور، دانشکده مهندسی عمران و حمل‌ونقل،  
دانشگاه اصفهان

سعید نادی\*

استادیار گروه مهندسی نقشه‌برداری، دانشکده مهندسی عمران و حمل‌ونقل، دانشگاه اصفهان

(تاریخ دریافت ۱۳۹۹/۶/۱۸ - تاریخ پذیرش ۱۳۹۹/۹/۹)

### چکیده

با توجه به تأثیر نامطلوب آلاینده‌ها بر محیط زیست و سلامت انسان، تجزیه و تحلیل داده‌های کیفیت هوا اهمیت زیادی در حفاظت از محیط زیست و رویارویی با مشکلات آلودگی هوا دارد. داده‌های گم‌شده در سری‌های زمانی به‌خصوص داده‌های مربوط به آلودگی هوا موجب بروز چالشی ویژه در برابر آنالیز این داده‌ها می‌شود که ضرورت استفاده از روش‌هایی با عنوان جانپه را برای مقابله با این پدیده نمایان می‌کند. مقادیر گم‌شده، موجب کاهش حجم داده و تغییر الگوهای زمانی موجود در داده‌ها و نتیجه‌گیری اشتباه در تجزیه و تحلیل داده‌ها می‌شود. در این پژوهش به منظور جانپه مقادیر از دست‌رفته در داده‌های سری زمانی غلظت آلاینده PM<sub>2.5</sub> از ۱۲ ایستگاه سنجش آلودگی شهر تهران، روشی ترکیبی بر مبنای رگرسیون جانپه با در نظر گرفتن وابستگی و شباهت‌های مکانی و زمانی بین ایستگاه‌ها توسط الگوریتم پیچش زمانی پویا معرفی شده است. داده‌هایی با مقادیر گم‌شده با الگویی مشابه با داده‌های اصلی در دامنه ۱۰، ۱۵ و ۲۰ درصد گم‌شدگی در داده‌ها با هدف ارزیابی عملکرد مدل‌های جانپه شبیه‌سازی شدند. سپس روش پیشنهادی در ترکیب با روش‌های مختلف جانپه چندگانه همانند روش طبقه‌بندی و رگرسیون درختی، نمونه تصادفی و میانگین تطابق پیش‌بینی کننده، اجرا و نتایج با روش‌های جانپه منفرد مقایسه شد. نتایج بیانگر برتری روش معرفی شده در ترکیب با رگرسیون درختی در مقایسه با دیگر روش‌های جانپه چندگانه و منفرد است.

**واژه‌های کلیدی:** آلاینده PM<sub>2.5</sub>، جانپه منفرد و چندگانه، داده‌های گم‌شده، مخاطرات، معیار شباهت DTW.

## مقدمه

با افزایش صنعتی شدن شهرها، آلودگی هوا به یکی از مخاطرات جدی زیست‌محیطی کلان‌شهرهای جهان از جمله تهران مبدل شده است که گذشته از ضررهای جبران‌ناپذیر برای سلامت شهروندان، تأثیرات اجتماعی و اقتصادی فراوانی بر شهر تحمیل می‌کند [۱، ۴]. با توجه به آثار نامطلوب آلاینده‌ها بر محیط زیست و سلامت انسان، تجزیه و تحلیل داده‌های کیفیت هوا اهمیت زیادی در حفاظت از محیط زیست و مخاطرات آن و رویارویی با مشکلات آلودگی هوا دارد [۱۲]. در طول دهه اخیر شمار زیادی از داده‌های کنترل کیفیت هوا که در بردارنده غلظت آلاینده‌های موجود در جو هستند، توسط ایستگاه‌های سنجش آلودگی در شهرهای مختلف کشور جمع‌آوری شده‌اند که به دلایل مختلف مانند تعمیرات و نگهداری سنجنده‌ها، خطاهای دستگاهی و خطاهای پردازشی دارای گم‌شدگی با فواصل مختلف‌اند [۲۱]. این مقادیر گم‌شده، سبب اشکالاتی در تجزیه و تحلیل داده‌ها می‌شود و تصمیم‌گیری براساس این داده‌ها را با چالش روبه‌رو می‌کند. برای مثال بسیاری از الگوریتم‌های یادگیری ماشین همانند شبکه‌های عصبی و بردار پشتیبان با استفاده از داده‌های کامل به کار گرفته می‌شوند [۷].

داده‌های گم‌شده، پدیده شایعی در مسائل سری زمانی است و معرفی مدل‌ها و روش‌های کارآمد برای مدیریت این خلأ در داده‌ها، گام مؤثری برای کاهش بایاس و افزایش قدرت مدل است. روش‌های رایج که قادر به جانهی مقادیر گم‌شده در سری زمانی هستند، به‌طور معمول به سه دسته حذف داده‌های گم‌شده، جانهی منفرد و جانهی چندگانه تقسیم می‌شوند [۲]. حذف داده‌های گم‌شده از سری زمانی با عنوان حذف کامل<sup>۱</sup> یکی از روش‌های قدیمی و متداول برای مقابله با این مشکل است [۲۲]. این روش سبب از بین رفتن اطلاعات بارزش و کاهش حجم داده می‌شود. تحلیل نتایج در این روش نیز به دلیل تغییر در توزیع داده‌ها و برآورد اریب پارامترها، با خطا مواجه می‌شود. روش جانهی منفرد یکی از رویکردهای رایج برای پر کردن داده‌های جاافتاده است. ساده‌ترین شیوه جانهی منفرد، برآورد مقادیر با روش‌هایی نظیر میانگین، میانه، مد، جانهی تصادفی و میانگین شرط دار برای نقطه‌ای است که در آن مقادیر گم‌شده است [۱۵]. بدین منظور مقداری ثابت برای همه مقادیر از دست‌رفته جایگزین می‌شود و این جانهی، به الگو و ماهیت سری زمانی وابستگی ندارد. این روش‌ها سریع و آسان‌اند، اما موجب تولید پاسخ‌های نامطمئن و بایاس می‌شوند. در روش‌های جانهی چندگانه، مقادیر

---

1. Listwise deletion

گم‌شده با چند مقدار ممکن جایگذاری می‌شوند [۱۶]. از روش‌های به‌نسبت ساده‌ی جانمایی چندگانه می‌توان به روش‌های رگرسیون خطی و غیرخطی [۱۴]، جانمایی چندمتغیره‌ی معادلات زنجیره‌ای<sup>۱</sup> شامل روش‌های تطابق میانگین پیش‌بینی‌کننده<sup>۲</sup> [۲۰]، رگرسیون درختی<sup>۳</sup> [۵] و رگرسیون لجستیک [۱۹]، همچنین روش‌های پیچیده و زمان‌بر مانند بیشینه‌سازی انتظار<sup>۴</sup> [۲۰] و بیزین [۸] اشاره کرد. افزون‌بر روش‌های جانمایی چندگانه، روش‌های نوین همانند روش‌های شبکه‌ی عصبی یادگیری عمیق [۹، ۱۷] نیز امروزه پرکاربرد است.

با توجه به اینکه در بسیاری از موارد الگوی گم‌شدگی در سری زمانی ایستگاه‌ها مشابه هم نیست، در دوره‌های زمانی مواجهه با گم‌شدگی مقادیر در یک ایستگاه، ممکن است در ایستگاه‌های دیگر این گم‌شدگی وجود نداشته باشد. بر این اساس، نوآوری اصلی این مقاله استفاده از الگوی موجود در داده‌های سری زمانی ایستگاه‌های مشابه برای پر کردن داده‌های جافتاده در سری زمانی یک ایستگاه است. بدین منظور در این پژوهش با استفاده از الگوریتم پیش‌زمانی پویا، ایستگاه‌هایی با شباهت بیشینه شناسایی و در مدل‌های جانمایی هم‌زمان پردازش شد تا از الگوی موجود در سری زمانی ایستگاه‌های مشابه برای پر کردن مقادیر جافتاده استفاده شود. نتایج نشان داد که استفاده از روش‌های جانمایی چندگانه بر مبنای روش بیان‌شده در این مقاله، در بازه‌های گم‌شدگی طولانی در مقایسه با روش‌های منفرد برتری دارد.

### منطقه تحقیق و داده‌های استفاده‌شده

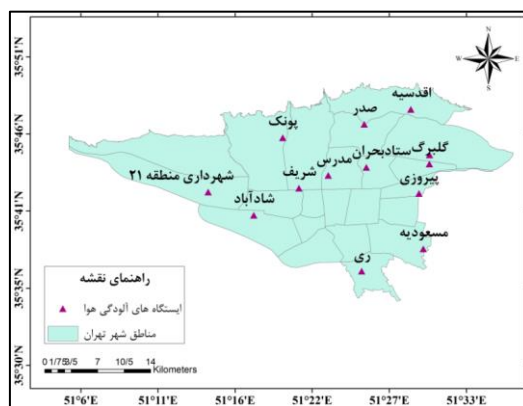
تهران با جمعیتی بالغ بر ۸/۵ میلیون نفر، از شمال به رشته‌کوه البرز و از جنوب به دشت کویر منتهی است. امروزه آلودگی هوای ناشی از ترافیک سنگین، افزایش جمعیت، صنایع و گردوغبار به یکی از مشکلات جدی تهران تبدیل شده است [۲۱]. در این تحقیق، از داده‌های غلظت آلاینده PM<sub>2.5</sub> ثبت‌شده در ۱۲ ایستگاه نظارت بر کیفیت هوا استفاده شده است. داده‌ها به‌صورت ساعتی در دامنه‌ی زمانی ۱۷ آذر ۱۳۹۵ تا ۸ اسفند ۱۳۹۷ از سایت کنترل کیفیت هوا جمع‌آوری شد. اطلاعات آماری ۱۲ ایستگاه سنجش آلودگی همانند میانگین، میانه، انحراف معیار، تعداد مقادیر گم‌شده، و بیشترین و کمترین مقدار داده‌های ساعتی در دامنه‌ی زمانی مورد نظر، در جدول ۱ نشان داده شده است.

1. Multivariate Imputation by Chained Equations (MICE)
2. Predictive mean matching (PMM)
3. Classification and regression trees (Cart)
4. Maximum Exception Maximum Exception

بیشترین و کمترین مقدار گم‌شدگی در داده‌های ایستگاه‌ها متعلق به ایستگاه شادآباد و ایستگاه شهرری به ترتیب با مقدار عددی ۲۳۵۲ و ۵۲۰ داده گم‌شده است. ایستگاه‌های شادآباد، پیروزی، شهرداری منطقه ۲۱ و مسعودیه بیشترین مقادیر غلظت آلاینده را داشتند. موقعیت مکانی این ایستگاه‌ها در شکل ۱ مشاهده می‌شود.

جدول ۱. مشخصات آماری غلظت آلاینده  $PM_{2.5}$  ثبت شده در ایستگاه‌های سنجش آلودگی

نام ایستگاه	انحراف معیار	مقادیر گم‌شده	میانگین	میانه	بیشترین	کمترین
اقدسیه (S1)	۱۳	۱۲۰۵	۲۲	۲۰	۱۸۹	۰
مسعودیه (S2)	۱۳	۱۵۸۷	۳۹	۳۶	۳۸۴	۱
پیروزی (S3)	۲۳	۲۰۲۷	۳۳	۳۸	۲۸۵	۱
ری (S4)	۲۴	۵۲۰	۳۹	۳۴	۳۱۹	۱
شهرداری منطقه ۲۱ (S5)	۲۰	۱۹۰۴	۳۵	۳۰	۳۵۹	۱
شادآباد (S6)	۲۴	۲۳۵۲	۴۰	۳۵	۳۹۳	۱
شریف (S7)	۲۲	۱۸۵۶	۳۸	۳۴	۲۷۸	۱
گلبرگ (S8)	۱۵	۱۴۲۶	۲۰	۲۲	۱۱۷	۰
صدر (S9)	۲۳	۱۰۴۱	۳۸	۳۲	۳۰۸	۱
پونک (S10)	۱۶	۱۹۷۱	۲۴	۲۱	۱۹۲	۱
مدرس (S11)	۲۳	۶۴۸	۳۵	۲۹	۲۰۴	۱
ستاد بحران (S12)	۱۸	۸۲۷	۳۰	۲۷	۱۵۸	۱



شکل ۱. موقعیت مکانی ایستگاه‌های سنجش آلودگی شهر تهران

## روش تحقیق

هدف این مقاله، معرفی روشی نوآورانه بر مبنای شامل کردن وابستگی‌های مکانی بین سری‌های زمانی مربوط به ایستگاه‌های مشابه از دیدگاه رفتار سری زمانی در جانهی اطلاعات گم‌شده مربوط به هر ایستگاه سنجش آلودگی است. در گام اول از طریق الگوریتم پیچش زمانی پویا، شباهت مکانی- زمانی میان سری‌های زمانی غلظت آلاینده PM<sub>2.5</sub> ایستگاه‌ها به صورت دوجه‌دو محاسبه شد. سپس برای جانهی در هر ایستگاه هدف، از وابستگی ایستگاه‌هایی که بیشترین شباهت را با ایستگاه مورد نظر داشتند استفاده می‌شود. در گام دوم، داده کامل اولیه با حذف مقادیر گم‌شده در هر ایستگاه تشکیل می‌شود. در گام بعدی با الگوی مشابه با داده‌های گم‌شده اصلی، داده‌های گم‌شده جدیدی با ۱۰، ۱۵ و ۲۰ درصد گم‌شدگی در داده‌ها حاصل می‌شود. گام چهارم، دربرگیرنده اجرا و مقایسه الگوریتم‌های مختلف جانهی چندگانه و منفرد برای پر کردن داده‌های از دست‌رفته است. در نهایت عملکرد روش‌های مختلف جانهی از طریق شاخص‌های معرفی شده بررسی می‌شود.

### استفاده از الگوریتم پیچش زمانی پویا (DTW)

ایده اصلی در این مقاله این است که با توجه به این موضوع که بسیاری از روش‌های جانهی چندگانه رگرسیون مینا هستند با کنار هم قرار دادن ایستگاه‌های مشابه از لحاظ مکانی- زمانی برای جانهی مقادیر گم‌شده در یک ایستگاه به دقت بیشتری می‌توان دست یافت. برای این منظور در این پژوهش از الگوریتم پیچش پویای زمانی<sup>۱</sup> برای تعیین وابستگی و شباهت مکانی- زمانی بین سری زمانی ثبت‌شده در ایستگاه‌های مختلف استفاده شد. این الگوریتم به منظور اندازه‌گیری شباهت بین دو دنباله از داده‌های سری زمانی استفاده می‌شود. به بیان دیگر این الگوریتم روش مناسبی برای یافتن مسیر بهینه در نقاط بهینه متناظر بین دو سری زمانی است به نحوی که کمترین فاصله مکانی یا بیشترین تشابه بین دو سری زمانی حاصل شود [۱۸].

### استفاده از روش‌های جانهی چندگانه

در این پژوهش از روش‌های جانهی چندگانه در ترکیب با معادلات زنجیره‌ای<sup>۲</sup> استفاده شده است. در الگوریتم معادلات زنجیره‌ای، هر مقدار گم‌شده با دو یا چند مقدار جانهی از طریق روش‌های رگرسیونی پر می‌شود و هر مجموعه از داده‌های به‌دست‌آمده را با استفاده از

1. Dynamic Time Wrapping (DTW)  
2. Multiple Imputation by Chained Equations

روش‌های مربوط به داده‌های کامل تحلیل می‌کنند. داده‌های گم‌شده با داده‌های به‌دست‌آمده در مدل جانمایی جایگزین می‌شود. با تکرار  $m$  بار این عمل،  $m$  مجموعه داده کامل به‌دست می‌آید. هر یک از این مجموعه‌ها جداگانه تحلیل شده و برآوردی از پارامترهای مدل حاصل می‌شود. مقادیر برآوردشده از  $m$  امین مجموعه داده جانمایی شده به‌صورت رابطه ۱ نشان داده می‌شود [۳].

$$\hat{\theta}_m = \frac{1}{m} \sum_{i=1}^m \hat{\theta} \quad (1)$$

که در آن  $\hat{\theta}$  برآوردگر  $\theta$  از مجموعه داده کامل است. برای به‌دست آوردن ضرایب رگرسیون از قوانین ادغام روبین استفاده می‌شود و برآورد پارامترها با میانگین‌گیری از  $m$  برآورد به‌صورت رابطه ۲ حاصل می‌شود [۳].

$$\hat{\beta} = \frac{1}{m} \sum_{k=1}^m \hat{\beta}^{(k)} \quad (2)$$

که در آن  $\hat{\beta}$  ضرایب رگرسیونی حاصل از مدل‌های جانمایی است. در پایان  $m$  نتیجه با محاسبه میانگین، واریانس و دامنه اطمینان متغیر مورد نظر در یک نتیجه تلفیق می‌شود. از مزایای معادلات زنجیره‌ای، امکان چند بار محاسبه و برآورد خطای معیار با استفاده از روشی یکسان است [۳].

### جانمایی منفرد

برخی الگوریتم‌های جانمایی منفرد نیز مبتنی بر ساختار و رفتار سری زمانی هستند. از الگوریتم‌های ساده این دسته می‌توان به روش آخرین مشاهده رو به جلو<sup>۱</sup>، روش میانگین وزن‌دار متحرک<sup>۲</sup> و روش‌های درونیابی<sup>۳</sup> اشاره کرد [۶]. افزون بر آن، در این روش‌ها دیگر یک مقدار ثابت برای همه مقادیر جافتاده جایگزین نمی‌شود. بسته‌های آماری متعددی به‌منظور پیاده‌سازی این الگوریتم‌ها در زبان برنامه‌نویسی R تعبیه شده است که در مقالات به‌کار رفته است [۲۳، ۱۰]. پیاده‌سازی این الگوریتم‌ها نیز سریع و ساده است، اما در مواردی که دامنه اعداد گم‌شده بزرگ باشد، موجب ناپایداری الگوریتم و عدم قطعیت در جواب‌ها می‌شود. [۱۱].

1. Last observation carried forward (LOCF)
2. Weighted moving average
3. Interpolation

### داده ها و بحث

به منظور اعتبارسنجی و بررسی نتایج روش های مختلف جهانی باید در داده های خام اصلی شبیه سازی صورت گیرد. واضح است که عملکرد جهانی افزون بر تعداد داده های از دست رفته به الگوی داده های گم شده نیز وابسته است در این پژوهش فرضیه کاملاً تصادفی بودن داده ها آزمایش شده است. همچنین نرمال و همگن بودن توزیع داده ها نیز بررسی شده است. نتایج حاکی از این است که فرضیه قبول "کاملاً تصادفی" بودن داده ها رد شده و بر این اساس سازوکار داده های مورد نظر "تصادفی" شناخته شده است [۱۴].

در گام بعدی به بررسی رابطه و شباهت مکانی زمانی سری زمانی آلاینده PM<sub>2.5</sub> ایستگاه های سنجش آلودگی هوا بر اساس الگوریتم DTW پرداخته شده است. نتایج اعمال الگوریتم DTW به صورت جدول ۲ مشاهده می شود. در این جدول فواصل DTW دویه دوی همه ایستگاه ها محاسبه شده است. بدین ترتیب فاصله هر دو سری زمانی تعیین می شود. ایستگاه هایی با فواصل DTW کمتر، شباهت بیشتری از منظر رفتار سری زمانی آلودگی ثبت شده در آنها دارند.

جدول ۲. الگوریتم DTW روی آلاینده PM<sub>2.5</sub> در ۱۲ ایستگاه سنجش آلودگی

S12	S11	S10	S9	S8	S7	S6	S5	S4	S3	S2	S1	
۰/۹۵	۱/۴۷	۰/۷۵	۱/۵۷	۰/۷	۱/۳۹	۱/۵۹	۱/۳۶	۱/۷۸	۱/۴۴	۰/۹۵	۰	(S1)
۰/۸۶	۱/۳۴	۰/۶۹	۱/۴۷	۰/۶۸	۱/۲۹	۱/۵۲	۱/۲۹	۱/۷	۱/۳۷	۰	۰/۹۵	(S2)
۰/۹۴	۰/۹۸	۱/۰۶	۰/۹	۱/۲۷	۰/۸۷	۰/۵۹	۱	۱/۱۹	۰	۱/۳۷	۱/۴۴	(S3)
۱/۲	۰/۹۹	۱/۵۵	۱	۱/۶۴	۰/۹۷	۰/۵۹	۱/۰۵	۰	۱/۱۹	۱/۷	۱/۷۸	(S4)
۰/۸۹	۰/۸۳	۱/۲	۰/۸۳	۱/۲۶	۰/۸۲	۰/۸	۰	۱/۰۵	۱/۰۲	۱/۲۹	۱/۳۶	(S5)
۱/۰۱	۰/۷۵	۱/۳۹	۰/۷۸	۱/۴۴	۰/۷۱	۰	۰/۸۲	۰/۹۵	۰/۹۵	۱/۵۲	۱/۵۹	(S6)
۰/۷۸	۰/۶	۱/۱۳	۰/۶۹	۱/۲۶	۰	۰/۷۱	۰/۸۱	۰/۹۸	۰/۸۷	۱/۲۹	۱/۳۹	(S7)
۰/۸۹	۱/۲۹	۰/۷۲	۱/۳۹	۰	۱/۲۶	۱/۴۴	۱/۲۶	۱/۶۴	۱/۲۷	۰/۶۸	۰/۷	(S8)
۰/۸۹	۰/۶۶	۱/۳	۰	۱/۳۹	۰/۶۹	۰/۷۸	۰/۸۳	۱	۰/۹	۱/۴۷	۱/۵۷	(S9)
۰/۷۴	۱/۱۸	۰	۱/۳	۰/۷۲	۱/۱۳	۱/۳۹	۱/۲	۱/۵۵	۱/۰۶	۰/۶۹	۰/۷۵	(S10)
۰/۷۳	۰	۱/۱۸	۰/۶۶	۱/۲۹	۰/۶	۰/۷۵	۰/۸۳	۰/۹۹	۰/۹۸	۱/۳۴	۱/۴۷	(S11)
۰	۰/۷۳	۰/۷۸	۰/۸۹	۰/۸۹	۰/۷۴	۱/۰۱	۰/۸۹	۱/۲	۰/۹۴	۰/۸۶	۰/۹۵	(S12)

مطابق با جدول ۲، ایستگاه های ری (S4) و اقدسیه (S1)، بیشترین فاصله DTW را با مقدار عددی ۱/۷۸ دارند. روی نقشه نیز این دو ایستگاه بیشترین فاصله و در نتیجه کمترین شباهت را با هم دارند. در مقابل ایستگاه های مدرس (S11) و شریف (S7) با فاصله DTW با مقدار ۰/۶ کمترین فاصله و بیشترین شباهت را دارند.

پس از تعیین ایستگاه‌های دارای بیشترین شباهت مکانی-زمانی، شبیه‌سازی داده‌های گم‌شده و تولید داده جدید از داده‌های کامل انجام گرفت. به‌منظور شبیه‌سازی الگوی تصادفی داده‌ها به‌صورت تصادفی بخشی از داده‌ها حذف شده است. با توجه به اینکه سری زمانی داده‌های ساعتی آلودگی دارای فواصلی از داده‌های گم‌شده در بیش از یک ماه نیز هستند، الگوهای فواصل مختلف باید در نظر گرفته شود. بدین منظور بخشی دیگر از داده‌ها به‌صورت فواصل چندروزه و چندساعته به‌صورت تصادفی از داده‌های کامل حذف شد تا الگویی مشابه با الگوی داده از دست‌رفته اصلی ایجاد شود. در نهایت داده‌ها با الگوهای داده‌های از دست‌رفته ۱۰، ۱۵ و ۲۰ درصد شبیه‌سازی شده‌اند. از طریق نتایج به‌دست‌آمده از جدول ۲ که بیانگر فواصل DTW دوبه‌دوی ایستگاه‌هاست، داده‌های سری زمانی ایستگاه هدف در کنار داده‌های سری زمانی سه ایستگاه مشابه آن به‌ترتیب با کمترین فاصله DTW قرار گرفته و ورودی مدل‌های جانمایی چندگانه را تشکیل داده است.

ارزیابی دقت روش‌های جانمایی فرایندی چالش‌انگیز است، زیرا استفاده از روش‌های جانمایی به‌معنای وجود داده‌های ناقص است و شبیه‌سازی این داده‌های گم‌شده با هر دو وابستگی زمانی و مکانی با مشکل مواجه است و حتی بهترین مدل نیز نمی‌تواند به قطعیت دینامیک فرایند تصادفی داده‌ها را به تصویر بکشد [۱۳]. در این مقاله به‌منظور ارزیابی دقت و مقایسه روش‌های مختلف جانمایی چندگانه و منفرد با در دست داشتن داده کامل اولیه و داده جانمایی شده با میزان گم‌شدگی ۱۰، ۱۵ و ۲۰ درصد، از شاخص‌های جذر میانگین مربعات خطا<sup>۱</sup>، ضریب تعیین<sup>۲</sup> و ضریب همبستگی پیرسون<sup>۳</sup> استفاده شد [۲۱].

در این پژوهش برای پیاده‌سازی الگوریتم‌های جانمایی چندگانه همانند روش تطابق میانگین پیش‌بینی‌کننده، طبقه‌بندی و رگرسیون درختی، نمونه تصادفی و همچنین پیاده‌سازی الگوریتم‌های مختلف جانمایی منفرد همانند روش‌های درون‌یابی و آخرین مشاهده رو به جلو از زبان برنامه‌نویسی R استفاده شد. عملکرد روش‌های جانمایی چندگانه طبقه‌بندی رگرسیون درختی، تطابق میانگین پیش‌بینی‌کننده و نمونه تصادفی با الگوی ۱۰، ۱۵ و ۲۰ درصد گم‌شدگی در داده‌ها در جدول ۳ نشان داده شده است. شایان ذکر است که معیارهای ارزیابی داده‌های هر ایستگاه با میانگین سه بار تکرار، جداگانه محاسبه شده و سپس میانگین همه ایستگاه‌ها در جداول مذکور ارائه شده است.

- 
1. RMSE
  2. R Squared
  3. R coefficient



جدول ۳. عملکرد روش‌های جانمایی چندگانه با الگوی داده‌گم‌شده ۱۰، ۱۵ و ۲۰ درصد

۱۰ درصد گم‌شدگی			۱۵ درصد گم‌شدگی			۲۰ درصد گم‌شدگی			روش‌های جانمایی
RMSE	r	R <sup>2</sup>	RMSE	r	R <sup>2</sup>	RMSE	r	R <sup>2</sup>	چندگانه
۱۱/۷۸	۰/۸	۰/۶۶	۱۰/۸۳	۰/۷۶	۰/۶	۱۱/۷۵	۰/۷۵	۰/۵۸	رگرسیون درختی
۱۱/۰۲	۰/۷۸	۰/۶۴	۱۱/۰۱	۰/۷۵	۰/۵۸	۱۱/۹۳	۰/۷۵	۰/۵۶	تطابق میانگین پیش‌بینی‌کننده
۱۱/۳۷	۰/۷۷	۰/۶۲	۱۱/۲	۰/۷۴	۰/۵۶	۱۱/۹۶	۰/۷۳	۰/۵۵	نمونه تصادفی

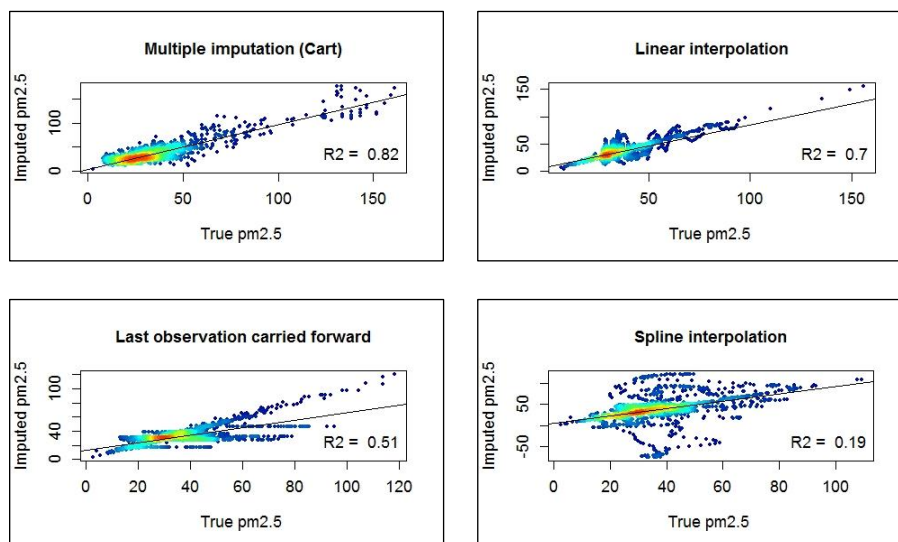
همان‌طور که مشاهده می‌شود، روش جانمایی چندگانه رگرسیون درختی با معیار ضریب تعیین ۰/۶۶ و ضریب همبستگی ۰/۸ در ۱۰ درصد گم‌شدگی داده، ضریب تعیین ۰/۶ و ضریب همبستگی ۰/۷۶ در ۱۵ درصد گم‌شدگی داده‌ها، ضریب تعیین ۰/۵۸ و ضریب همبستگی ۰/۷۵ در ۲۰ درصد گم‌شدگی داده‌ها، در بین روش‌های جانمایی چندگانه بهترین عملکرد را نشان داده است. واضح است که با افزایش درصد گم‌شدگی داده‌ها دقت معیارهای ارزیابی کاهش می‌یابد. در بین ۱۲ ایستگاه استفاده‌شده در این پژوهش، بهترین نتیجه از ترکیب روش جانمایی رگرسیون درختی مطابق با معیار ضریب تعیین، مربوط به ایستگاه مدرس با مقدار عددی ۰/۹۲، ۰/۸۴، ۰/۸۲ و بدترین عملکرد روش ذکرشده، مربوط به ایستگاه مسعودیه با مقدار عددی ضریب تعیین ۰/۳۲، ۰/۲۵، ۰/۲۳ به ترتیب با مقدار ۱۰، ۱۵ و ۲۰ درصد گم‌شدگی در داده‌هاست.

مطابق با نتایج ارائه‌شده، روش میانگین تطابق پیش‌بینی‌کننده و روش تصادفی، عملکرد مشابهی نشان داده‌اند و به نسبت روش رگرسیون درختی ضعیف‌تر عمل کرده‌اند. جدول ۴ عملکرد روش‌های جانمایی منفرد با روش‌های درون‌یابی خطی، اسپیلاین و روش آخرین مشاهده رو به جلو را نشان می‌دهد.

جدول ۴. خلاصه عملکرد جانمایی منفرد در سه داده‌گم‌شده ۱۰، ۱۵ و ۲۰ درصد

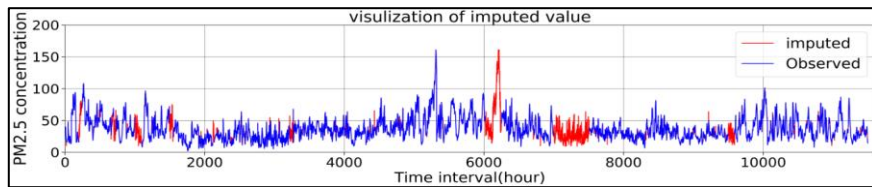
۱۰ درصد گم‌شدگی			۱۵ درصد گم‌شدگی			۲۰ درصد گم‌شدگی			روش‌های جانمایی
RMSE	r	R <sup>2</sup>	RMSE	r	R <sup>2</sup>	RMSE	r	R <sup>2</sup>	خطی
۸/۵۳	۰/۸۲	۰/۶۸	۱۰/۹	۰/۷۵	۰/۵۹	۱۳/۹۴	۰/۶۶	۰/۴۵	درون‌یابی خطی
۳۱/۴۲	۰/۳۴	۰/۱۴	۵۹/۶۳	۰/۱۹	۰/۰۶	۶۷/۰۶	۰/۱۷	۰/۰۳	درون‌یابی اسپیلاین
۱۰/۸۱	۰/۷۳	۰/۵۵	۱۲/۴۵	۰/۶۳	۰/۴۲	۱۶/۳۴	۰/۶	۰/۳۷	آخرین مشاهده رو به جلو

مطابق با جدول ۴، براساس هر سه معیار ارزیابی، روش درون‌یابی خطی از دیگر روش‌های ارائه‌شده بهتر بوده است. از این‌رو در میان روش‌های منفرد برای داده‌های ارائه‌شده، این روش مناسب‌تر است. همچنین روش درون‌یابی اسپیلاین ضعیف‌ترین عملکرد را در بین همه روش‌های جانهی چندگانه و منفرد از خود نشان داده است. اگرچه در داده‌هایی با ۱۰ درصد گم‌شدگی، روش درون‌یابی خطی در مقایسه با روش رگرسیون درختی در شاخص‌های ارزیابی بیشترین ضریب تعیین و همبستگی و کمترین خطا را کسب کرده است، اما باید به این نکته توجه داشت که روش درون‌یابی خطی برای مقادیر گم‌شده با بازه کم، عملکرد بسیار مناسبی نشان می‌دهد، اما هنگامی که بازه گم‌شدگی داده‌ها بیشتر شود، برای مثال در بازه مربوط به گم‌شدگی ۲۰ درصد این روش‌ها قادر به جانهی مناسب برای داده‌های گم‌شده نیستند و یک مقدار ثابت یا با تغییرات کم را برای تمام مقادیر گم‌شده در هر بازه در نظر می‌گیرد. شکل ۲ نمودار پراکندگی مقادیر مشاهداتی در مقابل مقادیر جانهی شده با روش‌های رگرسیون درختی، درون‌یابی خطی اسپیلاین و آخرین مشاهده رو به جلو در ۱۰ درصد گم‌شدگی در ایستگاه منطقه ۲۱ را نشان می‌دهد. این نمودارها ارتباط بین مقادیر مشاهداتی و مقادیر جانهی شده را به صورت یک معادله خطی درجه یک نشان می‌دهد.

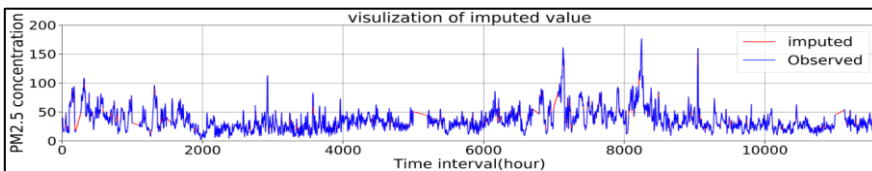


شکل ۲. نمودار مقدار مشاهداتی در مقابل مقدار جانهی‌شده ایستگاه شهرداری منطقه ۲۱

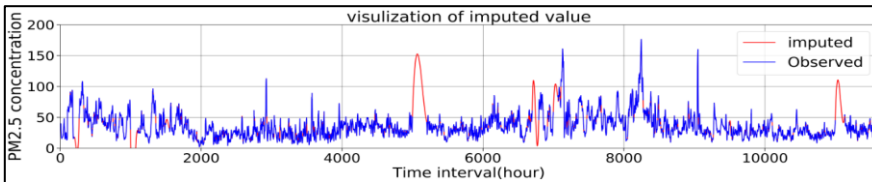
در شبیه‌سازی داده‌های گم‌شده به‌دلیل اینکه بخشی از داده‌ها به‌صورت تصادفی با بازه یک تا سه رکورد حذف شده‌اند، پر کردن این داده‌ها توسط روش درون‌یابی خطی با دقت زیادی انجام گرفته است از این‌رو در مجموع دقت روش درون‌یابی خطی در ۱۰ درصد گم‌شدگی در داده‌ها بیشتر از روش‌های جانپی چندگانه برآورد شده است. شکل ۳ نمایش بصری الگوی حاصل از اعمال روش‌های مختلف جانپی برای کامل کردن گم‌شدگی‌های ۱۰ درصد را نشان می‌دهد.



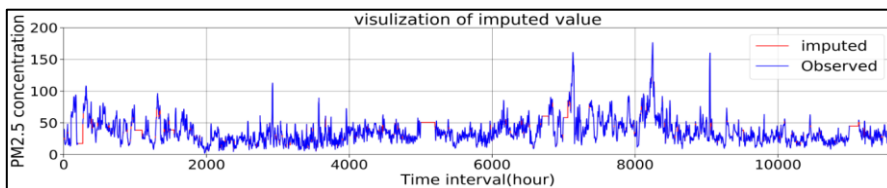
(الف)



(ب)



(پ)



(ج)

شکل ۳. جانپی چندگانه و منفرد مقادیر PM<sub>2.5</sub> ایستگاه شهرداری منطقه ۲۱ با گم‌شدگی ۱۰ درصد داده‌ها: (الف) جانپی چندگانه روش رگرسیون درختی؛ (ب) درون‌یابی خطی؛ (پ) درون‌یابی اسپیلاین؛ (ج) جانپی منفرد روش آخرین مشاهده رو به جلو

همان‌طور که در شکل ۳ نیز دیده می‌شود، الگوی حاصل از روش رگرسیون درختی بیش از روش‌های دیگر با الگوی موجود در سری زمانی تطابق دارد.

### نتیجه‌گیری

وجود داده‌های گم‌شده در سری زمانی غلظت آلاینده‌ها بر عملکرد تجزیه و تحلیل داده‌ها در الگوریتم‌های یادگیری ماشین تأثیر منفی می‌گذارد و موجب بایاس می‌شود. روش‌های متعددی برای جانهی داده‌های گم‌شده ارائه شده است که انتخاب بهترین روش نیازمند مطالعه است. نتایج این پژوهش نشان داده است که تعیین شباهت مکانی- زمانی ایستگاه‌ها و استفاده از الگوی ایستگاه‌های مشابه با استفاده از الگوریتم پیچش زمانی پویا در ترکیب با روش‌های رگرسیون مبنای، موجب بهبود عملکرد مدل با فواصل زیاد گم‌شدگی می‌شود و مدل رگرسیون درختی مناسب‌ترین روش برای جانهی چندگانه به‌شمار می‌رود. روش‌های جانهی منفرد، با وجود سرعت و سادگی، وابسته به طول بازه گم‌شدگی در زمان هستند و عملکرد آنها به متغیر بررسی شده بستگی دارد. بنابراین استفاده از روش‌های منفرد در داده‌های آلودگی هوا که فواصل گم‌شدگی زیاد دارند توصیه نمی‌شود. شایان ذکر است که در هر منطقه مطالعاتی که با استفاده از الگوریتم پیچش زمانی پویا بتوان شباهت بین ایستگاه‌های مختلف را تشخیص داد، مدل قابلیت استفاده برای جانهی مقادیر از دست‌رفته غلظت آلاینده‌ها در آن منطقه مطالعاتی را خواهد داشت. با توجه به تأثیر دیگر عوامل مانند پارامترهای هواشناسی بر آلودگی هوا، در پژوهش‌های آینده می‌توان با افزودن این پارامترها دقت مدل را افزایش داد.

## منابع

- [۱]. بازگیر، سعید؛ قدیری معصوم، مجتبی؛ شمسی‌پور، علی‌اکبر؛ و سیدی سرنجیانه، شیوا (۲۰۱۵). «تحلیل رابطه آلودگی هوای تهران با ترافیک و شرایط جو برای کاهش مخاطرات، مدیریت مخاطرات محیطی، دوره ۲، شماره ۱. ص ۳۵-۴۹.
- [۲]. باقی یزدل، رقیه؛ جمالی، احسان؛ خدایی، ابراهیم؛ و حبیبی مجتبی (۱۳۹۵). «روش‌های برخورد با داده‌های گمشده: مزایا، معایب، رویکردهای نظری و معرفی نرم‌افزارها». نامه آموزش عالی، دوره ۹، شماره ۳۳، ص ۳۷-۱۱.
- [۳]. عثمانی، فرشته؛ و راسخی، علی‌اکبر (۱۳۹۷). «روش‌های وزن‌دهی احتمال معکوس و جانپی چندگانه برای تحلیل پاسخ در حالت گم‌شدگی». علوم آماری، دوره ۱۲ شماره ۲، ص ۴۸۳-۴۶۹.
- [۴]. کرمانی، آذر؛ اکبری، مه‌ری؛ علیجانی، بهلول؛ و مفاخری، امید (۲۰۱۵). «تحلیل آماری-همدیدی غلظت آلاینده مونواکسیدکربن براساس سمت و سرعت باد و مخاطره آن در شهر تهران». مدیریت مخاطرات محیطی، دوره ۲ شماره ۴، ص ۴۵۰-۴۳۹.
- [5]. Burgette, L.F.; & Reiter, J.P. (2010). "Multiple imputation for missing data via sequential regression trees", *American journal of epidemiology*, 172(9), pp: 1070-1076. doi: <https://doi.org/10.1093/aje/kwq260>.
- [6]. Caillault, É.P.; Lefebvre, A.; & Bigand, A. (2017). "Dynamic time warping-based imputation for univariate time series data", *Pattern Recognition Letters*. doi:<https://doi.org/10.1016/j.patrec.2017.08.019>.
- [7]. Chen, X.; & Xiao, Y. (2018). "A novel method for air quality data imputation by nuclear norm minimization", *Journal of Sensors*. doi:<https://doi.org/10.1155/2018/7465026>.
- [8]. Erler, N.S.; Rizopoulos, D.; Jaddoe, V.W.; Franco, O.H.; & Lesaffre, E.M. (2019). "Bayesian imputation of time-varying covariates in linear mixed models", *Statistical methods in medical research*, 28(2), pp: 555-568. doi:<https://doi.org/10.1177/0962280217730851>.
- [9]. Fortuin, V.; Rättsch, G.; & Mandt, S. (2019). "Multivariate time series imputation with variational autoencoders", *arXiv preprint arXiv:1907.04155*. doi: <https://arxiv.org/abs/1907.04155>..
- [10]. Ghazali, S.M.; Shaadan, N.; & Idrus, Z. (2020). "Missing data exploration in air quality data set using R-package data visualisation tools", *Bulletin of Electrical Engineering and Informatics*, 9(2), pp: 755-763. doi:<https://doi.org/10.11591/eei.v9i2.2088>.
- [11]. Gómez-Carracedo, M.; Andrade, J.; López-Mahía, P.; Muniategui, S.; & Prada, D. (2014). "A practical comparison of single and multiple imputation methods to handle complex missing data in air quality datasets", *Chemometrics and Intelligent Laboratory Systems*, 134, pp: 23-33. doi:<https://doi.org/10.1016/j.chemolab.2014.02.007>.

- [12]. Hadeed, S.J.; O'Rourke, M.K.; Burgess, J.L.; Harris, R.B.; & Canales, R.A. (2020). "Imputation methods for addressing missing data in short-term monitoring of air pollutants", *Science of The Total Environment*, pp: 139140. doi:<https://doi.org/10.1016/j.scitotenv.2020.139140>.
- [13]. Junger, W.; & De Leon, A.P. (2015). "Imputation of missing data in time series for air pollutants", *Atmospheric Environment*, 102, pp: 96-104. doi:<https://doi.org/10.1016/j.atmosenv.2014.11.049>.
- [14]. Junninen, H.; Niska, H.; Tuppurainen, K.; Ruuskanen, J.; & Kolehmainen, M. (2004). "Methods for imputation of missing values in air quality data sets", *Atmospheric Environment*, 38(18), pp: 2895-2907. doi:<https://doi.org/10.1016/j.atmosenv.2004.02.026>.
- [15]. Lin, J.; Li, N.; Alam, M.A.; & Ma, Y. (2020). "Data-driven missing data imputation in cluster monitoring system based on deep neural network", *Applied Intelligence*, 50(3), pp: 860-877. doi:<https://doi.org/10.1007/s10489-019-01560-y>.
- [16]. Liu, X.; Wang, X.; Zou, L.; Xia, J.; & Pang, W. (2020). "Spatial imputation for air pollutants data sets via low rank matrix completion algorithm", *Environment International*, 139, pp: 105713. doi:<https://doi.org/10.1016/j.envint.2020.105713>.
- [17]. Ma, J.; Cheng, J.C.; Jiang, F.; Chen, W.; Wang, M.; & Zhai, C. (2020). "A bi-directional missing data imputation scheme based on LSTM and transfer learning for building energy data", *Energy and Buildings*, pp: 109941. doi:<https://doi.org/10.1016/j.enbuild.2020.109941>.
- [18]. Mishra, S.; Dwivedi, V.; Sarvanan, C.; & Pathak, K. (2013). "Pattern discovery in hydrological time series data mining during the monsoon period of the high flood years in Brahmaputra River basin", *International Journal of Computer Applications*, 67(6).
- [19]. Raghunathan, T.E.; Lepkowski, J.M.; Van Hoewyk, J.; & Solenberger, P. (2001). "A multivariate technique for multiply imputing missing values using a sequence of regression models", *Survey methodology*, 27(1), pp: 85-96.
- [20]. Rombach, I.; Gray, A.M.; Jenkinson, C.; Murray, D.W.; & Rivero-Arias, O. (2018). "Multiple imputation for patient reported outcome measures in randomised controlled trials: advantages and disadvantages of imputing at the item, subscale or composite score level", *BMC medical research methodology*, 18(1), pp: 87. doi:<https://doi.org/10.1186/s12874-018-0542-6>.
- [21]. Shahbazi, H.; Karimi, S.; Hosseini, V.; Yazgi, D.; & Torbatian, S. (2018). "A novel regression imputation framework for Tehran air pollution monitoring network using outputs from WRF and CAMx models", *Atmospheric Environment*, 187, pp: 24-33. doi:<https://doi.org/10.1016/j.atmosenv.2018.05.055>.
- [22]. Stead, A.D.; & Wheat, P. (2020). "The case for the use of multiple imputation missing data methods in stochastic frontier analysis with illustration using English local highway data", *European Journal of Operational Research*, 280(1), pp: 59-77. doi:<https://doi.org/10.1016/j.ejor.2019.06.042>.
- [23]. Zeileis, A.; Grothendieck, G.; Ryan, J.A.; Andrews, F.; & Zeileis, M.A. (2019). "Package "zoo"".