

استفاده از خوشه‌بندی و رویکردی ترکیبی برای پر کردن مقادیر جاافتاده عددی

امیرمسعود سفیدیان و نگین دانشپور

که در آن M تعداد صفات هر نمونه است.
(۲) جاافتاده به طور تصادفی $(MAR)^2$:

در این حالت احتمال این که برای هر نمونه مقدار یک صفت مثل A_j معلوم یا جاافتاده باشد ممکن است به دیگر صفات وابسته باشد، ولی به مقدار خود آن صفت وابسته نیست. بنابراین در این حالت مقدار جاافتاده از دیگر صفات قابل محاسبه است یعنی

$$P(A_j = \text{NULL} | A_1, \dots, A_M) = P(A_j = \text{NULL} | A_1, \dots, A_{j-1}, A_{j+1}, \dots, A_M) \quad (2)$$

که در آن M تعداد صفات هر نمونه است.
(۳) جاافتاده به طور غیر تصادفی $(MNAR)^3$:

در این حالت احتمال این که برای هر نمونه مقدار یک صفت مثل A_j معلوم یا جاافتاده باشد ممکن است به مقدار خود آن صفت وابسته باشد. بنابراین در این حالت یک صفت جاافتاده به طور مستقیم از روی صفات دیگر قابل محاسبه نیست. در عمل، این که داده جاافتاده در کدام دسته‌بندی ذکر شده قرار می‌گیرد معمولاً غیر ممکن است زیرا نحوه ایجاد شدن مقدار جاافتاده نامعلوم است [۳].

رویکردهای ساده متفاوتی برای رسیدگی به مقادیر جاافتاده وجود دارد:

- حذف رکوردهایی با مقادیر جاافتاده
- پر کردن دستی مقادیر جاافتاده
- پر کردن با مقدار ثابت
- پر کردن با مقدار میانگین (یا مد)

در رویکرد اول اگر تعداد مقادیر جاافتاده بالا باشد بسیاری از داده‌های موجود را بیهوده هدر می‌دهد. رویکرد دوم در عمل نشدنی و یا بسیار زمان‌بر است. در رویکرد سوم فرض می‌شود که همه مقادیر جاافتاده مقدار یکسانی دارند که اصولاً فرض غلطی است. روش چهارم استفاده از مقدار میانگین برای صفات عددی و مقدار مد برای داده‌های نمادین است. مقدار میانگین بسیار حساس به داده‌های پرت است. مشکل دیگر این روش‌ها این است که در این روش‌ها از روابط احتمالی موجود بین صفات استفاده نمی‌شود. بنابراین این رویکردهای ساده نمی‌توانند دقت کافی برای پر کردن مقادیر جاافتاده را فراهم کنند و بنابراین روش‌های پیچیده‌تری برای رفع مقادیر جاافتاده ارائه شده است. توجه به ارتباطات موجود احتمالی در زیرمجموعه‌های یک مجموعه داده و بهره‌گیری از روشی مناسب در صورت عدم وجود ارتباطات قوی به صورت ترکیبی برای تخمین مقادیر جاافتاده می‌تواند پیشنهاد مناسبی برای استفاده از اطلاعات داده‌های موجود برای تخمین مقادیر جاافتاده باشد.

در این پژوهش رویکردی جدید مبتنی بر خوشه‌بندی داده‌ها و استفاده

چکیده: تخمین مقادیر جاافتاده یک گام مهم در پیش‌پردازش داده‌ها است. در این مقاله یک رویکرد دومرحله‌ای برای پر کردن مقادیر جاافتاده عددی ارائه شده است. در مرحله اول داده‌ها خوشه‌بندی می‌شوند و در مرحله دوم داده‌های جاافتاده درون هر خوشه با استفاده از یک روش ترکیبی از k نزدیک‌ترین همسایه وزن‌دار و رگرسیون خطی تخمین زده می‌شوند. از معیار همبستگی بین صفات در هر خوشه برای تعیین روش پر کردن داده‌های جاافتاده استفاده می‌شود. کیفیت پر کردن مقادیر جاافتاده با استفاده از معیار میانگین مربعات خطا سنجیده می‌شود. تأثیر پارامترهای مختلف بر میزان خطای داده‌های تخمین زده شده بررسی می‌گردد. عملکرد روش ارائه‌شده برای تخمین داده‌های جاافتاده بر روی پنج مجموعه داده نیز بررسی می‌شود. در نهایت عملکرد روش ارائه‌شده با چهار روش پر کردن با مقدار میانگین، روش تخمین با شبکه عصبی پرسپترون چندلایه (MLP)، روش پر کردن با خوشه‌بندی c-means فازی و روش k خوشه و نزدیک‌ترین همسایه مبتنی بر دسته (CKNNI) مقایسه می‌شود. نتایج به دست آمده نشان داده که خطای تخمین مقادیر جاافتاده در روش ارائه‌شده کمتر از خطا در دیگر روش‌های مقایسه‌شده است.

کلیدواژه: رگرسیون، مقادیر جاافتاده، نزدیک‌ترین همسایگان، همبستگی.

۱- مقدمه

مجموعه داده‌های دنیای واقعی معمولاً کامل نیستند. ممکن است مقدار یک یا چند صفت از صفات چندین نمونه به دلایل مختلف جاافتاده و نامعلوم باشند. وجود مقادیر جاافتاده اجتناب‌ناپذیر است. پیش‌پردازش یک گام مهم پیش از داده‌کاوی و یادگیری ماشین و یکی از مراحل پیش‌پردازش رسیدگی به مقادیر جاافتاده است. وجود مقادیر جاافتاده در مجموعه داده می‌تواند دقت الگوریتم‌های یادگیری مانند خوشه‌بندی و دسته‌بندی را به شدت کاهش دهد. مقدار جاافتاده می‌تواند به دلایل مختلفی مثل خطای انسانی، در دسترس نبودن مقادیر و یا خطای تجهیزات ثبت اطلاعات رخ دهند.

سه دسته‌بندی معمولاً برای داده‌های جاافتاده انجام می‌شود [۱] تا [۴]:

(۱) جاافتاده به طور کاملاً تصادفی $(MCAR)^1$:

در این حالت احتمال این که برای هر نمونه مقدار یک صفت مثل A_j معلوم یا جاافتاده باشد (مقدار تهی یا NULL داشته باشد) وابسته به هیچ یک از دیگر صفات نیست و کاملاً مستقل از دیگر صفات است یعنی

$$P(A_j = \text{NULL} | A_1, \dots, A_M) = P(A_j = \text{NULL}) \quad (1)$$

این مقاله در تاریخ ۱۵ آذر ماه ۱۳۹۵ دریافت و در تاریخ ۲ مرداد ماه ۱۳۹۶ بازنگری شد.

امیرمسعود سفیدیان، دانشکده مهندسی کامپیوتر، دانشگاه تربیت دبیر شهیدرجایی، تهران، (email: amirmasoud.sefidian@srttu.edu)

نگین دانشپور، دانشکده مهندسی کامپیوتر، دانشگاه تربیت دبیر شهیدرجایی، تهران، (email: ndaneshpour@srttu.edu)

1. Missing Completely at Random

2. Missing at Random

3. Missing Not at Random

۲- پیشینه تحقیق

در این بخش خلاصه‌ای از رویکردهای موجود برای پرکردن مقادیر جاافتاده بیان می‌شود. روش‌های مختلفی مثل روش‌های آماری، استفاده از یادگیری ماشین و داده‌کاوی برای پرکردن مقادیر جاافتاده ارائه شده است. دسته‌ای از روش‌ها با استفاده از قوانین انجمنی^۶ داده‌های جاافتاده را پر می‌کنند. محدودیت اصلی این روش‌ها این است که تنها قابل استفاده برای داده‌های غیر عددی هستند. در [۸] از یک رویکرد مبتنی بر قوانین انجمنی برای پرکردن مقادیر جاافتاده استفاده شده است. ابتدا مجموعه داده رابطه‌ای به مجموعه داده تراکنشی تبدیل می‌شود تا بتوان از آن قوانین انجمنی را استخراج نمود. پس از یافتن قوانین یافت شده از مجموعه داده تراکنشی برای هر مقدار جاافتاده، مجموعه قوانینی که در سمت راست آنها همان صفت جاافتاده حضور دارند، انتخاب می‌شوند. برای هر یک از قوانین انتخاب شده امتیازی محاسبه می‌گردد. قانونی با بیشترین امتیاز به عنوان قانون برگزیده انتخاب شده و از مقدار سمت راست آن قانون به عنوان مقدار جایگزین صفت مقدار صفت جاافتاده استفاده می‌شود. از نقاط قوت این روش می‌توان به دخیل کردن معیار وابستگی دو طرف قوانین انجمنی و طول قوانین و حمایت از مقادیر جاافتاده چندگانه اشاره کرد.

در [۹] قوانین یافت شده رتبه‌بندی می‌شوند و قانونی با بهترین رتبه به عنوان قانون انتخاب شده برای پرکردن مقدار جاافتاده استفاده خواهد شد. قوانین گزینش شده به ترتیب بر اساس طول قانون، درصد اعتماد^۷ و درصد پشتیبانی^۸ رتبه‌بندی می‌شوند. از نقاط قوت این روش می‌توان به دخیل کردن طول قوانین و رتبه‌بندی قوانین اشاره کرد. در نظر گرفتن تنها بهترین قانون هنگام پرکردن مقادیر از نقاط ضعف این روش می‌باشد.

از شبکه‌های عصبی نیز برای پرکردن مقادیر جاافتاده استفاده شده است. در [۵] با استفاده از شبکه عصبی پرسپترون چندلایه (MLP) داده‌های جاافتاده برای هجده مجموعه داده تخمین زده شده‌اند. همچنین تأثیر روش‌های مختلف یادگیری و پارامترهای مختلف شبکه عصبی بر روی دقت نهایی نیز بررسی شده‌اند. در [۱۰] با استفاده از طراحی و آموزش یک شبکه عصبی نگاشت خود سازمان‌دهنده (SOM)^۹ مقادیر جاافتاده تخمین زده شده‌اند. مشکل اصلی شبکه‌های عصبی زمان‌بر بودن فرایند یادگیری و ساخت مدل بهینه برای یک مجموعه داده نسبت به دیگر روش‌ها و همچنین قابلیت تفسیر^{۱۰} پایین است [۱۱]. به علاوه، اصولاً آموزش مناسب شبکه‌های عصبی نیاز به تعداد نمونه‌های ورودی بالایی دارد و بنابراین هنگامی که تعداد نمونه‌های جاافتاده بالا است، دقت تخمین این روش‌ها کاهش می‌یابد.

روش k نزدیک‌ترین همسایه یکی از معروف‌ترین و پایه‌ای‌ترین روش‌ها برای تخمین مقادیر جاافتاده است. سادگی و نتایج قابل قبول آن نسبت به دیگر روش‌ها از ویژگی‌های آن است. از دیگر مزایای این روش قابلیت استفاده برای هم صفات عددی و هم صفات غیر عددی با تغییر تابع سنجش فاصله است. همچنین در این روش زمانی صرف ساخت مدل برای تخمین داده‌های جاافتاده نمی‌شود [۷]. این ویژگی‌ها سبب شده که در این مقاله نیز از حالت وزن‌دار این روش به عنوان یک روش تخمین

از ترکیب روش k نزدیک‌ترین همسایه وزن‌دار^۱ (WKNN) و رگرسیون خطی^۲ برای پرکردن مقادیر جاافتاده عددی ارائه شده است. ویژگی اصلی روش ارائه شده استفاده ترکیبی و توأم از همبستگی موجود بین داده‌ها و فاصله بین نمونه‌ها برای پرکردن مقادیر جاافتاده است.

در روش‌های موجود برای پرکردن مقادیر جاافتاده اصولاً دو نوع رویکرد وجود دارد. رویکرد اول روش‌هایی هستند که فرض می‌کنند روابط قوی بین صفات در مجموعه داده وجود دارد و بنابراین تلاش می‌کنند تا روابط موجود را مدل‌سازی کنند. از جمله این روش‌ها می‌توان به روش‌های آماری مانند رگرسیون اشاره کرد. مشکل اصلی استفاده از این روش‌ها به تنهایی این است که ممکن است در داده‌های ورودی روابط قوی وجود نداشته باشد و بنابراین مدل ساخته شده دقت پایینی را ارائه خواهد داد. به علاوه در اکثر روش‌های ارائه شده مبتنی بر روابط قوی از معیار مشخص و دقیقی برای سنجش روابط بین صفات، مانند ضریب همبستگی^۳ استفاده نشده است.

دسته دیگری از روش‌ها مانند روش‌های مبتنی بر k نزدیک‌ترین همسایگان، کاملاً از روابط بین صفات چشم‌پوشی می‌کنند و فقط از میزان نزدیکی نمونه‌ها استفاده می‌کنند. در این روش‌ها اطلاعات بسیار مفیدی مانند قدرت استنتاج یک مقدار جاافتاده از طریق روابط موجود با دیگر صفات نادیده گرفته می‌شود. بنابراین در این مقاله ترکیبی از دو دسته روش ذکر شده به منظور بهره‌وری از مزایای آنها ارائه می‌شود. با توجه به مطالعات انجام شده، پیش از این ترکیبی از این دو دسته روش ارائه نشده است. در روش ارائه شده ابتدا سعی می‌شود تا با استفاده از خوشه‌بندی بتوان روابط قوی موجود بین نمونه‌ها را پیدا کرد. علت این کار این است که ممکن است روابط در زیرمجموعه‌ای از مجموعه داده ورودی قوی‌تر از کل مجموعه داده باشد. در کارهای پیشین از خوشه‌بندی برای گروه‌بندی نمونه‌های نزدیک به هم استفاده شده است. در این مقاله از خوشه‌بندی برای یافتن نمونه‌هایی با همبستگی بالا استفاده شده است. پس از خوشه‌بندی، میزان ارتباط بین صفات با معیاری دقیق و کمی به نام ضریب همبستگی درون هر خوشه اندازه‌گیری می‌شود. لازم به ذکر است که این معیار پیش از این برای این منظور استفاده نشده است. در نهایت بسته به این که آیا روابط به حد کافی قوی هستند یا خیر، از یکی از دو دسته روش بالا استفاده می‌شود.

نتایج بررسی شده بر روی پنج مجموعه داده نشان داده است که رویکرد پیشنهادی، دقت بیشتری در مقایسه با چهار روش پرکردن با مقدار میانگین، روش تخمین با شبکه عصبی پرسپترون چندلایه^۴ (MLP) [۵]، روش تخمین با خوشه‌بندی c -means فازی [۶] و روش k خوشه و نزدیک‌ترین همسایه مبتنی بر دسته^۵ (CKNN) [۷] دارد.

بخش‌های بعدی این مقاله بدین شرح است. در بخش دوم مروری بر روش‌های گذشته بیان می‌شود. در بخش سوم جزئیات رویکرد پیشنهادی ارائه می‌گردد و پیچیدگی زمانی رویکرد پیشنهادی تحلیل می‌شود. در بخش چهارم نتایج و آزمایش‌ها مطرح شده است و در نهایت در بخش پنجم نتیجه‌گیری بیان می‌شود.

6. Association Rules
7. Confidence
8. Support
9. Self Organizing Map
10. Interpretability

1. Weighted K-Nearest Neighbor
2. Linear Regression
3. Correlation Coefficient
4. Multi-Layer Perceptron
5. Class-Based K-Clusters Nearest Neighbor Imputation

استفاده شود.

پرکردن مقادیر جاافتاده استفاده شده است. در [۱۹] یک رویکرد دومرحله‌ای برای پرکردن مقادیر جاافتاده ارائه شده است. این رویکرد شامل یک تخمین اولیه از مقادیر جاافتاده توسط نزدیک‌ترین مرکز خوشه به روش خوشه‌بندی k-means و یک گام پالایش و بهبود مقدار تخمین زده شده از گام قبلی، توسط شبکه MLP می‌باشد. در [۲۰] در گام اول به جای استفاده از خوشه‌بندی k-means از خوشه‌بندی c-means فازی استفاده شده که در آن هر نمونه با یک احتمال متعلق به هر خوشه است. این روش نتایج بهتری را به همراه داشته است. از نقاط قوت این روش‌ها می‌توان به بهبود k-means برای پرکردن مقادیر جاافتاده اشاره کرد. نیاز به آموزش یک شبکه MLP برای هر مقدار جاافتاده و آموزش مجدد روی کل مجموعه آموزشی، تعیین ساختار (به هم بندی و وزن‌ها) بهینه شبکه MLP برای هر بار آموزش برای پرکردن هر مقدار جاافتاده از نقاط ضعف این روش می‌باشد.

در [۲۱] یک رویکرد ترکیبی از قوانین انجمنی و k نزدیک‌ترین همسایه برای پرکردن مقادیر جاافتاده ارائه شده است. از نقاط قوت این روش می‌توان به بهره‌گیری از قدرت ترکیب دو روش متفاوت اشاره کرد. در واقع در این روش سعی در استفاده از ارتباطات و فواصل موجود برای پرکردن مقادیر جاافتاده شده و این روش قابلیت استفاده برای داده‌های عددی را به طور مستقیم ندارد (به دلیل استفاده از قوانین انجمنی).

در [۲۲] روشی مبتنی بر الگوریتم EM^Y و درخت تصمیم به نام DMI^A برای پرکردن مقادیر جاافتاده طراحی شده است. ایده اصلی این روش این است که وابستگی بین افزایشی مجموعه داده قوی‌تر از وابستگی بر کل مجموعه داده است و بنابراین الگوریتم EM بر روی بخش‌های متفاوت مجموعه داده بهتر عمل خواهد کرد و در نهایت نتیجه به دست آمده بهتر خواهد شد. افزایشی مجموعه داده توسط درخت تصمیم صورت می‌پذیرد. از نقاط قوت این روش توجه به این است که وابستگی بین صفات لزوماً در کل مجموعه داده برقرار نیست. مناسب بودن برای مجموعه داده‌هایی که وابستگی بین صفات آنها وجود دارد از نقاط ضعف این روش است.

رگرسیون نیز به عنوان یک روش پایه برای تخمین مقادیر جاافتاده می‌تواند مورد استفاده قرار گیرد. زمانی که رابطه بین دو متغیر قوی‌تر باشد، دقت رگرسیون نیز بیشتر خواهد بود. با استفاده از محاسبه معادله رگرسیون خطی می‌توان داده‌های جاافتاده را از طریق داده‌های موجود تخمین زد. در [۳] از یک رویکرد افزایشی و با استفاده از تولید دنباله‌ای از رگرسیونها تخمین داده‌های جاافتاده انجام می‌شود. مشکل اصلی روش‌های مبتنی بر رگرسیون، ضعف عملکرد آنها در صورت عدم وجود روابط قوی بین صفات پیش‌بینی کننده و صفت پیش‌بینی شده است.

در [۲۳] دو راهبرد برای بهبود دقت پرکردن مقادیر جاافتاده ارائه شده است. ایده اصلی در هر دو راهبرد، انجام عمل انتخاب نمونه^۹ است. در این این عمل سعی می‌شود تا نمونه‌های شامل داده‌های پرت و داده‌های مغشوش^{۱۰} از مجموعه داده حذف شوند تا کیفیت داده‌ها پیش از عملیات پرکردن مقادیر جاافتاده افزایش یابد. در واقع مجموعه داده از یک فیلتر حذف داده‌های پرت و مغشوش عبور داده می‌شود. در هر دو راهبرد

در [۱۲] از یک رویکرد مبتنی بر روش k نزدیک‌ترین همسایه برای پرکردن مقادیر جاافتاده استفاده شده است. در این رویکرد ابتدا درجه جافتادگی هر نمونه که برابر است با تعداد صفات جافتاده در هر نمونه، محاسبه می‌گردد. نمونه‌هایی که درجه جافتادگی آنها بزرگ‌تر از دو است از ادامه پردازش خارج می‌شوند و سپس مقادیر جافتاده با مقادیر نزدیک‌ترین نمونه‌های کامل پر می‌شوند. از نقاط قوت این روش می‌توان به سادگی رویکرد ارائه‌شده و قدرت تعمیم بالا برای کاربردهای مختلف با تغییر معیار شباهت اشاره کرد. عدم توانایی در پرکردن مقادیر جافتاده چندگانه، حساسیت به داده‌های پرت و حساسیت بالا به تابع شباهت از نقاط ضعف این روش می‌باشد. همچنین لزوماً بررسی فقط نزدیک‌ترین نمونه، همواره بهترین گزینه نیست.

در [۱۳] از یک رویکرد پرکردن k نزدیک‌ترین همسایه تکراری مبتنی بر گری^۱ به نام GBKII^۲ استفاده شده است. رویکرد مطرح‌شده در یک فرایند تکراری بر روی مقادیر جافتاده، تک به تک مقادیر جافتاده را پر می‌نماید. از مزیت‌های این روش امکان پرکردن مقادیر جافتاده چندگانه است. به‌علاوه چون سنجش شباهت فقط بین نمونه‌های با برچسب یکسان صورت می‌گیرد، پیچیدگی زمانی این روش پایین است.

در [۷] از یک نسخه بهبودیافته روش k نزدیک‌ترین همسایه و استفاده از خوشه‌بندی داده‌های قرارگرفته در یک دسته (کلاس) به نام CKNNI برای تخمین مقادیر جافتاده استفاده شده است. مشکل اصلی این روش نیاز به دانستن برچسب تمام نمونه‌ها است. در واقع اگر کلاس نمونه‌ها مشخص نباشد نمی‌توان از روش ذکرشده استفاده نمود.

در [۱۴] دقت رویکرد نزدیک‌ترین همسایه بر روی چهار مجموعه داده با دو معیار فاصله مَنهتن^۳ و اقلیدسی برای پرکردن مقادیر جافتاده بررسی و نتایج با روش پرکردن با میانگین مقایسه شده است. رویکرد نزدیک‌ترین همسایه دقت بیشتری در تخمین مقادیر جافتاده در برابر روش میانگین داشته است.

در [۱۵] به کمک جمع‌سپاری^۴ مقادیر جافتاده پر می‌شوند. مشکل اصلی این روش به طور کامل خودکار نبودن و تأثیرگذاری عامل انسانی در دقت نهایی داده‌های پرشده است.

در [۱۶] ابتدا داده‌ها خوشه‌بندی شده و سپس هر مقدار جافتاده با میانگین دو مقدار جایگزین می‌شود. مقدار اول فاصله نمونه جافتاده تا مرکز خوشه‌ای است که نمونه جافتاده به آن تعلق دارد. مقدار دوم برابر با مقدار صفتی از مرکز خوشه است که متناظراً در نمونه جافتاده، نامعلوم است. چشم‌پوشی از روابط احتمالی بین نمونه‌های موجود در یک مجموعه داده مشکل اصلی این روش است.

در [۱۷] با استفاده از یک خوشه‌بندی به نام خوشه‌بندی همسایگان مشترک برای هدف پرکردن داده‌های جافتاده استفاده شده است. از نقاط قوت این روش می‌توان به استفاده از معیار شباهت برای صفات ترکیبی و دخیل کردن معیار جدیدی برای سنجش شباهت نمونه‌ها اشاره کرد. نحوه توزیع داده‌ها بر تعداد همسایگان مشترک در این روش تأثیرگذار است و ممکن است ماتریس همسایگان مشترک تنگ^۵ شود.

در [۱۸] از خوشه‌بندی‌های k-means و c-means فازی برای

6. Topology
7. Expectation Maximization
8. Decision Tree Based Missing Value Imputation Technique
9. Instance Selection
10. Noisy Data

1. Grey
2. Grey-Based K-NN Iteration Imputation
3. Manhattan
4. Crowdsourcing
5. Sparse

جدول ۱: بررسی کلی مشکلات روش‌های موجود.

مشکل اصلی	رویکرد پرکردن
کاهش چشم‌گیر دقت تخمین هنگامی که تعداد نمونه‌های جافتاده بالا است	روش‌های مبتنی بر شبکه عصبی
غیر قابل اعمال مستقیم برای تخمین مقادیر جافتاده عددی	روش‌های مبتنی بر قوانین انجمنی
عملکرد ضعیف وقتی ارتباطات قوی بین صفات وجود ندارد	روش‌های مبتنی بر رگرسیون و الگوریتم EM
عدم توجه و استفاده از روابط موجود بین صفات	روش‌های مبتنی بر k نزدیک‌ترین همسایگان

۳- رویکرد پیشنهادی

در این بخش بیانی دقیق و رسمی از مراحل رویکرد پیشنهادی ارائه و سپس پیچیدگی رویکرد پیشنهادی تحلیل می‌شود.

۳-۱ مراحل رویکرد پیشنهادی

روش ارائه شده برای پرکردن داده‌های جافتاده شامل دو گام اصلی است. در گام اول داده‌های کامل (بدون مقادیر جافتاده) خوشه‌بندی می‌شوند. هدف اصلی از خوشه‌بندی این است که ممکن است همبستگی قوی بر روی کل مجموعه داده برقرار نباشد در حالی که بین زیرمجموعه‌ای از نمونه‌ها، ارتباطات قوی وجود داشته باشد. بنابراین سعی می‌شود با استفاده از خوشه‌بندی، همبستگی‌های احتمالی موجود در زیرمجموعه‌های مجموعه داده اصلی را پیدا نمود. در گام بعدی مقادیر جافتاده در هر خوشه با استفاده از روشی مرکب از k نزدیک‌ترین همسایه وزن‌دار و رگرسیون خطی پر خواهند شد. از یک حد آستانه $(0 \leq t \leq 1)$ بر روی مقادیر ماتریس همبستگی بین صفات برای تعیین روش تخمین مقدار جافتاده استفاده می‌شود. اگر صفتی با همبستگی کافی به مقدار صفت جافتاده پیدا شود از روش رگرسیون خطی و در غیر این صورت از روش k نزدیک‌ترین همسایه وزن‌دار استفاده می‌گردد. در واقع ایده اصلی این است که در صورت وجود همبستگی بین داده‌ها از این همبستگی تا حد ممکن استفاده شود و در غیر این صورت از نزدیک‌ترین نمونه‌ها برای پرکردن مقادیر جافتاده استفاده شود. خوشه‌بندی داده‌ها در این مقاله با روش خوشه‌بندی k -means انجام می‌شود. شکل ۱ روند کلی رویکرد پیشنهادی را نشان می‌دهد. شبه‌کد الگوریتم پیشنهادی در شکل ۲ نمایش داده شده است. در ادامه جزئیات مراحل مختلف رویکرد پیشنهادی شرح داده می‌شود.

فرض می‌شود مجموعه داده رابطه‌ای DS شامل N رکورد بر روی شمای (A_1, A_2, \dots, A_M) تعریف شده و هر نمونه شامل M صفت است. مقدار صفت i ام نمونه i ام را نمایش می‌دهد و تابع $d(x, y)$ بیانگر فاصله دو نمونه x و y است.

در ابتدا مجموعه داده اصلی (DS) به دو مجموعه داده داده‌های کامل $(DS_C = \{c_1, c_2, \dots, c_p\})$ (رکوردهایی بدون مقادیر جافتاده) و مجموعه داده‌های با مقادیر جافتاده $(DS_M = \{m_1, m_2, \dots, m_o\})$ تقسیم می‌شود (خط ۱ از شبه‌کد). سپس مجموعه DS_C با استفاده از خوشه‌بندی به c تا خوشه S_1, S_2, \dots, S_c تقسیم می‌شود و بنابراین c مرکز خوشه به دست می‌آید (خط ۲). هر یک از عناصر DS_C با توجه به میزان نزدیکی به مراکز خوشه به یک خوشه منتسب می‌شوند (خط ۳). هر یک از عناصر DS_M نیز با توجه به میزان نزدیکی به مراکز خوشه به یک خوشه منتسب می‌شود (خط ۴). قابل ذکر است که فاصله هر نمونه جافتاده بدون در نظر گرفتن صفات جافتاده محاسبه می‌گردد.

عملیات انتخاب نمونه از طریق الگوریتم $^{1}DROP$ [۲۴] انجام شده است. تفاوت دو راهبرد پیشنهادی در ترتیب انجام عملیات انتخاب نمونه و پرکردن داده‌ها است. انجام عملیات انتخاب نمونه منجر به بهبود دقت حاصل از دسته‌بندی داده‌های تخمین زده شده است.

در [۲۵] با استفاده از ترکیب نظریه مجموعه‌های رافِ فازی 2 و k نزدیک‌ترین همسایگان، الگوریتمی به نام پرکردن نزدیک‌ترین همسایگان رافِ فازی 3 (FRNNI) ارائه شده است. در این الگوریتم از تقریب‌های پایین و بالا 4 در نزدیک‌ترین نمونه‌ها، نسبت به نمونه جافتاده برای پرکردن مقادیر نامعلوم استفاده می‌شود. علت اصلی استفاده از نظریه مجموعه‌های رافِ فازی در این روش مناسب بودن این رویکردها در تعامل با محیط‌هایی با حضور عدم قطعیت 5 است و بنابراین انتخاب خوبی برای مسئله پرکردن مقادیر جافتاده است.

با مطالعات انجام‌گرفته نیاز به روشی برای پرکردن مقادیر جافتاده با توجه به همبستگی احتمالی موجود بین صفات و همچنین استفاده از روش جایگزینی مناسب برای حالتی که همبستگی کافی بین صفات وجود ندارد ضروری می‌رسد. بنابراین رویکردی در این مقاله پیشنهاد می‌شود که ابتدا با استفاده از رگرسیون خطی، سعی در استفاده از روابط موجود بین داده‌ها برای پرکردن مقادیر جافتاده می‌نماید. تلاش می‌شود که روابط در زیرمجموعه‌هایی از مجموعه داده اصلی یافت شوند. دلیل این امر این است که شاید روابط قوی بر روی کل مجموعه داده یافت نشود در حالی که بین زیرمجموعه‌ای از نمونه‌ها، ارتباطات قوی‌تری وجود داشته باشد. در صورت یافت‌نشدن میزان ارتباط کافی درون زیرمجموعه داده‌ها از روش پایه نزدیک‌ترین همسایگان وزن‌دار استفاده می‌شود. مرجع [۲۲] اگرچه با استفاده از درخت تصمیم سعی در افزایش همبستگی نمونه‌ها دارد با این حال از معیاری دقیق برای سنجش صحت فرضیه افزایش همبستگی درون هر افزاز از نمونه‌ها استفاده نمی‌کند. روش‌های مبتنی بر رگرسیون به تنهایی این نکته که ممکن است ارتباط بین صفات در یک مجموعه داده قوی نباشند را در نظر نمی‌گیرند و بنابراین دقت عملکرد آنها در این شرایط به شدت کاهش می‌یابد. در روش‌های مبتنی بر k نزدیک‌ترین همسایه، اطلاعات راجع به همبستگی بین صفات به طور کلی نادیده گرفته می‌شوند. در روش ترکیبی ارائه‌شده سعی گردیده تا تمام این موارد به طور توأم برای بهبود کیفیت تخمین در نظر گرفته شوند. جدول ۱ مشکلات اصلی در روش‌های موجود را به اختصار نمایش می‌دهد. روش پیشنهادی سعی دارد تا به طور هم‌زمان مشکلات دو دسته از روش‌های مبتنی بر رگرسیون و روش‌های مبتنی بر k نزدیک‌ترین همسایگان را برطرف سازد.

1. Incremental Reduction Optimization Procedure
2. Fuzzy Rough Sets Theory
3. Fuzzy-Rough Nearest Neighbor Imputation
4. Lower and Upper Approximations
5. Uncertainty

Input:

DS : A dataset with N instances defined on schema (A_1, A_2, \dots, A_M)

with missing values

c : number of clusters

p : comparison percentage

t : threshold for method selection

Output:

DS : Imputed dataset

Imputation-Method (c, p, t) :

- Split DS to $DS_C = \{c_1, c_2, \dots, c_p\}$ (contains complete instances) and $DS_M = \{m_1, m_2, \dots, m_Q\}$ (contains instances with missing value) ($P + Q = N$)
- Cluster DS_C to c clusters using k-means algorithm (S_1, S_2, \dots, S_c)
- Assign each instance in DS_C to a cluster
- Assign each instance in DS_M to a cluster
- $COR(S_1), COR(S_2), \dots, COR(S_c) \leftarrow$ Compute correlation matrix for each cluster
- foreach** cluster S_f
- foreach** $m_i^j \in S_f$
- $w^* \leftarrow$ find index of most correlated attribute to A_j using $COR(S_f)$
 $\rho \leftarrow COR(S_f)$
- if** $\rho_{j,w^*} \geq t$
- Create a linear regressor using A_j as dependent variable and A_{w^*} as independent variable
- $estimatedValue \leftarrow$ estimate m_i^j using created regressor
- end**
- else**
- $c_1, c_2, \dots, c_k \leftarrow$ Select $p\%(k)$ of nearest neighbors to m_i^j in S_f
- $estimatedValue \leftarrow$ estimate m_i^j using weighted k nearest neighbor algorithm
- end**
- $m_i^j \leftarrow estimatedValue$
- return** $DS_C \cup DS_M$

شکل ۲: شبه‌کد رویکرد پیشنهادی.

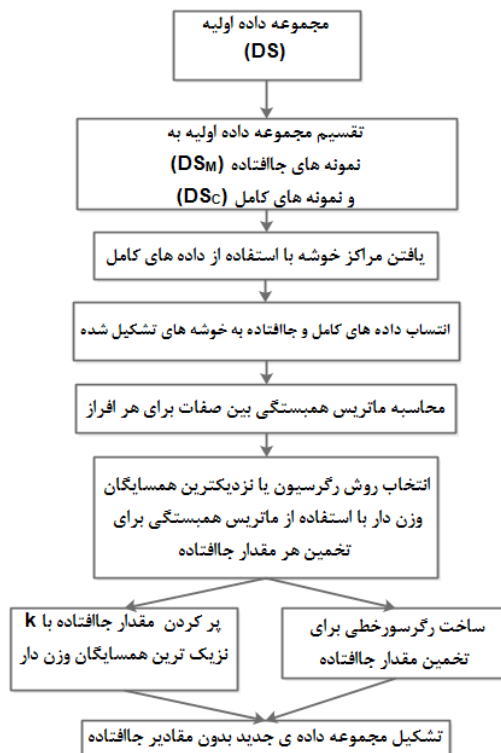
$$A_j = a.A_{w^*} + b \quad (۶)$$

که در آن a و b ضرایب رگرسیون خطی هستند. این ضرایب از طریق برازش^۴ بهترین خط بر روی داده‌ها به کمک روش حداقل مربعات^۵ محاسبه می‌شوند. در این حالت مقدار جاافتاده از (۷) محاسبه می‌گردد (خط ۱۱)

$$m_i^j = a.m_i^{w^*} + b \quad (۷)$$

اگر $\rho_{j,w^*} \leq t$ باشد آن گاه مقدار جاافتاده با استفاده از روش k نزدیک‌ترین همسایه وزن‌دار محاسبه می‌گردد. در این حالت، فاصله نمونه جاافتاده با تمام نمونه‌های کامل موجود در خوشه S_f محاسبه می‌شود و این فواصل به ترتیب صعودی مرتب می‌گردند (یعنی $d(c_1, m_i) < d(c_2, m_i) < \dots < d(c_{|S_f|}, m_i)$ که در آن $|S_f|$ اندازه

- Dependent Variable
- Fitting
- Least Square Method



شکل ۱: روند کلی رویکرد پیشنهادی.

ضریب همبستگی پیرسون^۱ برای جفت صفات (A_i, A_j) که $i \neq j$ است برای نمونه‌های هر خوشه محاسبه شده و نتایج در یک ماتریس همبستگی برای آن خوشه ذخیره می‌گردد (خط ۵). برای خوشه f ، ماتریس همبستگی به صورت (۳) خواهد بود

$$COR(S_f) = \begin{bmatrix} 1 & \rho_{1,2} & \dots & \rho_{1,M} \\ \rho_{2,1} & 1 & \dots & \rho_{2,M} \\ \vdots & \vdots & \ddots & \vdots \\ \rho_{M,1} & \rho_{M,2} & \dots & 1 \end{bmatrix} \quad (۳)$$

که در آن هر عنصر از (۴) محاسبه می‌گردد

$$\rho_{i,j} = \frac{|\text{cov}(A_i, A_j)|}{\sqrt{\sigma(A_i)\sigma(A_j)}} \quad , \quad 0 \leq \rho_{i,j} \leq 1 \quad (۴)$$

که $\text{cov}(A_i, A_j)$ کواریانس بین مقادیر صفت A_i و صفت A_j و $\sigma(A_i)$ و $\sigma(A_j)$ به ترتیب برابر با واریانس مقادیر صفت A_i و صفت A_j هستند. $|\cdot|$ بیانگر عملگر قدر مطلق است. چون میزان همبستگی بین صفات در این روش مهم است نه نوع همبستگی (همبستگی مثبت یا منفی) از قدر مطلق مقدار ضریب همبستگی در محاسبات استفاده می‌شود. حال برای هر مقدار جاافتاده مثل $m_i^j \in S_f$ اندیس مرتبط‌ترین صفت نسبت به صفت جاافتاده طبق (۵) از طریق ماتریس همبستگی افزاز مربوطه پیدا می‌شود (خط ۸)

$$w^* = \arg \max_w (\rho_{j,w}) \quad , \quad w^* \neq j \quad (۵)$$

اگر $\rho_{j,w^*} \geq t$ باشد که در آن t حد آستانه از پیش تعریف شده است، آن گاه یک رگرسیون خطی با صفت مستقل A_{w^*} ^۲ و صفت وابسته^۳ A_j به صورت (۶) ساخته می‌شود (خط ۱۰)

- Pearson Correlation Coefficient
- Independent Variable

جدول ۲: مشخصات مجموعه داده‌های استفاده شده در آزمایش‌ها.

مجموعه داده	تعداد نمونه	تعداد ویژگی
Iris	۱۵۰	۴
Wine	۱۷۸	۱۳
Glass	۲۱۴	۱۰
Haberman	۳۰۶	۳
Wholesale Customers	۴۴۰	۸

با در نظر گرفتن موارد گفته شده، پیچیدگی الگوریتم پیشنهادی در بدترین حالت از مرتبه $O(c.M^T.P \log P)$ خواهد بود. اگر فرض کنیم که در بدترین حالت $P \approx N$ ، آن گاه پیچیدگی رویکرد پیشنهادی برابر با $O(c.M^T.N \log N)$ است.

۴- آزمایش‌ها و نتایج

در این بخش بررسی‌های جامعی بر روی تأثیر مقادیر مختلف الگوریتم ارائه شده آمده و همچنین دقت رویکرد پیشنهادی با چهار روش دیگر مقایسه می‌شود. از پنج مجموعه داده موجود در پایگاه داده UCI [۲۶] برای آزمایش‌ها استفاده می‌شود. جدول ۲ مشخصات این مجموعه داده‌ها را نمایش می‌دهد که این مجموعه داده‌ها بدون مقادیر جافتاده می‌باشند و پیاده‌سازی‌های لازم با استفاده از زبان برنامه‌نویسی پایتون انجام شده است.

۴-۱ مشخصات آزمایش‌ها

در این زیربخش جزئیات آزمایش‌های انجام شده بیان می‌شود. برای بررسی تأثیر متغیرهای رویکرد پیشنهادی، ابتدا به تعداد ۱۰ درصد کل نمونه‌های هر مجموعه داده، مقادیر جافتاده به طور تصادفی در آن مجموعه داده ایجاد و سپس روش ارائه شده بر روی آن مجموعه داده اعمال می‌شود.

از معیار فاصله اقلیدسی برای سنجش فاصله دو نمونه استفاده شده است. مقدار این فاصله بین دو نمونه x_i و x_j از (۱۰) محاسبه می‌شود

$$d(x_i, x_j) = \sqrt{\sum_{k=1}^M (x_i^k - x_j^k)^2} \quad (10)$$

که در آن M تعداد صفات هر نمونه می‌باشد.

پیش‌پردازش نرمال‌سازی داده‌ها به بازه ۰ تا ۱ به روش Min-Max با استفاده از (۱۱) برای هر مقدار در مجموعه داده صورت گرفته است

$$x_i^j = \frac{x_i^j - x_i^{j,\min}}{x_i^{j,\max} - x_i^{j,\min}} \quad (11)$$

که در آن $x_i^{j,\max}$ و $x_i^{j,\min}$ به ترتیب مقادیر کمینه و بیشینه برای صفت j ام هستند. نرمال‌سازی از غلبه صفاتی با بازه مقادیر کوچک بر صفاتی با بازه مقادیر بزرگ در محاسبه فاصله جلوگیری می‌کند.

هر یک از آزمایش‌های انجام شده در ادامه، ۵ بار به صورت تصادفی انجام گردیده و میانگین خطای این ۵ بار به عنوان نتیجه نهایی، گزارش شده است. برای محاسبه خطا از معیار میانگین مربعات خطا^۲ مطابق (۱۲) استفاده می‌شود

$$RMSE = \sqrt{\frac{1}{n} \sum_{k=1}^n (y_k - y'_k)^2} \quad (12)$$

که در آن y_k مقدار واقعی یک مقدار جافتاده، y'_k مقدار تخمین زده شده برای آن مقدار جافتاده و n تعداد مقادیر جافتاده است.

از خوشه‌بندی k-means برای خوشه‌بندی داده‌ها استفاده شده است. اگرچه این روش یکی از قدیمی‌ترین روش‌های خوشه‌بندی است با این حال هنوز هم یکی از پرکاربردترین روش‌های خوشه‌بندی است [۲۷]. حداکثر تعداد ۱۰۰ تکرار برای روش خوشه‌بندی k-means تنظیم شده است. در ادامه تأثیر متغیرهای مختلف رویکرد ارائه شده یعنی تعداد

خوشه S_f است). p درصد نمونه‌های کامل در خوشه نمونه جافتاده به عنوان k نزدیک‌ترین نمونه کامل به نمونه جافتاده انتخاب می‌شوند $(\{c_1, c_2, \dots, c_k\} \in S_f)$ و در نهایت مقدار جافتاده از (λ) محاسبه می‌گردد (خطوط ۱۳ و ۱۴)

$$m_i^j = \frac{D_1 \times c_1^j + D_2 \times c_2^j + \dots + D_k \times c_k^j}{D_1 + D_2 + \dots + D_k} \quad (8)$$

که در آن D_e از (۹) به دست می‌آید

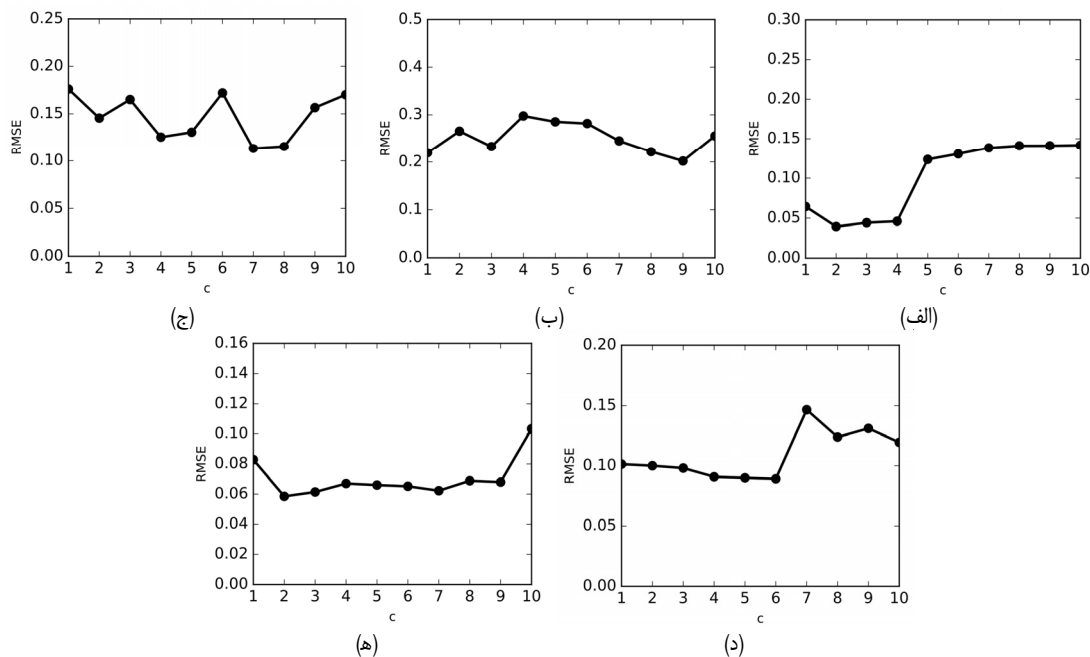
$$D_e = \frac{1}{d(c_e, m_i)} \quad , \quad e = 1, 2, \dots, k \quad (9)$$

در واقع (۸) رابطه میانگین وزن دار است که وزن‌ها در آن معکوس فاصله دو نمونه هستند. در نهایت مقدار محاسبه شده m_i^j به جای مقدار جافتاده جایگزین می‌شود و پس از پرکردن تمام مقادیر جافتاده مجموعه داده کامل بدون مقادیر جافتاده به دست می‌آید (خط ۱۶).

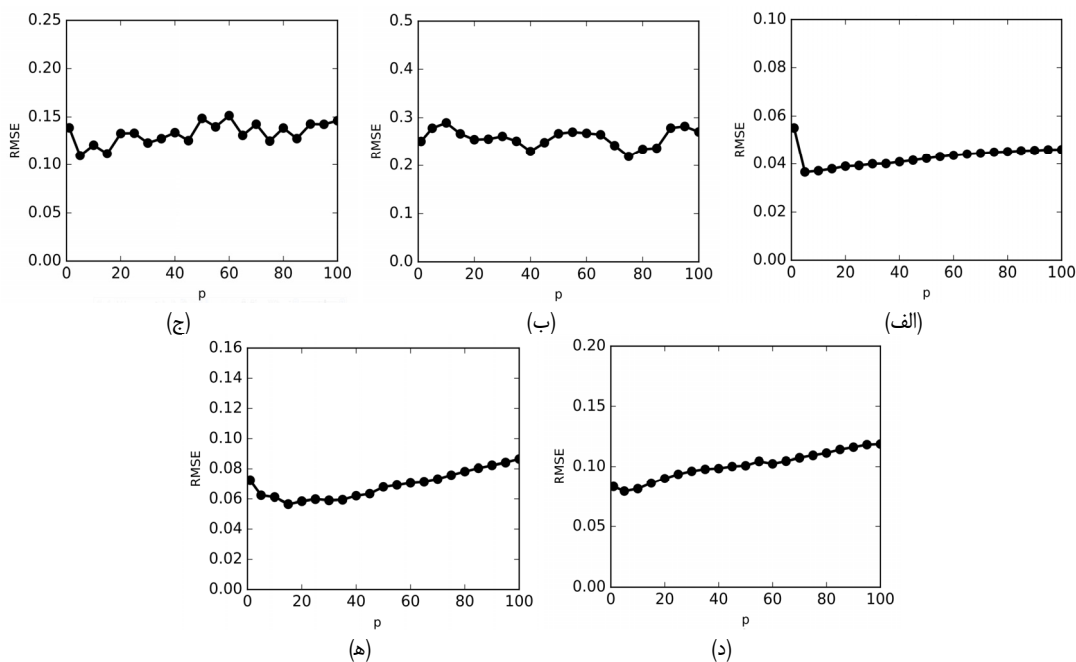
۳-۲ تحلیل پیچیدگی زمانی

عملیات پیش‌پردازش تنها یک بار و پیش از عملیات اصلی استخراج دانش از داده‌ها انجام می‌شود و بنابراین پیچیدگی زمانی عملیات پیش‌پردازش در اولویت پایین‌تری نسبت به کیفیت این امر قرار می‌گیرد. با این حال برای توصیف دقیق‌تر، پیچیدگی زمانی رویکرد پیشنهادی در این بخش تحلیل می‌شود.

فرض کنید مجموعه داده ورودی شامل N نمونه است. از این N نمونه P تای آنها کامل و Q تای آنها شامل مقادیر جافتاده است $(P+Q=N)$. همچنین فرض کنید هر نمونه شامل M ویژگی است. روش مطرح شده شامل سه گام اصلی است. در گام اول یک خوشه‌بندی k-means فقط بر روی داده‌های کامل انجام می‌شود (خطوط ۲ تا ۴). مرتبه زمانی این خوشه‌بندی برابر با $O(c.I.P.M)$ است که در آن c تعداد خوشه‌ها و I برابر با حداکثر تعداد تکرار در الگوریتم k-means است. در مرحله دوم به تعداد خوشه‌ها، ماتریس همبستگی محاسبه می‌شود (خط ۵) و هر ماتریس شامل M^T عنصر است. مرتبه زمانی محاسبه هر عنصر در بدترین حالت^۱ برابر با $O(P.\log P)$ است. از طرفی به تعداد خوشه‌ها یعنی c بار، نیاز به محاسبه ماتریس همبستگی است. بنابراین هزینه کلی این گام در بدترین حالت از مرتبه $O(c.M^T.P \log P)$ خواهد بود. آخرین گام، تخمین مقادیر جافتاده است (خطوط ۶ تا ۱۵). در این مرحله به تعداد مقادیر جافتاده (یعنی Q) یکی از روش‌های رگرسیون خطی با پیچیدگی $O(M^T.P)$ و یا k نزدیک‌ترین همسایگان با پیچیدگی $O(M.P)$ انجام می‌شود. بنابراین پیچیدگی مرحله دوم در بدترین حالت (حالتی که تمام تخمین‌ها از طریق رگرسیون انجام شوند) برابر با $O(Q.M^T.P)$ خواهد بود. با توجه به این که در بدترین حالت



شکل ۳: تأثیر تعداد خوشه‌ها (c) بر روی مقدار RMSE در مجموعه داده‌های (الف) Wholesale Customers، (ب) Haberman، (ج) Wine، (د) Glass و (ه) Iris.



شکل ۴: تأثیر درصد مقایسات (p) بر روی مقدار RMSE در مجموعه داده‌های (الف) Wholesale Customers، (ب) Haberman، (ج) Wine، (د) Glass و (ه) Iris.

بنابراین فرضیه کاهش خطای تخمین با خوشه‌بندی داده‌ها با توجه به نمودارها قابل درک است. مقدار c بهینه برای مجموعه داده‌های مختلف، متفاوت است. آزمایش‌ها نشان می‌دهند افزایش مقدار c از یک حد به بعد، به تدریج باعث افزایش مقدار خطای تخمین می‌شود.

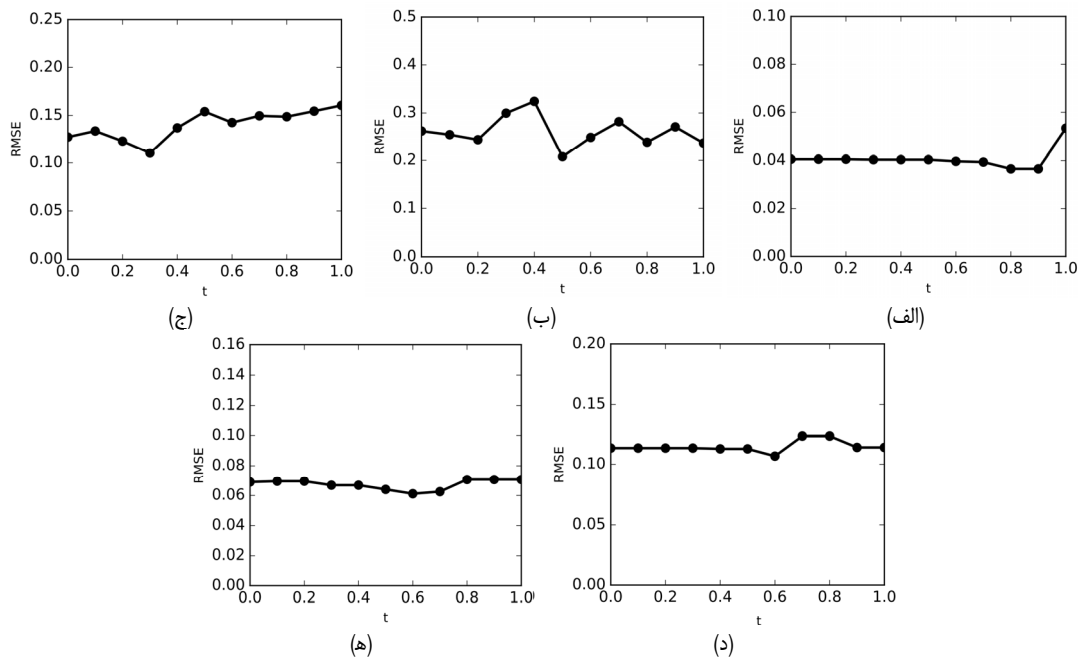
۳-۴ بررسی تأثیر درصد مقایسات (p) برای پرکردن مقدار جاافتاده در هر خوشه

در این قسمت نتایج آزمایش‌هایی برای بررسی تأثیر مقدار p در میزان خطای بررسی می‌شود. ابتدا به تعداد ۱۰٪ کل نمونه‌ها، مقدار جاافتاده به صورت تصادفی در هر مجموعه داده ایجاد شده است. سپس مقدار p بین ۱ تا ۱۰۰ درصد با گام ۵ درصد، برای مقادیر ثابت c و t ، تغییر می‌یابد و نتایج ثبت می‌شوند. شکل ۴ مقدار خطای RMSE را برای مقادیر مختلف p برای مجموعه داده‌های مختلف نمایش می‌دهد. همان

خوشه‌ها (c)، درصد انتخاب نزدیک‌ترین همسایه‌ها (p) و حد آستانه برای تعیین روش تخمین (t) بررسی می‌گردند.

۲-۴ بررسی تأثیر تعداد خوشه‌ها (c) بر کیفیت تخمین داده‌های جاافتاده

در این قسمت نتایج آزمایش‌هایی برای بررسی تأثیر مقدار c در میزان خطای بررسی می‌شود. ابتدا به تعداد ۱۰٪ کل نمونه‌ها، مقدار جاافتاده به صورت تصادفی در هر مجموعه داده ایجاد شده است. سپس مقدار c بین ۱ تا ۱۰ برای مقادیر ثابت p و t ، تغییر می‌یابد و نتایج ثبت می‌شوند. شکل ۳ مقدار خطای RMSE را برای مقادیر مختلف c برای مجموعه داده‌های مختلف نمایش می‌دهد. همان طور که از نمودارها پیداست برای مجموعه داده‌های مختلف خوشه‌بندی داده‌ها مقدار خطای تخمین را کاهش داده و در واقع مقدار خطا در یک مقدار $c \neq 1$ بهینه شده است.



شکل ۵: تأثیر حد آستانه (t) بر روی مقدار RMSE در مجموعه داده‌های، (الف) Wholesale Customers، (ب) Haberman، (ج) Wine، (د) Glass و (ه) Iris

۴-۵ مقایسه با روش پرکردن با دیگر روش‌ها

در این قسمت میزان خطای تخمین داده‌ها در مقایسه با روش پرکردن با مقدار میانگین (Mean)، روش تخمین با شبکه عصبی پرسپترون چندلایه (MLP) [۵]، روش تخمین با خوشه‌بندی c -means فازی (FCM) [۶] و همچنین روش CKNNI [۷] مقایسه می‌شود. ابتدا به تعداد ۵٪، ۱۰٪، ۱۵٪، ۲۰٪ و ۲۵٪ کل نمونه‌ها، مقدار جافتاده به صورت تصادفی در هر مجموعه داده ایجاد شده است. سپس مقادیر جافتاده با چهار روش ذکرشده و روش پیشنهادی در این مقاله، پر شده‌اند. مقدار خطای RMSE برای این رویکردها محاسبه شده و نتایج برای مجموعه داده‌های مختلف ثبت شدند. شکل ۶ مقدار خطای RMSE محاسبه‌شده برای مجموعه داده‌های مختلف برای درصد‌های مختلف مقادیر جافتاده تولیدشده را برای روش‌های ذکرشده نشان می‌دهد. همان طور که از شکل‌ها پیداست برای همه مجموعه داده‌ها به غیر از Haberman رویکرد پیشنهادی خطای تخمین کمتری نسبت به چهار روش دیگر داشته است. اگرچه در مجموعه داده Haberman دقت تخمین روش پیشنهادی نسبت به دیگرها کمتر است با این حال روش پیشنهادی عملکرد بسیار نزدیکی نسبت به بهترین روش در ۵٪، ۲۰٪ و ۲۵٪ جافتادگی داشته است. علت ضعف نسبی روش پیشنهادی در این مجموعه داده احتمالاً تعداد کم صفات پیش‌بینی کننده در این مجموعه داده است. به طور کلی با توجه به نتایج موجود می‌توان مشاهده کرد که خطای روش پیشنهادی در اکثر موارد (به غیر از مجموعه داده Haberman) از دیگر روش‌های مقایسه‌شده کمتر است.

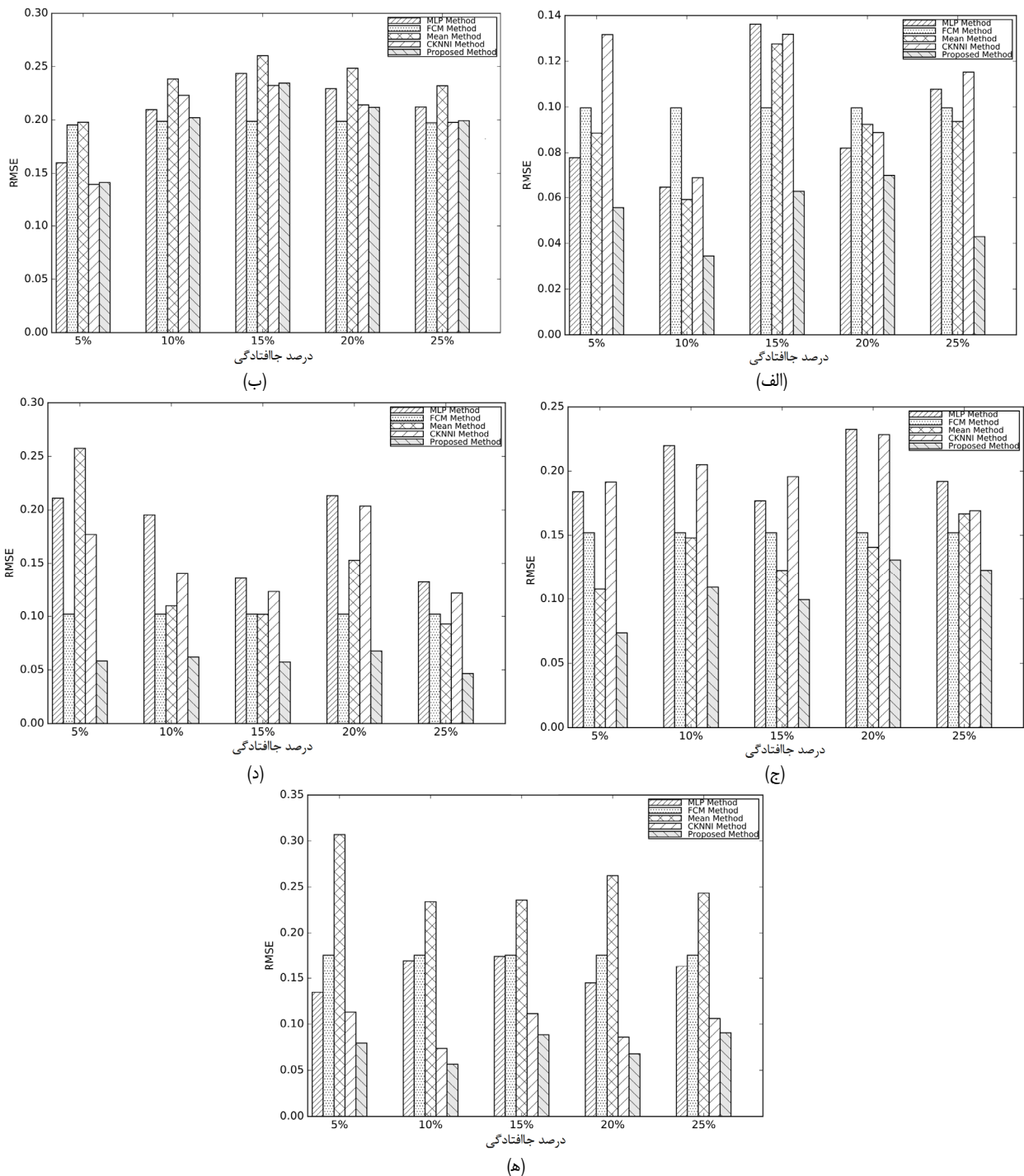
۵- نتیجه‌گیری

در این مقاله روشی جدید و دومارحله‌ای برای تخمین داده‌های جافتاده ارائه گردید. روش پیشنهادی شامل یک مرحله خوشه‌بندی و ترکیب دو روش پایه برای تخمین مقادیر جافتاده است. تأکید روش پیشنهادی بر استفاده از همبستگی‌های موجود بر مجموعه داده‌ها تا حد ممکن برای

طور که از نمودارها پیداست برای مجموعه داده‌های Iris، Glass و Wholesale Customers پس از کاهش مقدار خطا در مقادیر کوچک p ، با افزایش مقدار p ، مقدار خطا نیز یک سیر تقریباً صعودی یکنوا داشته است. برای مجموعه داده Wine نیز اگرچه نوسان برای مقدار خطا وجود داشته است با این حال پس از یک مقدار بهینه ($p = 5\%$)، مقدار خطا یک سیر به طور کلی صعودی را داشته است. برای مجموعه داده Haberman مقدار بهینه p در مقادیر بالاتر حاصل شده است. در واقع نمودارها این مفهوم را بیان می‌کنند که مقدار بهینه k در روش k نزدیک‌ترین همسایه نه بسیار کوچک است و نه بسیار نزدیک به تعداد کل نمونه‌های موجود در هر خوشه. مقدار بهینه p برای هر مجموعه داده را می‌توان طی آزمایش‌ها تعیین نمود.

۴-۴ بررسی تأثیر حد آستانه (t) برای تعیین روش تخمین

در این قسمت نتایج آزمایش‌هایی برای بررسی تأثیر مقدار t در میزان خطا بررسی می‌شود. ابتدا به تعداد ۱۰٪ کل نمونه‌ها، مقدار جافتاده به صورت تصادفی در هر مجموعه داده ایجاد شده است. سپس مقدار t بین ۰ تا ۱ با گام ۰/۱ برای مقادیر ثابت p و c ، تغییر می‌یابد و نتایج ثبت می‌شوند. شکل ۵ مقدار خطای RMSE را برای مقادیر مختلف t برای مجموعه داده‌های مختلف نمایش می‌دهد. همان طور که از نمودارها پیداست برای مجموعه داده‌های مختلف یک مقدار بهینه $0 < t < 1$ توانسته است مقدار خطا را کمینه کند. این نمودارها بیانگر این هستند که ترکیب دو روش k نزدیک‌ترین همسایگان وزن‌دار و رگرسیون خطی توانسته است بهتر از این دو روش به تنهایی عمل کند. همان طور که پیداست در $t = 0$ که بیانگر روش k نزدیک‌ترین همسایگان وزن‌دار به تنهایی و در $t = 1$ که بیانگر روش رگرسیون خطی به تنهایی است، مقادیر خطا بیش از مقدار کمینه برای خطاست. مختلف بودن مقادیر بهینه برای t به این معناست که برای مجموعه داده‌های مختلف یکی از روش‌ها بر دیگری غلبه دارد. در واقع میزان همبستگی درون نمونه‌های افزایشی تشکیل شده از مجموعه داده‌های مختلف، متفاوت است.



شکل ۶: مقایسه مقدار RMSE برای رویکرد ارائه‌شده، روش CKNNI و روش میانگین، برای درصد‌های مختلف جاافتادگی در مجموعه داده‌های (الف) Wholesale Customers، (ب) Haberman، (ج) Wine، (د) Glass و (ه) Iris.

معیارهای سنجش همبستگی مانند معیار اطلاعات متقابل^۱ در کیفیت پرکردن مقادیر جاافتاده در پژوهش‌های آینده بررسی خواهد شد.

مراجع

[1] R. J. A. Little and D. B. Rubin, *Statistical Analysis with Missing Data*, Second Edition, John Wiley Sons, Inc., pp. 11-15, 2002.
 [2] I. B. Aydılek and A. Arslan, "A hybrid method for imputation of missing values using optimized fuzzy c-means with support vector regression and a genetic algorithm," *Information Sciences*, vol. 233, no. 2, pp. 25-35, Jun. 2013.

تخمین مقادیر جاافتاده است. نتایج بررسی‌شده بر روی پنج مجموعه داده نشان داد که رویکرد پیشنهادی میانگین مربعات خطای کمتری در اکثر موارد در مقایسه با چهار روش مختلف دیگر دارد. تأثیر پارامترهای رویکرد پیشنهادی در میزان خطا نیز بررسی شدند.

از مشکلات اصلی رویکرد پیشنهادی می‌توان به چالش انتخاب بهینه برای پارامترهای رویکرد پیشنهادی اشاره کرد اگرچه این چالش نیز با عملیات زمان‌بر جستجوی کامل برای انتخاب بهترین ترکیب پارامترها قابل برطرف شدن است. تأثیر دیگر روش‌های خوشه‌بندی در مرحله اول و ترکیب دیگر روش‌ها در دست بررسی است. به علاوه، تأثیر دیگر

- [18] D. Li, J. Deogun, W. Spaulding, and B. Shuart, "Towards missing data imputation: a study of fuzzy k-means clustering method," in *Rough Sets and Current Trends in Computing*, vol. 3066, pp. 573-579, Jun. 2004.
- [19] N. Ankaiah and V. Ravi, "A novel soft computing hybrid for data imputation," in *Proc. of the 7th Int. Conf. on Data Mining, DMIN'11*, pp. 65-71, Jul. 2011.
- [20] S. Azim and S. Aggarwal, "Hybrid model for data imputation: using fuzzy c means and multi layer perceptron," in *Proc. IEEE Int. Advance Computing Conf., IACC'14*, vol. 4, pp. 1281-1285, Feb. 2014.
- [21] S. Bashir, S. Razzaq, U. Maqbool, S. Tahir, and A. R. Baig, *Using Association Rules for Better Treatment of Missing Values*, arXiv preprint arXiv: 0904.3320, 2009.
- [22] G. Rahman and Z. Islam, "A decision tree-based missing value imputation technique for data pre-processing," in *Proc. of the Ninth Australasian Data Mining Conf., AusDM'11*, vol. 121, pp. 41-50, Dec. 2011.
- [23] C. F. Tsai and F. Y. Chang, "Combining instance selection for better missing value imputation," *J. of Systems and Software*, vol. 122, no. 1, pp. 63-71, Dec. 2016.
- [24] D. R. Wilson and T. R. Martinez, "Reduction techniques for instance-based learning algorithms," *Machine Learning*, vol. 38, no. 3, pp. 257-286, Mar. 2000.
- [25] M. Amiri and R. Jensen, "Missing data imputation using fuzzy-rough methods," *Neurocomputing*, vol. 205, no. 1, pp. 152-164, Sept. 2016.
- [26] M. Lichman, *UCI Machine Learning Repository School of Information and Computer Science*, Irvine, CA: University of California, 2013.
- [27] A. K. Jain, "Data clustering: 50 years beyond k-means," *Pattern Recogn. Lett.*, vol. 31, no. 8, pp. 651-666, Jun. 2010.
- [3] B. van Stein and W. Kowalczyk, "An incremental algorithm for repairing training sets with missing values," in *Proc. Int. Conf. on Information Processing and Management of Uncertainty in Knowledge-Based Systems*, vol. 611, pp. 175-186, Jun. 2016.
- [4] A. A. Chavan and V. K. Verma, "Treatment of missing values for association rules: a recent survey," *International J. of Computer Applications*, vol. 70, no. 26, pp. 1-4, May 2013.
- [5] E. L. Silva-Ramrez, R. Pino-Mejas, and M. Lopez-Coello, "Single imputation with multilayer perceptron and multiple imputation combining multilayer perceptron and k-nearest neighbours for monotone patterns," *Applied Soft Computing*, vol. 29, no. 1, pp. 65-74, Apr. 2015.
- [6] P. Raja and K. Thangavel, "Soft clustering based missing value imputation," in *Proc. Annual Convention of the Computer Society of India: Digital Connectivity-Social Impact*, vol. 679, pp. 119-133, Dec. 2016.
- [7] C. Jiang and Z. Yang, "CKNNI: an improved knn-based missing value handling technique," in *Proc. Int. Conf. on Intelligent Computing*, pp. 441-452, Aug. 2015.
- [8] C. H. Wu, C. H. Wun, and H. J. Chou, "Using association rules for completing missing data," in *Proc. IEEE Fourth Int. Conf. on Hybrid Intelligent Systems, HIS'04*, pp. 236-241, 5-8 Dec. 2004.
- [9] J. Wu, Q. Song, and J. Shen, "An novel association rule mining based missing nominal data imputation method," in *Proc. IEEE Eighth ACIS Int. Conf. on Software Engineering, Artificial Intelligence, Networking, and Parallel/Distributed Computing, SNPD*, vol. 3, pp. 244-249, Jul. 2007.
- [10] N. Singh, A. Javeed, S. Chhabra, and P. Kumar, "Missing value imputation with unsupervised Kohonen self organizing map," in *Emerging Research in Computing, Information, Communication and Applications*, pp. 61-76, Jul. 2015.
- [11] J. Han, M. Kamber, and J. Pei, *Data Mining: Concepts and Techniques*, San Francisco, CA, USA: Morgan Kaufmann Publishers Inc., 3rd Ed., pp. 398-408, 2011.
- [12] R. Krishnamoorthy, S. Sreedhar Kumar, and B. Neelagund, "A new approach for data cleaning process," in *Proc. IEEE Recent Advances and Innovations in Engineering, ICRAIE'14*, 5 pp. Jul. 2014.
- [13] C. Zhang, X. Zhu, J. Zhang, Y. Qin, and S. Zhang, "GBKII: an imputation method for missing values," in *Advances in Knowledge Discovery and Data Mining*, vol. 11, pp. 1080-1087, May 2007.
- [14] E. R. Hruschka, E. R. Hruschka, and N. F. F. Ebecken, "Evaluating a nearest-neighbor method to substitute continuous missing values," in *Proc. Australasian Joint Conf. on Artificial Intelligence*, vol. 16, pp. 723-734, Dec. 2003.
- [15] C. Ye and H. Wang, "Capture missing values based on crowdsourcing," in *Proc. of the 9th Int. Conf. on Wireless Algorithms, Systems, and Applications, WASA'14*, pp. 783-792, Jun. 2014.
- [16] B. M. Patil, R. C. Joshi, and D. Toshniwal, "Missing value imputation based on k-mean clustering with weighted distance," in *Proc. Int. Conf. on Contemporary Computing*, vol. 3, pp. 600-609, Aug. 2010.
- [17] V. V. Ayuyev, J. Jupin, P. W. Harris, and Z. Obradovic, "Dynamic clustering-based estimation of missing values in mixed type data," in *Proc. Int. Conf. on Data Warehousing and Knowledge Discovery*, vol. 11, pp. 366-377, Aug. 2009.

امیرمسعود سفیدیان فارغ التحصیل مقطع کارشناسی ارشد در رشته‌ی مهندسی کامپیوتر- نرم افزار از دانشگاه تربیت دبیر شهید رجایی می‌باشد. نام‌برده تحصیلات خود را در مقطع کارشناسی مهندسی کامپیوتر- نرم‌افزار در سال ۱۳۹۴ با کسب رتبه اول در دانشگاه تربیت دبیر شهید رجایی به اتمام رسانده است. وی در سال ۱۳۹۴ بدون کنکور و با شرایط استعداد درخشان به‌طور مستقیم در مقطع کارشناسی ارشد دانشگاه تربیت دبیر شهید رجایی در رشته‌ی مهندسی کامپیوتر- نرم‌افزار پذیرفته شد. وی در تابستان ۱۳۹۶ تحصیلات خود را در مقطع کارشناسی ارشد با کسب رتبه‌ی اول به پایان رساند. زمینه‌های تحقیقاتی ایشان عبارتند از: جایگذاری مقادیر جاافتاده در داده‌ها، پیش‌پردازش داده‌ها و داده‌کاوی.

نگین دانشپور استادیار دانشکده مهندسی کامپیوتر دانشگاه تربیت دبیر شهید رجایی می‌باشد. نام‌برده تحصیلات خود را در مقطع کارشناسی مهندسی کامپیوتر- سخت‌افزار در سال ۱۳۷۸ با کسب رتبه اول در دانشگاه شهید بهشتی، و کارشناسی ارشد مهندسی کامپیوتر- نرم‌افزار در سال ۱۳۸۱ در دانشگاه صنعتی امیرکبیر به پایان رسانده است، و در سال ۱۳۸۹ دکتری خود در رشته مهندسی کامپیوتر- نرم‌افزار را از دانشگاه صنعتی امیرکبیر اخذ کرده است. زمینه‌های تحقیقاتی مورد علاقه ایشان عبارتند از: پایگاه داده تحلیلی، سیستم‌های تصمیم‌یار، پیش‌پردازش داده‌ها و داده‌کاوی.