

تکمیل ماتریس گراف در حضور داده پرت

علیرضا احمدی^۱ سینا مجیدیان^۲ محمدحسین کهائی^۳

۱- فارغ التحصیل کارشناسی ارشد- دانشکده مهندسی برق - دانشگاه علم و صنعت ایران - تهران - ایران -

ar_ahmadi71@yahoo.com

۲- فارغ التحصیل دکتری- دانشکده مهندسی برق - دانشگاه علم و صنعت ایران - تهران - ایران -

s_majidian@elec.iust.ac.ir

۳- دانشیار- دانشکده مهندسی برق - دانشگاه علم و صنعت ایران - تهران - ایران -

kahaei@iust.ac.ir

چکیده: در سال‌های اخیر موضوع تکمیل ماتریس بسیار مورد توجه محققان قرار گرفته است. در مسأله‌ی تکمیل ماتریس، هدف بازیابی کامل یک ماتریس رتبه پایین است که با استفاده از مشاهداتی تعداد محدودی از درایه‌های آن ماتریس انجام می‌شود. همچنین، مدل‌سازی ارتباط بین سطرها و ستون‌های ماتریس به صورت یک گراف موجب معرفی زمینه پژوهشی تکمیل ماتریس گراف شده است. در مسأله‌ی تکمیل ماتریس گراف، بازیابی ماتریس با استفاده از داده‌های مشاهده شده از طریق افزودن عبارات تغییرات کل گراف به تابع هدف مسأله‌ی تکمیل ماتریس انجام می‌پذیرد. اما در عمل داده‌ها، آغشته به نویز و شامل داده‌های پرت می‌باشند. به داده‌هایی که با سایر داده‌های مشاهده شده متفاوت باشند و از ساختار کلی آن‌ها پیروی نکنند، داده‌ی پرت گفته می‌شود. در این مقاله، روشی جدید برای تکمیل ماتریس گراف در شرایط وجود همزمان نویز و داده‌ی پرت در مشاهدات ارائه شده است. روش پیشنهادی با نام GMCO-DL از ماتریس لاپلاسیان جهت‌دار برای تعریف تغییرات کل گراف استفاده می‌نماید. نتایج شبیه‌سازی روش پیشنهادی حاکی از بهبود قابل ملاحظه‌ای از نظر خطا می‌باشد.

واژه‌های کلیدی: تکمیل ماتریس، پردازش سیگنال گراف، گراف جهت‌دار، داده پرت

نوع مقاله: پژوهشی

DOI: 10.52547/jiaeee.20.1.89

تاریخ ارسال مقاله: ۱۳۹۹/۱۱/۲۳

تاریخ پذیرش مشروط مقاله: ۱۴۰۰/۱۲/۲۷

تاریخ پذیرش مقاله: ۱۴۰۱/۷/۲۷

نام نویسنده‌ی مسئول: دکتر محمد حسین کهائی

نشانی نویسنده‌ی مسئول: ایران- تهران- دانشگاه علم و صنعت ایران- دانشکده مهندسی برق

۱- مقدمه

۱-۱- معرفی تکمیل ماتریس

پس از افزایش توجه پژوهشگران به موضوع تحقیقاتی حسگری فشرده، تعمیم آن به ابعاد بالاتر مورد بررسی قرار گرفت [۱]. این امر موجب معرفی مبحث بازیابی ماتریس یا به بیان دیگر مسأله‌ی تکمیل ماتریس شده است. تکمیل ماتریس به دنبال کامل نمودن یک ماتریس رتبه پایین^۱ از روی مشاهده‌ی تعداد محدودی از درایه‌های آن ماتریس می‌باشد [۲-۶]. در واقع تنک بودن ماتریس، معادل با تنک بودن بردار مقادیر تنکین^۲ ماتریس تعریف می‌شود. فرض رتبه پایین بودن ماتریس، شرط اصلی برای تکمیل ماتریس می‌باشد. در حالت کلی، بازیابی ماتریس بوسیله‌ی مشاهده‌ی برخی از درایه‌های آن امکان‌پذیر نمی‌باشد و اعمال شرط رتبه پایین بودن موجب کاهش تعداد جواب‌های ممکن می‌شود. اما شرط رتبه پایین بودن نیز شرط کافی برای یافتن پاسخ مسأله‌ی تکمیل ماتریس نمی‌باشد. برای تضمین بازیابی ماتریس از روی مشاهده‌ی تعداد محدودی از درایه‌های آن از معیار همدوسی استفاده می‌شود. به کمک معیار همدوسی ماتریس می‌توان کران پایین حداقل تعداد درایه مورد نیاز جهت بازیابی ماتریس را به دست آورد [۶].

ماتریس‌های رتبه پایین نقش مهمی در مسائل مربوط به پردازش سیگنال ایفا می‌کنند. در بسیاری از زمینه‌های تحقیقاتی همچون تخمین هاپلوتا پ [۲-۷ و ۸]، تخمین جهت منابع [۴ و ۱۴]، یادگیری ماشین^۳ [۱۵] و برش‌نگاری حالت کوانتوم^۴ [۱۶]، با این دسته از ماتریس‌ها مواجه هستیم. به بیان دیگر بسیاری از داده‌های موجود در جهان واقعی از مدل ماتریس‌های رتبه پایین تبعیت می‌کنند.

یکی از مسائل معروف تکمیل ماتریس، سیستم پیشنهادکننده‌ی سایت اجاره فیلم نت‌فلیکس است. در این مسأله، کاربران مختلف به عنوان سطرها و فیلم‌ها به عنوان ستون‌های ماتریس X در نظر گرفته می‌شوند. ماتریس $X \in \mathbb{R}^{N \times L}$ در حالت کلی یک ماتریس مستطیلی است. درایه‌های آن نشان‌دهنده‌ی نحوه‌ی نمره‌دهی کاربران مختلف به فیلم‌های گوناگون می‌باشند. به عنوان مثال درایه‌ی x_{ij} بیانگر نمره‌ای است که کاربر i به فیلم j ام اختصاص داده است. هدف این است با استفاده از تکمیل ماتریس و نحوه‌ی نمره‌دهی کاربران به فیلم‌هایی که قبلاً مشاهده نموده‌اند، میزان نمره‌دهی آن‌ها به فیلم‌هایی که مشاهده نکرده‌اند را پیش‌بینی نماییم. در واقع قصد داریم، با فرض دانستن تعداد کمی از درایه‌ها، ماتریس را به صورت کامل بازیابی نماییم [۱۱-۱۳ و ۱۷].

مثال دیگر برای سیستم پیشنهادکننده مسأله‌ی تخمین دمای پایگاه‌های هواشناسی می‌باشد. هدف این مسأله تخمین دمای

پایگاه‌های هواشناسی در طول سال، با دانستن دمای پایگاه‌های هواشناسی در برخی از ایام سال می‌باشد. دمای پایگاه‌های هواشناسی به عنوان سطرها و روزهای مختلف سال به عنوان ستون‌های ماتریس در نظر گرفته می‌شوند. به عبارت دیگر درایه‌ی x_{ij} ، دمای است که پایگاه i ام در روز j ام ثبت کرده است. در این مسأله بازیابی ماتریس با استفاده از در نظر گرفتن فاصله‌ی جغرافیایی بین پایگاه‌های مختلف انجام می‌پذیرد [۹-۱۰].

۲-۱- تکمیل ماتریس با مشاهدات آغشته به داده

پرت

وجود داده‌ی پرت^۵ در بسیاری از زمینه‌های تحقیقاتی گزارش شده است [۷-۹، ۲۳-۲۶]. در اصطلاح به داده‌هایی که با سایر داده‌های مشاهده شده متفاوت باشند و از ساختار کلی آن‌ها پیروی نکنند داده‌ی پرت گفته می‌شود. در مسأله‌ی تکمیل ماتریس نیز امکان وجود داده‌ی پرت در مشاهدات وجود دارد. به عنوان مثال در مسأله‌ی نت‌فلیکس، کاربر می‌تواند به صورت سهوی یا عمدی نمره‌ی اشتباهی به یک فیلم اختصاص دهد. به عبارت دیگر وجود داده‌ی پرت می‌تواند ویژگی رتبه پایین بودن ماتریس را با مشکل روبه‌رو سازد. در صورت وجود داده‌ی پرت در مجموعه مشاهدات، بازیابی صحیح ماتریس از روی مشاهده‌ی تعداد محدودی از درایه‌های آن امکان‌پذیر نمی‌باشد. معرفی روشی که علاوه بر بازیابی ماتریس از میان مشاهدات ناقص آغشته شده به نویز، توانایی شناسایی داده‌ی پرت را داشته باشد، عملکرد تکمیل ماتریس را بهبود می‌بخشد.

۳-۱- پردازش سیگنال گراف

پردازش داده‌های حجیم و کاربرد آن در بسیاری از زمینه‌های مختلف علمی و مهندسی موجب ایجاد چالش‌های جدیدی در حوزه‌ی پردازش سیگنال‌های گسسته در زمان شده است. محاسبات زیاد و پیچیده‌ی روش‌های پردازش سیگنال کلاسیک موجب شده است تا اینگونه روش‌ها کارایی خود را برای داده‌های با ابعاد بالا از دست بدهند. بدین منظور شناسایی روش‌هایی که منجر به کاهش محاسبات ناشی از ابعاد بالای داده‌ها شود در دستور کار محققین این حوزه قرار گرفته است. این امر موجب معرفی زمینه‌ی تحقیقاتی پردازش سیگنال گسسته به روی گراف یا DSP_G شده است. DSP_G از ادغام مفاهیم تئوری گراف و تحلیل‌های محاسباتی پردازش سیگنال بوجود آمده است. در این مبحث ارتباط بین نمونه‌های سیگنال گسسته، به صورت گراف مدل‌سازی می‌شود. تعمیم مفاهیم پردازش سیگنال گسسته در زمان

[۱۹].

سیگنال گراف را می‌توان به صورت نگاشتی از مجموعه‌ی رئوس V به مجموعه اعداد حقیقی R و یا در حالت کلی مختلط در نظر گرفت و به صورت زیر تعریف نمود:

$$s : \begin{matrix} V \rightarrow R, \\ v_n \rightarrow s_n \end{matrix} \quad (۱)$$

در واقع به هر رأس v_i یک مقدار حقیقی (یا مختلط) s_i اختصاص داده شده است. به بیان دیگر سیگنال گراف s تحت یک گراف مشخص را به صورت بردار نمایش می‌دهند:

$$s = [s_0 \ s_1 \ \dots \ s_{N-1}]^T \in R^N \quad (۲)$$

که در آن مقدار N برابر با تعداد رئوس گراف می‌باشد [۱۰].

۲- مدل مساله

۲-۱- تکمیل ماتریس با استفاده از نرم هسته‌ای

برای توصیف مساله تکمیل ماتریس نیازمند به تعریف اپراتور مشاهدات هستیم. مجموعه‌ی M به عنوان زیر مجموعه‌ای از اندیس‌های درایه‌های ماتریس $X \in R^{N \times L}$ به صورت زیر در نظر گرفته شده است [۵]:

$$M \subseteq \{1, 2, \dots, N\} \times \{1, 2, \dots, L\} \quad (۳)$$

در این صورت اپراتور مشاهدات $P : R^{N \times L} \rightarrow R^{N \times L}$ به صورت زیر در نظر گرفته می‌شود:

$$[P(X)]_{ij} = \begin{cases} X_{ij} & (i, j) \in M \\ 0 & \text{otherwise} \end{cases} \quad (۴)$$

فرض می‌شود $T \in R^{N \times L}$ ماتریس رتبه پایینی باشد که قصد بازیابی آن را داشته باشیم. مساله تکمیل ماتریس به صورت زیر معرفی می‌شود:

$$\begin{aligned} &\text{minimize} \quad \text{rank}(X) \\ &\text{subject to} \quad P(X) = P(T) \end{aligned} \quad (۵)$$

مساله فوق به دلیل وجود عبارت رتبه‌ی ماتریس در تابع هدف نامحذب است. همچنین این مساله به عنوان یک مساله‌ی دشوار ۷ شناخته می‌شود. برای حل این مشکل از عبارت نرم هسته‌ای ۸ به عنوان جایگزین محذب رتبه‌ی ماتریس استفاده می‌شود. بدین ترتیب مساله‌ی غیرمحذب به یک مساله‌ی بهینه‌سازی محذب تبدیل می‌گردد [۱۲].

$$\begin{aligned} &\text{minimize} \quad \|X\|_* \\ &\text{subject to} \quad P(X) = P(T) \end{aligned} \quad (۶)$$

برای سیگنال‌هایی با ساختار پیچیده، از فواید DSP_G محسوب می‌شود [۱۲ و ۱۹-۲۰]. بدین منظور شناسایی روش‌هایی که منجر به کاهش محاسبات ناشی از ابعاد بالای داده‌ها شود در دستور کار محققین این حوزه قرار گرفته است. این امر موجب معرفی زمینه‌ی تحقیقاتی پردازش سیگنال گسسته به روی گراف یا DSP_G شده است. DSP_G از ادغام مفاهیم تئوری گراف و تحلیل‌های محاسباتی پردازش سیگنال بوجود آمده است. در این مبحث ارتباط بین نمونه‌های سیگنال گسسته به صورت گراف مدل‌سازی می‌شود. تعمیم مفاهیم پردازش سیگنال گسسته در زمان برای سیگنال‌هایی با ساختار پیچیده، از فواید DSP_G محسوب می‌شود [۲۱].

۴-۱- معرفی سیگنال گراف

گراف یک مدلی ریاضی برای نمایش ارتباط بین اعضای یک مجموعه می‌باشد. برای نمایش گراف از مجموعه رئوس و یال‌هایی که رئوس مختلف را به هم متصل می‌کنند، استفاده می‌شود. نظریه‌ی گراف یکی از موضوعات مهم در ریاضیات گسسته است که به مطالعه‌ی گراف‌ها و مدل‌سازی مسائل تحقیقاتی به وسیله‌ی ساختارهای گراف می‌پردازد. لئونارد اویلر در سال ۱۷۳۶ میلادی و با حل مسئله‌ی پل‌های کونیسبرگ، اقدام به معرفی نظریه‌ی گراف نمود و واژه‌ی گراف برای اولین بار در سال ۱۸۷۸ میلادی توسط جوزف سیلستر برای این مدل ریاضی استفاده شد. نظریه‌ی گراف به بررسی مفاهیم پایه گراف مانند درجه، مسیر، قدم، درخت، گراف‌های کامل، زیرگراف و بسجاری از مسائل دیگر می‌پردازد. روش‌های مختلفی برای بیان ریاضی ساختار گراف وجود دارد. استفاده از ماتریس مجاورت^۹ گراف یکی از پرکاربردترین این روش‌ها می‌باشد. گراف را می‌توان به صورت $G = (V, A)$ نمایش داد که در آن مجموعه‌ی $V = \{v_0, v_2, \dots, v_{N-1}\}$ یک مجموعه‌ی N عضو نشان دهنده‌ی رئوس گراف و ماتریس A بیانگر ماتریس مجاورت گراف می‌باشد. درایه‌های ماتریس مجاورت را به صورت a_{nm} نشان می‌دهیم. در ابتدا و برای سادگی، ماتریس مجاورت را باینری در نظر می‌گیریم. وجود یا عدم وجود یال بین رئوس مختلف گراف توسط درایه‌های ماتریس مجاورت بیان می‌شود. درایه‌ی a_{nm} بیانگر وجود یک یال از رأس v_m به رأس v_n است؛ در واقع اگر مقدار a_{nm} برابر ۱ باشد، یک یال از رأس v_m به رأس v_n وجود دارد. به طور مشابه $a_{nm} = 0$ بیانگر عدم وجود یال از رأس v_m به رأس v_n می‌باشد. همچنین می‌توان از ماتریس مجاورت وزن‌دار برای بیان ارتباط دقیق‌تر میان رئوس گراف استفاده نمود. درایه‌های ماتریس مجاورت وزن‌دار مقادیر مختلف حقیقی و یا مختلط را اختیار می‌کنند

مجموعه همسایگی آن رأس بازیابی شود. در واقع اضافه شدن عبارت نرم تغییرات کل موجب می شود تا سیگنال گراف متناظر با ستون های ماتریس بازیابی شده \mathbf{X} هموار باشند. به عبارت دیگر ستون های ماتریس \mathbf{X} در زیر فضای فرکانس های کم گراف قرار دارند. عدد اسکالر β نیز میزان تاثیرگذاری دو عبارت تابع هدف در بازیابی \mathbf{X} را کنترل می نماید. در واقع همانند روش SVT [۸]، هر چه β بزرگتر باشد، تاثیرگذاری نرم هسته ای نسبت به نرم تغییرات کل افزایش می یابد. مسئله ذکر شده با کمینه سازی مجموع نرم تغییرات کل و نرم هسته ای متغیر \mathbf{X} ، ماتریس \mathbf{T} را بازیابی می نماید. ماتریس بازیابی شده ماتریسی رتبه پایین و با ستون های هموار خواهد بود. می توان به جای استفاده از اپراتور مشاهدات $P: \mathbb{R}^{N \times L} \rightarrow \mathbb{R}^{N \times L}$ از ماتریس مشاهدات ناقص برای بیان ساده تر مسئله استفاده نمود. ماتریس مشاهدات ناقص، برای مجموعه ای مشاهدات $M \subseteq \{1, 2, \dots, N\} \times \{1, 2, \dots, L\}$ ، به صورت زیر تعریف می شود:

$$(T_M)_{ij} = \begin{cases} (T)_{ij} & (i, j) \in M \\ 0 & \text{o.w} \end{cases} \quad (9)$$

در این صورت مسئله ای رابطه ای (۸) به صورت زیر بازنویسی و GMCM^۹ نامیده می شود [۱۸]:

$$\begin{aligned} & \text{minimize} \quad S_2(\mathbf{X}) + \beta \|\mathbf{X}\|_* \\ & \text{subject to} \quad \mathbf{X}_M = \mathbf{T}_M \end{aligned} \quad (10)$$

در حالتی که داده با نویز آغشته شده باشند، \mathbf{T}_M به صورت زیر تعریف می گردد:

$$\mathbf{T}_M = \mathbf{X}_M + \mathbf{Z}_M = (\mathbf{X} + \mathbf{Z})_M \quad (11)$$

ماتریس $\mathbf{Z} \in \mathbb{R}^{N \times L}$ ماتریس نویز می باشد. بدین ترتیب، مسئله تکمیل ماتریس گراف در حالت وجود نویز در مشاهدات، به صورت زیر تعریف می گردد:

$$\begin{aligned} & \text{minimize} \quad S_2(\mathbf{X}) + \beta \|\mathbf{X}\|_* \\ & \text{subject to} \quad \|\mathbf{X} - \mathbf{T}_M\|_F^2 \leq \varepsilon^2 \end{aligned} \quad (12)$$

که پارامتر ε مقدار خطای مجاز برای درایه های مشاهده شده را کنترل می نماید. مسئله ای بهینه سازی مقید فوق را می توان به یک مسئله بهینه سازی غیرمقید تبدیل نمود. برای این کار از روش تنظیم پارامترها استفاده می شود. مساله بهینه سازی به دست آمده روش GMCR^{۱۰} نامگذاری شده است [۹]:

که در آن $\|\cdot\|_*$ نماد نرم هسته ای ماتریس است و به صورت حاصل جمع مقادیر تکین یک ماتریس تعریف می شود. تحت شرایط مناسب مسئله کمینه سازی رتبه و مسئله محدب آن، معادل هستند و دارای پاسخ یکتا می باشند. در بحث فوق، مسئله تکمیل ماتریس در حالت عدم وجود نویز بررسی شد. هنگام مشاهده ی نویزی درایه های \mathbf{T} ، مسئله ی تکمیل ماتریس به صورت زیر بیان می شود:

$$\begin{aligned} & \text{minimize} \quad \|\mathbf{X}\|_* \\ & \text{subject to} \quad \|\mathbf{P}(\mathbf{X}) - \mathbf{P}(\mathbf{T})\|_F \leq \varepsilon \end{aligned} \quad (13)$$

پارامتر ε وظیفه ی کنترل مقدار نویز مجاز در درایه ها را بر عهده دارد. همچنین، $\|\cdot\|_F$ بیانگر نرم فربینیوس بوده و مقدار آن برابر با ریشه ی دوم حاصل جمع مجذور درایه های ماتریس می باشد. در ادامه بررسی مساله تکمیل ماتریس با استفاده از مفاهیم گراف می پردازیم.

۲-۲- تکمیل ماتریس گراف

هر سیگنال گراف مدل سازی شده توسط گراف $G = (V, A)$ ، به صورت یک بردار در نظر گرفته می شود، $\mathbf{x}^{(l)} \in \mathbb{R}^N$ ها ستون های ماتریس $\mathbf{X} \in \mathbb{R}^{N \times L}$ هستند و به عنوان سیگنال گراف در نظر گرفته می شود. همه ی $\mathbf{x}^{(l)}$ ها از ساختار یکسان گرافی $G = (V, A)$ تبعیت می کنند. به عبارت دیگر ارتباط بین تمام درایه های $\mathbf{x}^{(l)}$ ها توسط یک ساختار گراف مشخص، مدل سازی می شود.

یکی از مشکلاتی که در روش های ارائه شده برای مسئله ی تکمیل ماتریس با آن رو به رو هستیم، عدم در نظر گرفتن ارتباط موجود بین درایه های سطرها و یا ستون های ماتریس است. در حالت کلی در نظر گرفتن چنین ارتباطی بین درایه های ماتریس، بهبود در عملکرد بازسازی آن ماتریس را به دنبال خواهد داشت. با در نظر گرفتن ماتریس $\mathbf{T} \in \mathbb{R}^{N \times L}$ به عنوان ماتریس داده ها، مسئله ی تکمیل ماتریس گراف در حالت بدون نویز به صورت زیر تعریف می شود [۱۶]:

$$\begin{aligned} & \text{minimize} \quad S_2(\mathbf{X}) + \beta \|\mathbf{X}\|_* \\ & \text{subject to} \quad \mathbf{P}(\mathbf{X}) = \mathbf{P}(\mathbf{T}) \end{aligned} \quad (14)$$

که $S_2(\mathbf{X}) = \sum_{l=1}^L S_2(\mathbf{x}^{(l)}) = \|\mathbf{X} - \mathbf{A}^{\text{norm}} \mathbf{X}\|_F^2$ برابر با $S_2(\mathbf{X})$ نرم تغییرات کل گراف است.

همچنین، \mathbf{A}^{norm} ماتریس مجاورت نرمالیزه گراف $\mathbf{A} = \frac{1}{|\lambda_{\max}|} \mathbf{A}$ است. افزودن عبارت نرم تغییرات کل گراف موجب بهبود عملکرد مسئله ی تکمیل ماتریس می شود. کمینه سازی این عبارت موجب می شود، مقدار سیگنال در هر رأس گراف یا همان درایه های ستون های ماتریس \mathbf{X} ، به صورت ترکیب خطی از مقادیر سیگنال در رؤس

$$\begin{aligned} & \text{minimize} \quad \alpha \|L_{di}X\|_F^2 + \beta \|X\|_* + \gamma \|W\|_0 \\ & X, W, Z \\ & \text{subject to} \quad \|Z_M\|_F^2 \leq \varepsilon^2 \\ & T_M = (X + Z + W)_M \end{aligned} \quad (17)$$

در رابطه‌ی فوق، $\| \cdot \|_0$ بیانگر نرم l_0 است و مقدار آن با تعداد درایه‌های غیر صفر ماتریس برابر می‌شود. نرم l_0 محدب نمی‌باشد و برای تعریف مسأله‌ی بهینه‌سازی محدب از نرم l_1 استفاده می‌کنیم. دلیل استفاده از نرم l_0 آگاهی از صفر بودن تعداد زیادی از درایه‌های ماتریس داده‌ی پرت می‌باشد. بدین ترتیب روش پیشنهادی GMCO-DL با حل مساله بهینه‌سازی زیر انجام می‌پذیرد.

$$\begin{aligned} & \text{minimize} \quad \alpha \|L_{di}X\|_F^2 + \beta \|X\|_* + \gamma \|W\|_1 \\ & X, W, Z \\ & \text{subject to} \quad \|Z_M\|_F^2 \leq \varepsilon^2 \\ & T_M = (X + Z + W)_M \end{aligned} \quad (18)$$

که در آن $\| \cdot \|_1$ نرم l_1 است و برای ماتریس $W \in R^{N \times L}$ از رابطه‌ی زیر محاسبه می‌شود:

$$\|W\|_1 = \sum_{i=1}^N \sum_{j=1}^L w_{ij} \quad (19)$$

۴- نتایج شبیه‌سازی

داده‌ی مورد استفاده در این مقاله مربوط به دمای ۱۵۰ پایگاه هواشناسی ایالات متحده در طول یک سال می‌باشد که برای بررسی عملکرد روش GMCR استفاده شده است [۱۰ و ۱۷]. برای تولید داده، ماتریس دمای پایگاه‌ها را با $T \in R^{N \times L}$ نمایش می‌دهیم. هر سطر ماتریس متناظر یک پایگاه هواشناسی و همچنین هر ستون ماتریس نماینده یک روز از سال می‌باشد. بنابراین برای ابعاد ماتریس T مقادیر $N=150$ و $L=365$ حاصل می‌شود. برای دستیابی به ساختار گراف مسأله از فاصله‌ی جغرافیایی پایگاه‌ها نسبت به هم استفاده شده است. هر رأس گراف به یک پایگاه اختصاص داده شده است و هر رأس به ۸ رأس دیگری که پایگاه‌های متناظر با آن‌ها از نظر فاصله‌ی جغرافیایی فاصله‌ی کمتری دارند، وصل شده است. به عبارت دیگر هر رأس توسط ۸ یال جهت‌دار به سایر رؤس وصل شده است؛ بنابراین ماتریس مجاورت گراف، جهت‌دار می‌باشد. به طور شهودی می‌توان گفت که مناطقی که از نظر جغرافیایی به هم نزدیکتر هستند دمای مشابهی دارند. ساختار گراف حاکم بر داده‌های در نظر گرفته شده به صورت شکل ۱ نمایش داده می‌شود. این شکل نشان‌دهنده‌ی سیگنال گراف

$$\begin{aligned} & \text{minimize} \quad \| (X - T)_M \|_F^2 + \alpha S_2(X) + \beta \|X\|_* \\ & X \end{aligned} \quad (13)$$

در این روش از پارامترهای تنظیم α و β استفاده شده است.

۳- روش پیشنهادی برای تکمیل ماتریس گراف در حضور داده پرت (GMCO-DL)

در این بخش حالت وجود نویز و داده‌ی پرت را در درایه‌های مشاهده‌ی شده در نظر می‌گیریم. ابتدا مسأله را برای مشاهدات نویزی معرفی می‌کنیم. از مدل نویز سفید گوسی جمع‌شونده ($AWGN$)^{۱۱} برای آغشته کردن مشاهدات به نویز استفاده می‌کنیم. برای در نظر گرفتن نویز، ماتریس مشاهدات ناقص را به صورت زیر تعریف می‌کنیم:

$$T_M = X_M + Z = X_M + Z_M \quad (14)$$

در رابطه‌ی فوق $Z \in R^{N \times L}$ ماتریس نویز و $Z_M \in R^{N \times L}$ ماتریس مشاهدات ناقص نویز می‌باشند. به دلیل اینکه تنها به درایه‌هایی که مشاهده می‌شوند می‌توان نویز اضافه نمود، لذا ماتریس‌های Z و Z_M معادل یکدیگر هستند. مسأله تکمیل ماتریس در حالت وجود نویز در درایه‌های مشاهده شده در نظر می‌گیریم.

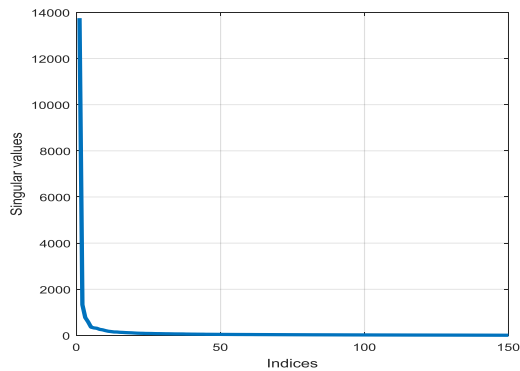
$$\begin{aligned} & \text{minimize} \quad \|L_{di}X\|_F^2 + \beta \|X\|_* \\ & X, Z \\ & \text{subject to} \quad \|Z\|_F^2 \leq \varepsilon^2 \end{aligned} \quad (15)$$

که در آن L_{di} ماتریس لاپلاسین جهت‌دار گراف است که به صورت $L_{di} = D_{in} - A$ تعریف می‌شود و D_{in} ماتریس درجه‌ی ورودی رؤس گراف است [۲۰]. ماتریس مشاهدات ناقص با وجود داده‌ی پرت به صورت زیر تعریف می‌نماییم:

$$T_M = X_M + Z_M + W = X_M + Z_M + W_M \quad (16)$$

که $W \in R^{N \times L}$ ماتریس داده‌ی پرت و $W_M \in R^{N \times L}$ ماتریس مشاهدات ناقص داده‌ی پرت می‌باشند. ماتریس‌های W و W_M معادل یکدیگرند؛ زیرا وجود داده‌ی پرت را می‌توان تنها در داده‌های مشاهده شده در نظر گرفت. ساختار ماتریس داده‌ی پرت را به صورت تنک در نظر می‌گیریم. به این معنا که اکثر درایه‌های $W \in R^{N \times L}$ مقدار صفر را اختیار می‌کنند. مسأله‌ی تکمیل ماتریس در حضور نویز و داده‌ی پرت را به صورت زیر بیان می‌کنیم:





شکل (۱): (بالا) ساختار گراف پایگاه‌های هواشناسی ایالات متحده

[۱۹] (پایین) نمودار مقادیر تکین ماتریس داده T

برای هر درصد مشاهدات مشخص ۱۰ بار این مرحله را تکرار می‌کنیم. همچنین، در هر مرحله دمای ۵۰ روز از روز های سال را به صورت تصادفی انتخاب کرده و به عنوان ستون‌های T در نظر می‌گیریم. در واقع برای ابعاد ماتریس مقادیر $N = 150$ و $L = 50$ حاصل می‌شود. در این بخش عملکرد روش GMCO-DL را در شناسایی داده‌های پرت مورد بررسی قرار می‌دهیم. مسأله‌ی بهینه‌سازی GMCO-DL رابطه‌ی (۱۸) بیان شده است. برای پیاده‌سازی این ماتریس، پارامتری را به عنوان درصد داده‌ی پرت تعریف می‌نماییم. این پارامتر بیانگر، تعداد درایه‌های آغشته شده به داده‌ی پرت نسبت به تعداد کل درایه‌های مشاهده شده می‌باشد. با در نظر گرفتن یک مقدار ثابت برای SNR و در درصد مشاهدات متفاوت اقدام به شناسایی درایه‌هایی که با داده‌های پرت آغشته شده‌اند می‌نماییم. جدول ۱ توضیحات مربوط به شبیه‌سازی این بخش را نشان می‌دهد:

جدول (۱): شرایط شبیه‌سازی روش GMCO-DL در مشاهدات

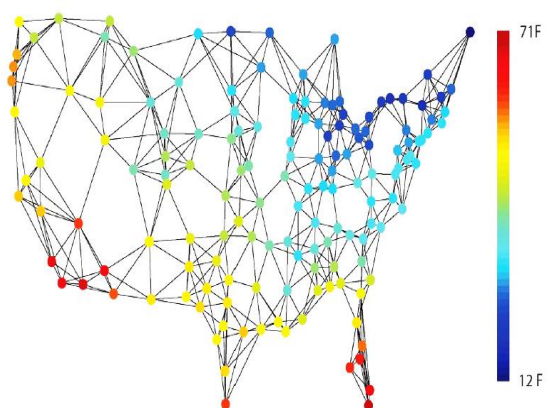
یکنواخت

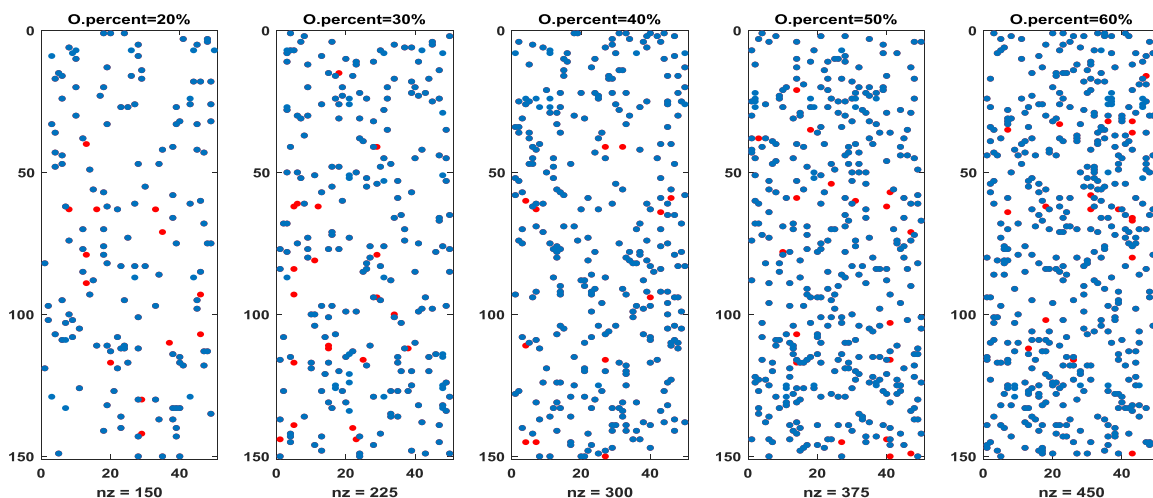
| ابعاد ماتریس داده | درصد مشاهدات | درصد داده‌ی پرت | SNR(dB) |
|---------------------------|--------------|-----------------|---------|
| $T \in R^{150 \times 50}$ | ۲۰٪ تا ۶۰٪ | ۱۰٪ | ۵ |

نتایج حاصل از شناسایی داده‌های پرت در شکل ۲ نشان داده شده است که از پنج نمودار تشکیل شده است و هر نمودار، نمایی از ماتریس در درصد مشاهدات مختلف را نشان می‌دهد. در این نمودارها، نقاط قرمز بیانگر درایه‌هایی است که به آن‌ها داده‌ی پرت اضافه شده است و GMCO-DL نتوانسته آن‌ها را تشخیص دهد. همچنین نقاط آبی نشان‌دهنده‌ی درایه‌هایی است که GMCO-DL به عنوان داده‌ی پرت شناسایی کرده است. تعداد کل درایه‌های آغشته شده به داده‌ی

دمای پایگاه‌های هواشناسی در یک روز می‌باشد. سیگنال‌های گراف شکل ۱ (بالا) را به صورت $t^{(l)} \in R^N$ و به عنوان ستون ماتریس T در نظر گرفته می‌شود. هر $t^{(l)}$ بیانگر دمای پایگاه‌های مختلف در یک روز مشخص می‌باشد. در شکل ۱ (پایین) نمودار مقادیر تکین ماتریس داده T رسم شده است. از آنجایی که ستون‌های ماتریس داده‌ی T ، سیگنال‌های گراف هستند و همگی از ساختار گراف نشان داده شده در شکل ۱ تبعیت می‌کنند، نتیجه می‌شود که ماتریس داده تقریباً رتبه‌پایین است؛ پنج مقدار تکین ماتریس نسبت به بقیه آن‌ها، مقادیر بزرگتری را اختیار نموده‌است.

شبیه‌سازی‌ها با استفاده از نرم افزار MATLAB و بسته نرم افزاری CVX به منظور حل مسائل پیچیده‌ی بهینه‌سازی محدب استفاده شده است. در این بخش میزان خطای بازیابی ماتریس را بر حسب درصد درایه‌های مشاهده شده رسم می‌نماییم. به عنوان مثال فرض کنید که درصد مشاهدات درایه‌ها برابر ۲۰٪ باشد. در این صورت، در هر مرحله از اجرای روش‌های تکمیل ماتریس، ۲۰٪ از درایه‌ها را به صورت تصادفی انتخاب می‌شود. سپس، ماتریس مشاهدات ناقص با استفاده از رابطه‌ی (۹) تشکیل داده می‌شود. در شبیه‌سازی‌ها از مشاهدات یکنواخت استفاده شده است. به عبارت دیگر فرض شده است که از هر سطر و ستون حداقل یک درایه مشاهده شده است و احتمال انتخاب سایر درایه‌ها یکسان است. سپس، با استفاده از درایه‌های مشاهده شده آغشته شده به نویز و داده‌ی پرت، روش پیشنهادی GMCO-DL برای بازیابی سایر درایه‌های ماتریس استفاده شده است.



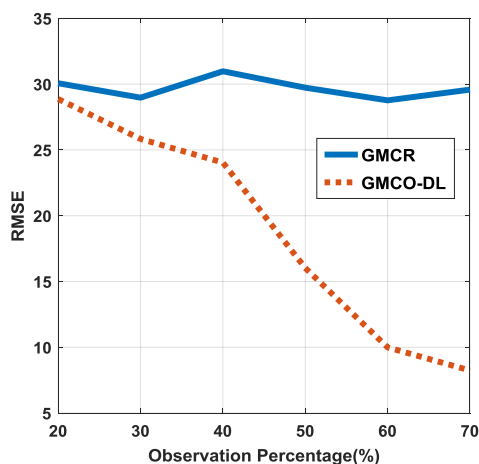


شکل (۲): عملکرد GMCO-DL در شناسایی پرت برای درصد مشاهدات متفاوت

نشده است.

۵- نتیجه گیری

در این مقاله، مساله تکمیل ماتریس گراف در حضور داده پرت مورد بررسی قرار گرفته است. بدین منظور روش GMCO-DL معرفی شده است. در این روش از ماتریس لاپلاسین جهت دار برای تعریف عبارت تغییرات کل گراف استفاده نمودیم. نتایج شبیه سازی نشان می دهند که روش GMCO-DL می تواند بیش از ۹۰ درصد درایه هایی را که به داده ی پرت آغشته شده اند، شناسایی نماید. این امر سبب می شود تا الگوریتم پیشنهادی، ماتریس مورد نظر را با خطای کمتری نسبت به روش GMCR بازیابی نماید.



شکل (۳): خطا بازیابی در معیار RMSE بر حسب درصد مشاهدات

پرت در زیر هر نمودار نشان داده شده است. نتایج شبیه سازی بیانگر عملکرد مطلوب روش پیشنهادی در شناسایی داده های پرت برای اطلاعات بیان شده در جدول ۱ می باشد. دقت شناسایی داده های پرت در این روش در حدود ۹۵ درصد می باشد. به عبارت دیگر روش GMCO-DL در شناسایی داده های پرت موفق تر عمل می کند. بدین ترتیب با اجرای متعدد الگوریتم برای داده های مختلف نتیجه می شود که GMCO-DL در برابر داده ی پرت مقاوم است. این بدان علت است که عبارت نرم l_1 ماتریس و عبارت تغییرات کل برای مدل سازی ارتباط بین درایه ها به مساله اضافه شده است. همچنین برای مقایسه میزان کارایی GMCO-DL در بازیابی ماتریس از معیار RMSE استفاده می نماییم.

$$RMSE = \sqrt{\frac{1}{N * L} \sum_{i=1}^N \sum_{j=1}^L (T_{ij} - X_{ij})^2} \quad (20)$$

میزان خطای بازیابی ماتریس از طریق روش GMCO-DL در شرایط مشاهدات آغشته شده به نویز و داده ی پرت در شکل ۳ قابل مشاهده است. کارایی این روش نسبت به روش GMCR سنجیده می شود. این شکل بر اساس اطلاعات جدول ۱ اقدام به بازیابی ماتریس می نمایند. کارایی روش پیشنهادی GMCO-DL در مقایسه با روش مرجع GMCR در شکل ۳ مقایسه شده است. وجود داده ی پرت در مشاهدات موجب شده است تا روش GMCR توانایی بازیابی ماتریس را نداشته باشد. به عبارت دیگر، حتی در هنگامی که درصد مشاهدات افزایش می یابد، امکان بازیابی ماتریس از روش GMCR امکان پذیر نیست. زیرا در تابع هدف و مدل ریاضی ماتریس مشاهده، داده پرت در نظر گرفته

مراجع

- [19] Sandryhaila, Aliaksei, and José MF Moura. "Discrete signal processing on graphs." *IEEE transactions on signal processing* 61.7 (2013): 1644-1656.
- [20] Sandryhaila, Aliaksei, and Jose MF Moura. "Big data analysis with signal processing on graphs: Representation and processing of massive data sets with irregular structure." *IEEE Signal Processing Magazine* 31.5 (2014): 80-90.
- [21] Candès, Emmanuel J., and Terence Tao. "The power of convex relaxation: Near-optimal matrix completion." *IEEE Transactions on Information Theory* 56.5 (2010): 2053-2080.
- [22] Singh, Rahul, Abhishek Chakraborty, and B. S. Manoj. "Graph Fourier transform based on directed Laplacian." *Signal Processing and Communications (SPCOM), 2016 International Conference on.* IEEE, 2016.
- [23] Li, Yongmou, et al. "A Graph-Based Method for Active Outlier Detection With Limited Expert Feedback." *IEEE Access* 7 (2019): 152267-152277.
- [24] Li, Ji, Jian-Feng Cai, and Hongkai Zhao. "Robust Inexact Alternating Optimization for Matrix Completion with Outliers." *Journal of Computational Mathematics* 38.2: 337-354. 2020.
- [25] Wang, Qianqian, et al. "Anomaly-Aware Network Traffic Estimation via Outlier-Robust Tensor Completion." *IEEE Transactions on Network and Service Management* 17.4: 2677-2689. 2020.
- [26] Tan, Teng, Lingwen Zhang, and Qiumei Li. "An Efficient Fingerprint Database Construction Approach Based on Matrix Completion for Indoor Localization." *IEEE Access* 8: 130708-130718. 2020.
- [۱] مودتی سمیرا. "تشخیص آریتمی‌های قلبی براساس تبدیل بسته موجک و الگوریتم فاکتورگیری ماتریس غیرمنفی تنک." *نشریه مهندسی برق و الکترونیک ایران* ۱۳۹۹؛ ۱۷ (۳): ۱۱۹-۱۲۸.
- [۲] مجیدیان سینا، کهائی محمدحسین. "تخمین هاپلوتاوپ با استفاده از فاکتورسازی ماتریس رتبه پایین در حضور داده پرت." *نشریه مهندسی برق و الکترونیک ایران* ۱۴۰۰؛ ۱۸ (۳): ۱۱۳-۱۲۰.
- [۳] مجیدیان سینا، کهائی محمدحسین. "تخمین هاپلوتاوپ با استفاده از ریلکس‌سازی بهینه‌سازی چندجمله‌ای." *نشریه مهندسی برق دانشگاه تبریز*. ۱۳۹۹، ۲، ۵۰، ۸۵۵-۸۶۳.
- [۴] مجیدیان سینا، حدادی فرزانه. "تخمین جهت منابع با استفاده از زیرفضای ختری-راو." *نشریه مهندسی برق و الکترونیک ایران*. ۱۳۹۶؛ ۱۴ (۲): ۳۷-۴۷.
- [5] Chatterjee, Sourav. "A deterministic theory of low rank matrix completion." *IEEE Transactions on Information Theory* 66.12 : 8046-8055..2020.
- [6] Candès, Emmanuel J., and Yaniv Plan. "Matrix completion with noise." *Proceedings of the IEEE* 98.6 (2010): 925-936.
- [7] Majidian, Sina, Mohamad Mahdi Mohades, and Mohammad Hossein Kahaei. "Matrix completion with weighted constraint for haplotype estimation." *Digital Signal Processing* 108 (2021): 102880.
- [8] Majidian, Sina, and Mohammad Hossein Kahaei. "NGS based haplotype assembly using matrix completion." *PLoS One* 14.3 (2019): e0214455.
- [9] Nie, Feiping, et al. "Robust Matrix Completion With Column Outliers." *IEEE Transactions on Cybernetics* (2021).
- [10] Akrami, Neda, Koorush Ziarati, and Soumyabrata Dev. "Graph-based local climate classification in Iran." *International Journal of Climatology* 42.3 (2022): 1337-1353.
- [11] Fathi, Hamid, Emad Rangriz, and Vahid Pourahmadi. "Two Novel Algorithms for Low-Rank Matrix Completion Problem." *IEEE Signal Processing Letters* 28 (2021): 892-896.
- [12] Ahmadi, Alireza, Sina Majidian, and Mohammad Hossein Kahaei. "Matrix Completion Using Graph Total Variation Based on Directed Laplacian Matrix." *Circuits, Systems, and Signal Processing* 40.6 (2021): 3099-3106.
- [13] Daei, Sajad, Farzan Haddadi, and Arash Amini. "Distribution-aware block-sparse recovery via convex optimization." *IEEE Signal Processing Letters* 26.4 (2019): 528-532.
- [14] Garg, Vaibhav, et al. "DOA Estimation via Shift-Invariant Matrix Completion." *Signal Processing*, 107993. 2021.
- [15] Chen, Yaru, et al. "A novel hierarchical deep matrix completion method." *IEEE Access*, 2021.
- [16] Kyriallidis, Anastasios, et al. "Provable compressed sensing quantum state tomography via non-convex methods." *Nature Quantum Information* 4.1 : 1-7. 2018.
- [17] Nguyen, Luong Trung, Junhan Kim, and Byonghyo Shim. "Low-rank matrix completion: A contemporary survey." *IEEE Access*. 7 (2019): 94215-94237.
- [18] Chen, Siheng, et al. "Signal recovery on graphs: Variation minimization." *IEEE Transactions on Signal Processing* 63.17 (2015): 4609-4624.

زیر نویس ها

¹ Low rank² Singular Values³ Machine learning⁴ Quantum state topography⁵ Outlier⁶ Adjacency matrix⁷ Non-deterministic Polynomial-time Hard (NP-Hard)⁸ Nuclear norm⁹ Graph signal Matrix Completion via total variation Minimization¹⁰ Graph signal Matrix Completion via total variation Regularization¹¹ Additive White Gaussian Noise