

## مقاوم سازی روش های یادگیری متریک در مقابل داده های پرت

داود ذبیح زاده خواجوی<sup>۱</sup>، رضا منصفی<sup>۲</sup>

<sup>۱</sup> استادیار گروه فناوری اطلاعات دانشکده مهندسی، دانشگاه فناوری نوین سبزوار، سبزوار، ایران، [d.zabihzadeh@gmail.com](mailto:d.zabihzadeh@gmail.com)

<sup>۲</sup> دانشیار گروه مهندسی کامپیوتر دانشکده مهندسی، دانشگاه فردوسی مشهد، مشهد، ایران، [monsefi@um.ac.ir](mailto:monsefi@um.ac.ir)

**چکیده:** در بسیاری از الگوریتم های یادگیری ماشین، شناسایی الگو و داده کاوی نیازمند آن هستیم که شباهت یا فاصله بین داده ها را به روش مناسب اندازه گیری کنیم. بعنوان مثال کارایی الگوریتم های خوشه بندی و یا طبقه بندی  $k$  نزدیک ترین همسایه به معیار فاصله/شباهت بستگی دارد. معیارهایی عمومی نظیر فاصله اقلیدسی و شباهت کسینوسی که بدون توجه به مفهوم داده ها میزان شباهت یا فاصله آن ها را مشخص می کنند، در بسیاری از کاربردها کارایی مناسبی ندارند. این مساله ضرورت یادگیری متریک را نشان می دهد. در یادگیری متریک هدف این است که با توجه به داده ها معیار شباهت یا فاصله بهینه به دست آید بطوری که داده هایی که از نظر مفهومی و منطقی شبیه به هم هستند، به یکدیگر نزدیک می شوند و داده هایی که از نظر مفهومی و منطقی شبیه نیستند از یکدیگر دور شوند. در این زمینه روش های زیادی ارائه شده است، اما همچنان یکی از چالش های مهم و جذاب، کاهش تأثیر داده پرت یا برجسب نویزی می باشد. در روش ارائه شده، مجموعه داده ورودی همزمان با یادگیری متریک، به دو بخش داده بدون خطا و بخش داده های پرت تقسیم می شود و یادگیری معیار فاصله تنها بر روی بخش بدون خطا انجام می شود. آزمایشات انجام شده بر روی داده های واقعی (در حضور و عدم حضور داده پرت و برجسب نویزی) کارایی الگوریتم ارائه شده را تایید می کند و برتری آن را نسبت به روش های همتا در مرزهای دانش در محیط های دارای نویز برجسب نشان می دهد.

**کلمات کلیدی:** یادگیری متریک، یادگیری متریک مقاوم، یادگیری متریک نزدیک ترین همسایه با حاشیه بزرگ، طبقه بندی  $k$  نزدیک ترین همسایه، داده های پرت، برجسب نویزی

### مقدمه

در این معیار فاصله  $M \in S_d^q$  (مخروط ماتریس های متقارن  $d \times d$  نیمه معین مثبت) است. این ویژگی  $M$  باعث می شود که متریک ماهالونوبیس چهار ویژگی نامفنی بودن، انطباق، تقارن و نامساوی مثلث را داشته باشد.

اکثر روش های یادگیری متریک فرضی بر وجود داده پرت و برجسب های نویزی قائل نشده اند. با این وجود آزمایشات مختلف تایید می کند که وجود این داده ها که باعث افت کارایی قابل ملاحظه آن ها می شود [۴]. از این رو کاهش تأثیر داده های پرت و برجسب های نویزی یکی از چالش های مهم در یادگیری متریک می باشد.

در روش ارائه شده، مجموعه داده ورودی همزمان با یادگیری متریک، به دو بخش داده بدون خطا و بخش داده های پرت تقسیم می شود و یادگیری معیار فاصله بر روی بخش بدون خطا انجام می شود. بطور دقیقتر، داده ورودی  $i$  ام  $\mathbf{x}_i \in \mathbb{R}^d$  به دو بخش  $\mathbf{x}_i^0$  (داده بدون خطا) و  $\mathbf{e}_i$  (داده پرت) به صورت  $\mathbf{x}_i = \mathbf{x}_i^0 + \mathbf{e}_i$   $i = 1, 2, \dots, n$  تجزیه می شود. برای اکثر داده ها، بخش  $\mathbf{x}_i^0$  مقدار دارد و بخش  $\mathbf{e}_i$  برابر با  $\mathbf{0}$  است، بنابراین ماتریس  $E = [\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n]$  تا حد ممکن باید تنگ باشد. روش پیشنهادی برخلاف روش های معمول، یادگیری متریک را بر روی ماتریس  $\mathbf{X}_0 = [\mathbf{x}_1^0, \mathbf{x}_2^0, \dots, \mathbf{x}_n^0]$  که حاوی داده های بدون نویز است، انجام می دهد. بدین ترتیب می توان انتظار داشت که متریک یادگرفته شده در مقابل نویز برجسب و داده های پرت مقاوم خواهد بود. نتیجه آزمایشات بر روی مجموعه داده های مختلف حاوی نویز برجسب، صحت این ادعا را تایید می کند.

ادامه مقاله بصورت زیر سازماندهی شده است. در بخش ۲ کارهای انجام شده در زمینه مقاوم سازی یادگیری متریک مرور می شوند. بخش ۳ به روش پیشنهادی می پردازد. بخش ۴ به آزمایشات انجام شده برای بررسی کارایی روش پیشنهادی و تحلیل نتایج اختصاص

یادگیری متریک<sup>۱</sup> یکی از مسائل بنیادی در برنامه های کاربردی می باشد. عملکرد بسیاری از الگوریتم های یادگیری ماشین از جمله  $k$  نزدیک ترین همسایه<sup>۲</sup> (kNN) [۱] و خوشه بندی<sup>۳</sup> به معیار فاصله ای که برای سنجش ارتباط بین داده های ورودی است، بستگی دارد [۲]. هدف روش های یادگیری متریک، نزدیک کردن داده های مشابه به یکدیگر و افزایش فاصله بین داده های غیرمشابه می باشد، این امر باعث بهبود عملکرد الگوریتم یادگیری مبتنی بر متریک می شود.

اغلب روش های یادگیری متریک از داده های آموزشی بشکل محدودیت های زوج و یا سه گانه استفاده می کنند، که بصورت زیر تعریف می شود:

$$S = \{(\mathbf{x}_i, \mathbf{x}_j) \mid \mathbf{x}_i \text{ and } \mathbf{x}_j \text{ should be similar}\}$$

$$D = \{(\mathbf{x}_i, \mathbf{x}_j) \mid \mathbf{x}_i \text{ and } \mathbf{x}_j \text{ should be dissimilar}\}$$

$$T = \{(\mathbf{x}_i, \mathbf{x}_j, \mathbf{x}_k) \mid \mathbf{x}_i \text{ should be more closer to } \mathbf{x}_j \text{ than to } \mathbf{x}_k\}$$

به طور کلی این روش های، مسئله بهینه سازی به صورت زیر را حل می کنند [۳]:

$$\underset{M}{\text{minimize}} \quad l(M, C) + \lambda r(M) \quad (1)$$

در رابطه بالا  $M$  نشان دهنده متریک (مجموعه پارامترهای یادگیری) است.  $C$  مجموعه محدودیت های مساله را نشان می دهد که اغلب بصورت محدودیت های زوج و سه گانه است.  $l$  تابع ضرر و  $r$  تابع تنظیم کننده هست که برای جلوگیری از بیش برآش<sup>۴</sup> تعریف شده است. یکی از مهم ترین روش های یادگیری متریک که تاکنون بیشتر توجه را به خود اختصاص داده است یادگیری فاصله ماهالونوبیس است که به صورت زیر تعریف می شود:

$$d_M(\mathbf{x}_i, \mathbf{x}_j)^2 = (\mathbf{x}_i - \mathbf{x}_j)^T M (\mathbf{x}_i - \mathbf{x}_j) \quad (2)$$

<sup>۴</sup> Over Fitting

<sup>۱</sup> Metric Learning

<sup>۲</sup> k-nearest neighbor classifier

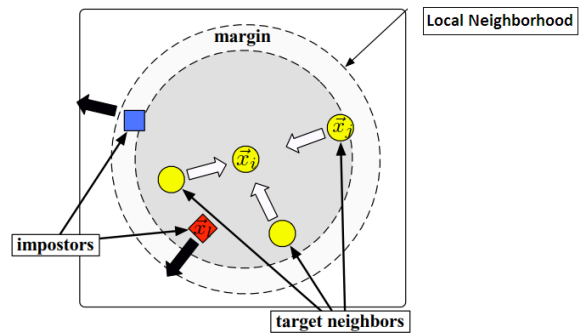
<sup>۳</sup> Clustering

یافته است. سرانجام، جمع بندی و نتیجه گیری به همراه چند پیشنهاد برای کارهای آتی در بخش ۵ ارائه شده است.

### مروری بر کارهای مرتبط

اولین کار در زمینه یادگیری فاصله ماهالونوبیس، الگوریتم MMC است [۵]. در این روش هدف ماکزیمم کردن فاصله بین داده‌های نامشابه است درحالی که مجموع فاصله داده‌های مشابه کمتر از یک مقدار معین باشد.

از کارهای مطرح در زمینه یادگیری متریک به روش LMNN می‌توان اشاره کرد که برای بهبود طبقه بند kNN طراحی شده است. با توجه به اینکه روش پیشنهادی بر روی این الگوریتم بمنظور مقاوم سازی آن اعمال شده است، این روش را با جزئیات بیشتر بررسی می‌کنیم. برای هر داده  $x_i$ ،  $k$  نزدیکترین همسایه به آن با برچسب یکسان را همسایه های هدف (Target Neighbors) آن می‌نامیم. همچنین به داده هایی که در محدوده فاصله  $x_i$  با دورترین همسایه هدف آن با در نظر گرفتن یک حاشیه (margin) مناسب قرار می‌گیرند، داده های بدل (impostors) گفته می‌شود. در LMNN، هدف کاهش فاصله هر داده با همسایه‌های هدف آن و درعین حال افزایش فاصله آن با داده های بدل است بطوری که داده های بدل از محدوده همسایگی تعریف شده توسط همسایه های هدف خارج شوند. در شکل زیر مفاهیم همسایه های هدف، همسایه های بدل و حاشیه نشان داده شده است.



شکل ۱- مفهوم همسایه های هدف و بدل داده  $x_i$  به‌همراه همسایگی تعریف شده و حاشیه

به‌صورت کلی تابع هزینه LMNN از دو بخش اصلی بصورت زیر تشکیل می‌شود:

$$\varepsilon(M) = (1 - \mu)\varepsilon_{pull}(M) + \mu\varepsilon_{push}(M) \quad (3)$$

$\varepsilon_{pull}(M)$  برای نزدیک کردن هر داده به همسایه های هدف آن، بصورت زیر تعریف می‌شود:

$$\varepsilon_{pull}(M) = \sum_{j \sim i} d_M(x_i, x_j)^2 \quad (4)$$

در رابطه بالا،  $j \sim i$  همسایه‌های هدف  $x_i$  را نشان می‌دهند. همچنین برای بیرون راندن داده های بدل از حاشیه، تابع زیان زیر تعریف شده است:

$$\varepsilon_{push}(M) = \sum_{i, j \in \mathcal{L}} \sum_l (1 - y_{il}) [1 + d_M(x_i, x_j)^2 - d_M(x_i, x_l)^2]_+ \quad (5)$$

در رابطه بالا  $y_{il} = 1$  اگر  $y_i = y_l$  باشد در غیراینصورت برابر 0 است. با ترکیب روابط (4) و (5) به مسئله بهینه‌سازی زیر می‌رسیم.

$$\begin{aligned} \text{minimize} & (1 - \mu) \sum_{i, j \in \mathcal{L}} d_M(x_i, x_j)^2 \\ & + \mu \sum_{i, j \in \mathcal{L}} \sum_l (1 - y_{il}) \xi_{ijl} \end{aligned} \quad (6)$$

subject to

$$1 + d_M(x_i, x_j)^2 - d_M(x_i, x_l)^2 \leq \xi_{ijl} \quad \xi_{ijl} \geq 0 \quad M \succcurlyeq 0$$

در این رابطه  $0 \leq \mu \leq 1$  امکان مصالحه بین دو جمله را فراهم می‌کند. این الگوریتم در عمل خوب کار می‌کند اما به دلیل عدم وجود تابع تنظیم‌کننده در برخی مسائل مشکل بیش‌برازش دارد.

بسیاری از روش‌های یادگیری متریک فرضی بر وجود داده‌های پرت و همچنین برچسب نویزی ندارند. زمانی که داده‌ها شامل ویژگی‌های نامربوط یا نویز برچسب باشند، این روش‌ها با مشکل مواجه می‌شوند. در ادامه این بخش به چند روش که در سال‌های اخیر برای مقاوم سازی متریک در برابر داده های پرت و نویز برچسب انجام شده است، می‌پردازیم.

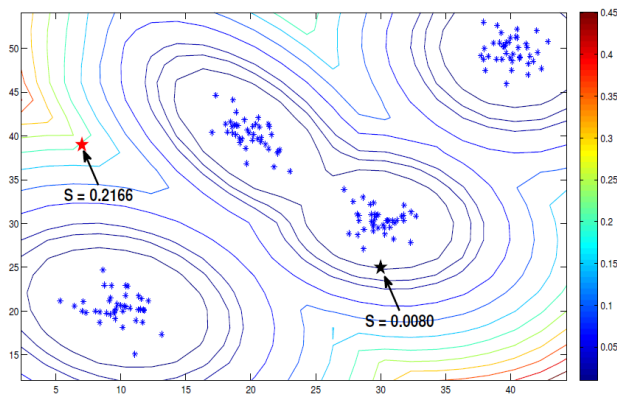
در مقاله [۶] برای بهبود عملکرد FCM<sup>۵</sup> نرم دوم اقلیدسی با تابع فاصله نرمالیزه شده  $1 - \exp(-\beta \|x_i - v_c\|^2)$  جایگزین شده است. متریک فاصله که در روش FCM استفاده می‌شود فقط فاصله اقلیدسی بین هر داده و مرکز خوشه را در نظر می‌گیرد و از پراکندگی فاصله بین داده‌ها در یک خوشه صرف‌نظر می‌کند. این روش تأثیر داده‌های دورافتاده را کاهش می‌دهد ولی از خود داده‌ها و میزان پراکندگی آن‌ها استفاده نمی‌کند و زمانی که خوشه‌های ناهمگن داشته باشیم، نتایج مطلوبی نمی‌دهد.

در برخی از روش‌ها برای مقاوم سازی متریک بر روی کاهش رتبه ماتریس  $M$  تأکید شده است. از جمله این روش‌ها می‌توان به  $\gamma$  اشاره کرد که از تابع trace بر روی ماتریس  $M$  برای تُنک کردن خروجی و از تُرم ترکیبی یک و دو برای تُنک کردن ورودی استفاده می‌کند. مقاله [۸] هم که توسعه روش LMNN است، از تُرم ترکیبی یک و دو ماتریس  $M$  برای جلوگیری از بیش‌برازش داده‌ها و همچنین تُنک کردن آن بهره می‌برد.

[۹] از  $KL^\gamma$  و همچنین احتمالات پیشین و پسین برای تشخیص داده های پرت استفاده می‌کند. بدین ترتیب که اگر تفاوت بین توزیع پیشین و پسین بر اساس داده ورودی  $x_n$  و مجموعه مدل‌های  $M \in \mathcal{M}$  از یک حد آستانه ( $T_{out}$ ) بیشتر باشد آن داده را داده دورافتاده در نظر می‌گیرد. رابطه (7) و شکل ۹ بیانگر این موضوع می‌باشد. در این رابطه  $P(M|x_n)$  نشان‌دهنده احتمال مدل  $M$  با توجه به داده ورودی  $x_n$  می‌باشد.

$$S(x_n, \mathcal{M}) = H_m(x_n) = KL(P(M|x_n), P(M)) \quad (7)$$

$$= \int_{M \in \mathcal{M}} P(M|x_n) \log \frac{P(M|x_n)}{P(M)} dM$$



شکل ۲- داده قرمز رنگ، نشان‌دهنده داده پرت و داده مشکلی، نشان‌دهنده داده مناسب می‌باشد که با استفاده از KL تشخیص داده شده‌اند [۹].

در مقاله [۱۰] برخلاف روش‌های دیگر که فرض را بر درستی انتخاب مجموعه‌های مشابه  $S$  و مجموعه‌های غیرمشابه  $D$  می‌گذارند، این احتمال را در نظر می‌گیرد که احتمال نادرست بودن مجموعه‌ها وجود دارد. در این روش به‌نوعی وزن دهی به قیدها انجام شده است و در نهایت تبدیل به مسئله محدب شده و با بهینه‌سازی هموار شده نیستروف حل می‌شود.

[۱۱] از ایده استفاده از داده‌های کمکی و بدون برچسب و ترکیب کردن آن با داده‌های آموزش بمنظور جلوگیری از بیش‌برازش بهره می‌برد. داده‌های کمکی به داده‌هایی گفته می‌شود که از منبع و تابع توزیع متفاوتی نسبت به هدف به دست می‌آید. در [۱۲] نشان داده شده است که استفاده از این داده‌ها در یادگیری متریک مفید است. روش کار، ساختن

<sup>۵</sup> Kullback-Leibler divergence

<sup>Δ</sup> Fuzzy C-Means Clustering

<sup>ε</sup> rank

متریک‌های دسته‌جمعی<sup>۸</sup> از داده‌های محدود آموزش و داده‌های کمکی می‌باشد. سپس متریک تجمیع شده<sup>۹</sup> از ترکیب داده‌های آموزش و کمکی بدست می‌آید.

در [۱۴] فرایند شناسایی مقاوم به دو فاز تشخیص داده‌های پرت و فاز شناسایی تقسیم شده است. بخش اول از رگرسیون خطی وزن‌دار برای یادگیری متریکی که داده‌های پرت و نویز را تشخیص می‌دهد، استفاده می‌کند. در بخش دوم با استفاده از متریک یادگرفته شده، مجموعه داده بزرگ به مجموعه کوچک با استفاده از معیار نزدیکترین همسایه فیلتر می‌شود، سپس یک بازنمایی تنک<sup>۱۰</sup> با استفاده از تکنیک حداقل مربعات نامنفی<sup>۱۱</sup> محاسبه می‌شود.

در [۱۴] برای یادگیری متریک مقاوم در برابر نویز، تابعی هدفی مشابه [۱۵] ارائه شده است که بجای  $\mu$  از  $\mu$  استفاده می‌کند. با این وجود بعلاوه نامحدود بودن مساله بهینه سازی، جواب مساله پیشنهاد شده به مقداره‌ی اولیه متریک بسیار حساس است. همچنین بدلیل اینکه در تابع  $\mu$  ضرر بصورت خطی رشد می‌کند، این روش چندان برای مجموعه داده‌های حاوی نویز برجسب و داده‌های پرت مناسب نیست.

در [۱۶] روش معروف NCA<sup>۱۲</sup> در محیط‌های حاوی نویز برجسب مقاوم شده است. برای این منظور ابتدا تاثیر برجسب نویزی بر روی مشتق تابع درست‌نمایی (likelihood) نشان داده شده و این روش سعی می‌کند با مدل سازی تابع احتمال درست بودن برجسب، تاثیر داده‌هایی که با احتمال بالا نویزی هستند را کاهش دهد.

در [۱۷] یادگیری متریک با استفاده از استنتاج بیزی برای مقاوم سازی ارائه شده است. دلیل این مطلب حساسیت بیشتر ماتریس فاصله در روش‌های تخمین نقطه‌ای به مثال‌های آموزشی عنوان شده است. این روش در واقع توسعه بیزی الگوریتم معروف LMNN است و باعث مقاوم‌پذیری بهتر آن در برابر نویز برجسب شده است.

### روش پیشنهادی

اگر داده‌های ورودی را به صورت  $\mathbf{X} \in \mathbb{R}^{d \times n}$  که  $n$  تعداد داده‌های ورودی و  $d$  ابعاد آن می‌باشد، در نظر بگیریم. می‌توان داده‌های ورودی  $\mathbf{X}$  را به دو بخش  $\mathbf{X}_0, \mathbf{E} \in \mathbb{R}^{d \times n}$  طبق رابطه (9) تقسیم کرد که  $\mathbf{X}_0$  را بخش داده‌های بدون خطا و  $\mathbf{E}$  را بخش مقادیر پرت یا خطا در نظر گرفت.

$$\mathbf{X} = \mathbf{X}_0 + \mathbf{E} \quad (8)$$

بدلیل این‌که بیشتر داده‌ها حاوی نویز برجسب نیستند، ماتریس  $\mathbf{E}$  باید تنک باشد. برای تنک سازی این ماتریس می‌توان  $\mu$  یک آن را مینیمم کرد که بصورت زیر تعریف می‌شود:

$$\|\mathbf{E}\|_{1,1} = \|\text{vec}(\mathbf{E})\|_1 = \sum_{i=1}^d \sum_{j=1}^n |e_{ij}| \quad (9)$$

با توجه به رابطه (1) مسئله بهینه‌سازی متریک معمولاً از دو بخش تابع ضرر که شامل ماتریس متریک  $\mathbf{M}$  و مجموعه محدودیت‌های  $\mathbf{C}$  است، و جمله تنظیم‌کننده برای جلوگیری از بیش‌برازش تشکیل شده است. در روش پیشنهادی برای آموزش ماتریس متریک  $\mathbf{M}$  داده‌های خالص  $\mathbf{X}_0 \in \mathbb{R}^{d \times n}$  جایگزین داده‌های ورودی  $\mathbf{X} \in \mathbb{R}^{d \times n}$  می‌شود. همچنین تنک سازی ماتریس  $\mathbf{E}$  را بعنوان جمله تنظیم‌کننده در نظر میگیریم. به این ترتیب به مساله بهینه‌سازی زیر می‌رسیم:

$$\min_{\mathbf{M} \succeq \mathbf{0}, \mathbf{X}_0} l(\mathbf{M}, \mathbf{C}, \mathbf{X}_0) + \lambda r(\mathbf{M}, \mathbf{E}) \quad (10)$$

به دلیل آنکه می‌توان ماتریس مقادیر پرت را طبق رابطه (8) به صورت  $\mathbf{E} = \mathbf{X} - \mathbf{X}_0$  نوشت، بنابراین با محاسبه مقدار ماتریس  $\mathbf{X}_0$  مقدار ماتریس  $\mathbf{E}$  محاسبه می‌شود و این متغیر را به‌عنوان مجهول در نظر نمی‌گیریم. ماتریس متریک  $\mathbf{M}$  نیمه معین مثبت متقارن است و می‌توان آن را به صورت  $\mathbf{M} = \mathbf{W}\mathbf{W}^T$  تجزیه کرد که  $\mathbf{W} \in \mathbb{R}^{d \times r}$  و  $r$  رتبه ماتریس  $\mathbf{M}$  را نشان می‌دهد. همچنین طبق رابطه (2) و با استفاده از انتقال خطی  $\mathbf{x}'_i = \mathbf{W}^T \mathbf{x}_i$ ، داده‌ها در فضای جدید به دست می‌آیند. رابطه (10) بیانگر این موضوع است که برای به دست آوردن ماتریس متریک  $\mathbf{M}$  و درنهایت ماتریس انتقال خطی  $\mathbf{W}$  از مقادیر بدون خطا  $\mathbf{X}_0$  استفاده می‌شود. نتایج آزمایشات تایید می‌کند که این

کار باعث بهبود قابل ملاحظه در کارایی متریک یادگرفته شده می‌شود. مسئله پیشنهادی نسبت به متغیرهای  $\mathbf{X}_0$  و  $\mathbf{M}$  نامحدوب می‌باشد و برای حل مسئله ارائه شده از چارچوب BCD<sup>۱۳</sup> معرفی شده استفاده شده است. این روش متغیرها را به  $\mu$  زیرگروه مجزا  $\chi = \square(\mathbf{M}, \mathbf{X}_0)$  تقسیم می‌کند بطوریکه هر زیر مسئله بطور مجزا محذب است. مساله اصلی به صورت تکراری و با در نظر گرفتن تنها متغیرهای هر گروه در یک‌زمان، حل می‌شود.

### اعمال روش پیشنهادی بر روی الگوریتم LMNN

با اعمال روش پیشنهادی بر روی الگوریتم LMNN به مساله بهینه سازی زیر می‌رسیم:

$$\begin{aligned} \text{minimize} \quad & (1 - \mu) \sum_{i,t \sim j} d_{\mathbf{M}}(\mathbf{x}_{0i}, \mathbf{x}_{0j})^2 \\ & + \mu \sum_{i,j \sim i} \sum_l (1 - y_{il}) \xi_{ijl} + \lambda_e \|\mathbf{E}\|_{1,1} \end{aligned} \quad (11)$$

subject to

$$\begin{aligned} 1 + d_{\mathbf{M}}(\mathbf{x}_{0i}, \mathbf{x}_{0j})^2 - d_{\mathbf{M}}(\mathbf{x}_{0i}, \mathbf{x}_{0l})^2 & \leq \xi_{ijl} \\ \mathbf{X} = \mathbf{X}_0 + \mathbf{E}, \xi_{ijl} \geq 0 \quad \mathbf{M} \succeq \mathbf{0} \end{aligned}$$

همانطور که مشاهده می‌شود در روش پیشنهادی برای یادگیری متریک از بخش بدون خطای داده ( $\mathbf{X}_0$ ) بهره می‌برد. همچنین ابر پارامتر  $\lambda_e$  برای تنظیم میزان اهمیت تنک بودن ماتریس  $\mathbf{E}$  به مساله اضافه شده است. برای حل رابطه نامحدوب (11) با بهره گرفتن از روش BCD می‌توان متغیرها  $\theta\chi = (\mathbf{M}, \mathbf{X}_0)$  را به دو زیر مسئله محذب تقسیم کرد و به صورت تکراری حل نمود. اثبات شده است که در هر بروزسانی، تابع هزینه هیچ‌وقت افزایش نمی‌یابد [۱۸]. تابع هزینه (11) به صورت بلاکی محذب می‌باشد در نتیجه در هر تکرار تمامی زیر مسائل محذب می‌باشند [۱۹]. مقادیر اولیه  $\mu$  و  $\lambda_e$  با اعتبار سنجی متقاطع و  $\mathbf{X}_0$  به دست می‌آید. گام‌های زیر به ترتیب برای حل تابع (11) انجام می‌پذیرد.

۱- در مرحله  $k$ ، ماتریس نیمه مثبت معین  $\mathbf{M}_s$  با استفاده از  $\mathbf{X}_{0,s-1}$  محاسبه می‌شود.

۲- محاسبه ماتریس داده‌های بدون خطا  $\mathbf{X}_{0,s}$  با ثابت در نظر گرفتن ماتریس  $\mathbf{M}_s$

۳- تکرار مراحل ۱ و ۲ تا زمانی که مساله به مقدار بهینه همگرا شود یا به تعداد تکرار موردنظر برسد.

### یادگیری ماتریس متریک $\mathbf{M}_s$

در این مرحله ماتریس  $\mathbf{X}_{0,s}$  ثابت فرض می‌شود و تنها متغیر  $\mathbf{M}$  را لحاظ می‌کنیم. مشابه الگوریتم LMNN می‌توان این مرحله را با استفاده از روش زیر گردایان نزولی بشکل زیر حل کرد:

$$\mathbf{G}_t = \frac{\partial \varepsilon}{\partial \mathbf{M}_t} = (1 - \mu) \sum_{i,j \sim i} \mathbf{C}_{ij} + \mu \sum_{(i,j,l) \in \mathcal{N}_t} (\mathbf{C}_{ij} - \mathbf{C}_{il}) \quad (12)$$

$$\mathbf{M}_{t+1} = \mathbf{M}_t - \tau \mathbf{G}_t \quad (13)$$

در رابطه بالا  $\mathbf{C}_{ij} = (\mathbf{x}_{0i,s-1} - \mathbf{x}_{0j,t-1})(\mathbf{x}_{0i,s-1} - \mathbf{x}_{0j,s-1})^T$  است و  $\tau$  نشان دهنده نرخ یادگیری یا گام حرکت است.

ماتریس متریک  $\mathbf{M}_t$  باید نیمه مثبت معین باقی بماند. برای محقق شدن این امر، تجزیه مقادیر ویژه  $\mathbf{M}_t = \mathbf{V}\mathbf{\Delta}\mathbf{V}^T$  را به دست می‌آوریم.  $\mathbf{V}$  ماتریس بردارهای ویژه می‌باشد و  $\mathbf{\Delta}$  ماتریس قطری مقادیر ویژه که  $\mathbf{\Delta} = \mathbf{\Delta}^- + \mathbf{\Delta}^+$  می‌باشد. با صفر در نظر گرفتن مقادیر ویژه منفی می‌توان تصویر  $\mathbf{M}_t$  را بر روی منحنی نیمه مثبت معین به صورت زیر به دست آورد:

$$\rho_s(\mathbf{M}_t) = \mathbf{V}\mathbf{\Delta}^+\mathbf{V}^T \quad (14)$$

### یادگیری داده‌های بدون خطا $\mathbf{X}_{0,s}$

فرض کنید در تکرار  $s$  قرار داریم. در مرحله قبل ماتریس متریک  $\mathbf{M}_s$  محاسبه شد و با ثابت در نظر گرفتن آن می‌توان ماتریس داده‌های بدون خطا  $\mathbf{X}_{0,s}$  را به دست آورد. ماتریس

<sup>۱۳</sup> Neighborhood Components Analysis

<sup>۱۲</sup> Block-Coordinate Descent

<sup>۸</sup> ensemble

<sup>۹</sup> aggregated

<sup>۱۰</sup> Sparse Representation

<sup>۱۱</sup> Non-negative Least Squares

که در رابطه بالا  $\alpha$  نرخ یادگیری را نشان می دهد. در الگوریتم ۱ مراحل روش پیشنهادی بطور خلاصه نشان داده شده است.

#### الگوریتم ۱- شبه کد روش پیشنهادی

<p><b>ورودی:</b> مجموعه آموزشی شامل <math>n</math> نمونه <math>\mathbf{x}_i \in \mathbb{R}^d</math> و برچسب‌های آن در حالت اولیه داریم: <math>\mathbf{x}_{0i} = \mathbf{x}_i, \forall i = 1, \dots, n</math> <b>خروجی:</b> ماتریس نیمه معین مثبت <math>\mathbf{M}</math>.</p>
<p><b>الگوریتم:</b></p> <ol style="list-style-type: none"> <li>۱- انتخاب <math>\lambda_e</math> و <math>\mu</math> مناسب.</li> <li>۲- محاسبه <math>\mathbf{G}_t</math> و با استفاده از رابطه (12)</li> <li>۳- محاسبه <math>\mathbf{M}_t</math> از رابطه (14) و ثابت بودن ماتریس <math>\mathbf{X}_0</math> در تمامی تکرارها</li> <li>۴- تکرار مراحل ۲ و ۳ تا زمانی که به کمینه سراسری نسبت به <math>\mathbf{M}_S</math> همگرا شود یا <math>t</math> به حداکثر تکرار برسد.</li> <li>۵- محاسبه <math>\frac{\partial \epsilon(\mathbf{x}_{0ij})}{\partial \mathbf{x}_{0ij}}</math> از رابطه (18) و <math>\mathbf{X}_{0,t}</math> از رابطه (19) با ثابت بودن ماتریس <math>\mathbf{M}_S</math></li> <li>۶- تکرار مرحله ۵ تا زمانی که به کمینه سراسری نسبت به <math>\mathbf{X}_S</math> همگرا شود یا <math>t</math> به حداکثر تکرار برسد.</li> <li>۷- تکرار مراحل ۲ تا ۶ و بررسی شرایط همگرایی و برگشت به گام ۲ در صورت همگرا نشدن</li> </ol>

#### آزمایش‌ها و تحلیل نتایج

جهت ارزیابی کارایی الگوریتم ارائه شده و مقایسه آن با سایر روش‌ها از ۸ پایگاه داده با اندازه و پیچیدگی‌های مختلف استفاده می‌کنیم. مشخصات این مجموعه داده‌ها در جدول گزارش شده است.

در پایگاه داده‌هایی که تعداد ویژگی‌ها زیاد می‌باشد، با استفاده از PCA ابعاد داده‌ها را قبل از اعمال الگوریتم کاهش می‌دهیم. نتایج از میانگین چندین بار اجرا با تخصیص ۷۰٪ از داده‌ها به داده‌های آموزش و ۳۰٪ داده‌های تست، به دست می‌آید. تعداد همسایه‌های هدف برای یادگیری الگوریتم پیشنهادی و LMNN در تمامی گزارش‌ها  $k=3$  و پارامتر  $\mu=0.5$  می‌باشد. همسایه‌های هدف در فضای ورودی (در صورت نیاز پس از کاهش ابعاد با PCA) بر اساس فاصله اقلیدسی به دست می‌آید. مقدار مناسب پارامتر  $\lambda_e$  برای هر پایگاه داده با استفاده از روش اعتبار سنجی متقاطع 5-fold حاصل می‌شود. طبقه بند KNN از فاصله اقلیدسی و LMNN از متریک حاصل شده و روش ارائه شده از داده‌های بدون خطا  $\mathbf{X}_0$  داده‌های آموزش و ماتریس متریک برای مقایسه استفاده می‌کند. همچنین از روش RNCA برای مقایسه با روش ارائه شده استفاده شده است. نتایج گزارش شده در این بخش بر اساس درصد خطای طبقه‌بندی است. طبقه بند استفاده شده، KNN با  $k=1$  می‌باشد. نویز برچسب با نرخ‌های متفاوت و همچنین داده پرت برای ارزیابی روش‌های مورد مقایسه، استفاده شده است.

$\mathbf{M}_S$  ثابت است، بنابراین قید مثبت بودن آن از رابطه (11) حذف می‌گردد و می‌توان تابع هزینه را بشکل ساده تر زیر نوشت:

$$\epsilon(\mathbf{X}_0) = L(\mathbf{X}_{0,S}) + \lambda_e \|\mathbf{X} - \mathbf{X}_0\|_1 \quad (15)$$

$$L(\mathbf{X}_0) = (1 - \mu) \sum_{i,j \sim i} d_M(\mathbf{x}_{0i}, \mathbf{x}_{0j})^2 + \mu \sum_{i,j \sim i} \sum_l (1 - y_{il}) \left[ 1 + d_M(\mathbf{x}_{0i}, \mathbf{x}_{0j})^2 - d_M(\mathbf{x}_{0i}, \mathbf{x}_{0l}) \right]_+ \quad (16)$$

اگر در هر مرحله برای محاسبه بردار گرادینان، مجموعه سه تایی  $\mathcal{N}_k = (i, j, l)$  که شامل قیده‌های فعال می‌باشند، را در نظر بگیریم آن‌گاه  $[Z]_+ = [Z]$  و در نتیجه بردار گرادینان  $\nabla_i L(\mathbf{x}_{0i})$  را می‌توان با مشتق‌گیری بصورت زیر بدست آورد:

$$\nabla_i L(\mathbf{x}_{0i}) = (1 - \mu) \left( \sum_{j \sim i} 2\mathbf{M}\mathbf{x}_{0i} - 2\mathbf{M}\mathbf{x}_{0j} + \sum_{l \sim j} -2\mathbf{M}\mathbf{x}_{0j} + 2\mathbf{M}\mathbf{x}_{0i} \right) + \mu \left( \sum_{(j \sim i) \in \mathcal{N}} 2\mathbf{M}\mathbf{x}_{0i} - 2\mathbf{M}\mathbf{x}_{0j} + \sum_{(l \sim j) \in \mathcal{N}} -2\mathbf{M}\mathbf{x}_{0j} + 2\mathbf{M}\mathbf{x}_{0i} + \sum_{(i \rightarrow l) \in \mathcal{N}} 2\mathbf{M}\mathbf{x}_{0i} - 2\mathbf{M}\mathbf{x}_{0l} + \sum_{(l \rightarrow i) \in \mathcal{N}} -2\mathbf{M}\mathbf{x}_{0l} + 2\mathbf{M}\mathbf{x}_{0i} \right) \quad (17)$$

در رابطه بالا،  $l \rightarrow i$  بیان‌کننده بدل بودن داده  $l$  برای  $i$  می‌باشد. همچنین بردار گرادینان  $\nabla_i \in (\mathbf{x}_{0i})$  را می‌توان از رابطه زیر محاسبه کرد:

$$\frac{\partial \epsilon(\mathbf{x}_{0ij})}{\partial \mathbf{x}_{0ij}} = \begin{cases} \frac{\partial L(\mathbf{x}_{0ij})}{\partial \mathbf{x}_{0ij}} + \lambda_e \text{sign}(x_{i,j} - x_{0i,j}), & |x_{i,j} - x_{0i,j}| > 0 \\ \frac{\partial L(\mathbf{x}_{0ij})}{\partial \mathbf{x}_{0ij}} + \lambda_e, & x_{i,j} - x_{0i,j} = 0, \frac{\partial L(\mathbf{x}_{0ij})}{\partial \mathbf{x}_{0ij}} < -\lambda_e \\ \frac{\partial L(\mathbf{x}_{0ij})}{\partial \mathbf{x}_{0ij}} - \lambda_e, & x_{i,j} - x_{0i,j} = 0, \frac{\partial L(\mathbf{x}_{0ij})}{\partial \mathbf{x}_{0ij}} > \lambda_e \\ 0, & x_{i,j} - x_{0i,j} = 0, -\lambda_e \leq \frac{\partial L(\mathbf{x}_{0ij})}{\partial \mathbf{x}_{0ij}} \leq \lambda_e \end{cases} \quad (18)$$

در حالت اولیه  $\mathbf{x}_{ij} = \mathbf{x}_{0ij}$  است و در رابطه (18) فقط  $\mathbf{x}_{0ij}$  که  $\left| \frac{\partial L(\mathbf{x}_{0ij})}{\partial \mathbf{x}_{0ij}} \right| > \lambda_e$  دارند، می‌تواند مقادیر غیر صفر داشته باشد. بنابراین ماتریس خطا ( $\mathbf{E}$ ) تنگ می‌باشد. درنهایت با توجه به گرادینان محاسبه شده در رابطه، ماتریس  $\mathbf{X}_{0,t}$  با روش نزول گرادینان از رابطه زیر به دست می‌آید:

$$\mathbf{X}_{0,t} = \left( \mathbf{X}_{0,t-1} - \alpha \nabla \epsilon(\mathbf{X}_{0,t}) \right) \quad (19)$$

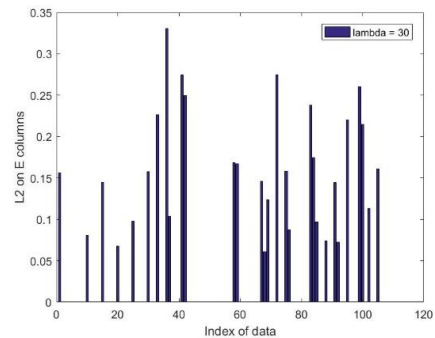
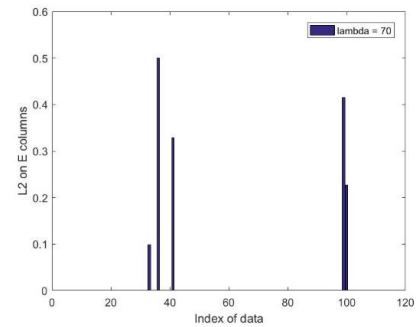
جدول ۱- ویژگی‌های مجموعه داده‌های مورداستفاده در آزمایش‌ها.

مجموعه داده	تعداد داده‌ها (آزمون / آموزش)	تعداد ویژگی‌ها	تعداد ویژگی‌های کاهش یافته	تعداد کلاس‌ها
ionosphere	۲۴۵/۱۰۷	۳۴	۳۴	۲
Heart	۱۸۹/۸۲	۱۳	۱۳	۲
balance	۴۳۷/۱۸۹	۴	۴	۳
wine	۱۰۶/۴۶	۱۳	۱۳	۳
iris	۱۰۵/۴۵	۴	۴	۳
Olivetti face	۲۸۰/۱۲۰	۲۰۰	۲۰۰	۴۰
MNIST	۲۱۰۰/۹۰۰	۷۸۴	۱۶۴	۱۰

### نتایج آزمایش‌ها

#### تأثیر پارامتر $\lambda$ بر تنگ بودن ماتریس خطا

در حالت اولیه  $X = X_0$  می‌باشد و طبق رابطه (18) مقادیر گرادیانی که از قدر مطلق  $\lambda_e$  بیشتر باشد، می‌تواند در جهت خلاف بردار گرادیان پیش برود. به عبارتی با تنظیم دقیق این پارامتر، بخش‌هایی از داده  $X_i$  که شامل مقادیر پرت است به  $e_i$  منتقل می‌شود. شکل ۱۹ تأثیر اندازه پارامتر  $\lambda$  بر انتخاب بخش‌هایی که شامل مقادیر پرت می‌باشند در مجموعه داده Olivetti face را نشان می‌دهد. در این شکل از نرم  $l_2$  روی هر ستون  $e_i \forall i = 1, 2, \dots, n$  برای نمایش آن استفاده شده است. نتایج این آزمایش تایید می‌کند اندازه پارامتر  $\lambda$  در تعیین تنگ بودن ماتریس  $E$  نقش مهمی دارد. هر چه  $\lambda$  مقدار بیشتری داشته باشد، میزان تنگ‌سازی بیشتر خواهد بود. اگر  $\lambda$  به صورتی اختیار شود که ماتریس  $e_{ij} = 0 \forall i = 1, \dots, n, j = 1, \dots, d$  یا به عبارتی تمامی مقادیر آن برابر صفر باشد، مدل پیشنهادی مشابه روش LMNN می‌باشد. با انتخاب  $\lambda$  مناسب می‌توان از بیش برآزش بر روی داده‌های آموزش نیز جلوگیری کرد.



شکل ۳- تأثیر اندازه پارامتر  $\lambda$  در تنگ بودن ماتریس خطای  $E$ .

### آزمایش بر روی مجموعه داده‌ها

ارزیابی روش‌ها بر روی داده‌های آزمون، بر مدل ارائه‌شده،  $k$  نزدیک‌ترین همسایه انجام شده است. در روش پیشنهادی نقطه شروع  $X_0 = X$  در نظر می‌گیریم. در تولید داده پرت از dd\_tools toolbox [۲۰] که از توصیف داده‌های هر کلاس و استفاده از شعاع و مرکز هر کلاس و اضافه کردن مقدار تصادفی در هر بعد داده پرت بهره می‌برد، استفاده می‌کنیم. از نرخ داده پرت ۱۰٪ در تمامی مجموعه داده‌ها جهت مقایسه روش‌ها استفاده می‌کنیم. نویز برچسب با نرخ ۱۰٪، ۲۰٪ و ۳۰٪ بر روی داده‌های آموزش استفاده شده است. هدف اصلی در آزمایش‌ها، بررسی مقاوم بودن و کارایی روش معرفی شده در حضور داده پرت و برچسب نویزی است.

در جدول‌های زیر نتایج آزمایشات بر روی مجموعه داده‌های مختلف نشان داده شده است. همانطور که مشاهده می‌شود روش ارائه‌شده در اکثر مجموعه‌های داده، کارایی بهتری از 1-NN و LMNN دارد. در مقایسه با روش RNCA، در نویز برچسب با نرخ کمتر روش ارائه‌شده کارایی بهتری دارد. روش معرفی شده، برخلاف تابع هزینه روش LMNN که در برابر داده با برچسب‌های نویزی کارایی مناسبی ندارد، با اعمال جمله تنظیم کننده مناسب در تابع هزینه نسبت به حضور برچسب نویزی مقاوم می‌باشد.

جدول ۲- بررسی و مقایسه روش ارائه‌شده با سایر روش‌ها در حضور نویز برچسب با نرخ‌های متفاوت در مجموعه داده‌های UCI.

مجموعه داده	1-NN	LMNN	RNCA	Proposed method	نویز برچسب (%)
ionosphere	۱۷,۹۸	۱۶,۸۹	۱۶,۰۰	۱۳,۸۰	۱۰
	۱۹,۵۸	۱۹,۰۱	۱۷,۵۰	۱۴,۹۵	۲۰
	۲۲,۸۶	۲۲,۴۳	۱۹,۰۰	۱۸,۹۸	۳۰
Heart	۲۳,۰۶	۲۰,۳۸	۱۹,۹	۲۰,۲۱	۱۰
	۲۵,۰۷	۲۴,۶۷	۲۴,۰۰	۲۰,۷۸	۲۰
	۲۸,۹۸	۲۶,۴۱	۲۸,۵۰	۲۱,۳۸	۳۰
balance	۲۸,۳۵	۱۷,۴۳	۱۶,۰۰	۱۶,۳۷	۱۰
	۳۲,۵۵	۲۷,۵۵	۲۲,۵	۲۴,۹۶	۲۰
	۳۶,۶۷	۳۵,۴۵	۳۲,۰۰	۳۴,۶۷	۳۰
wine	۲۴,۹۷	۵,۸۳	۲,۱۰	۲,۰۴	۱۰
	۲۹,۰۰	۸,۷۸	۴,۵۰	۶,۵۵	۲۰
	۳۸,۰۱	۱۴,۸۷	۷,۵۰	۹,۸۰	۳۰
iris	۸,۴۵	۸,۸۶	۳,۵۰	۳,۴۰	۱۰
	۱۱,۸۷	۱۲,۲۲	۷,۹۰	۷,۸۸	۲۰
	۱۸,۵۶	۱۸,۴۰	۱۴,۰۰	۱۴,۵۱	۳۰

جدول ۳- مقایسه روش ارائه‌شده با سایر روش‌ها در حضور نویز برچسب با نرخ‌های متفاوت در مجموعه داده‌های تصویر.

مجموعه داده‌ها	1-NN	LMNN	Proposed method	نویز برچسب (%)
Olivetti face	۱۳,۲۴	۱۲,۹۸	۸,۰۲	۱۰
	۱۸,۸۶	۲۵,۳۶	۱۰,۲۵	۲۰
	۲۲,۵۰	۳۸,۴۵	۱۵,۹۵	۳۰
MNIST	۲۰,۲۲	۱۹,۵۷	۱۶,۹۶	۱۰
	۲۲,۶۷	۲۱,۸۲	۱۷,۲۳	۲۰
	۳۳,۷۸	۲۶,۱۲	۲۱,۲۸	۳۰

در آزمایش بعدی میزان مقاوم بودن روش پیشنهادی در برابر داده پرت بررسی شده است. در جدول نتایج این آزمایش گزارش شده است. همانطور که مشاهده می‌شود مدل ارائه‌شده دارای تخمین خطای کمتری است و در مقایسه با روش LMNN و 1-NN کارایی بیشتری دارد.

جدول ۴- بررسی و مقایسه روش ارائه شده با مدل LMNN و طبقه بند K نزدیکترین

همسایه، شامل ۱۰٪ داده پرت.

مجموعه داده	1-NN	LMNN	Proposed method
ionosphere	۱۷,۳۵	۱۵,۴۵	۱۳,۴۵
Heart	۲۷,۰۰	۲۴,۰۱	۲۲,۴۸
balance	۲۹,۹۸	۱۵,۶۶	۱۵,۶۶
wine	۲۶,۱۳	۶,۹۰	۶,۱۰
iris	۵,۸۷	۶,۳۴	۴,۹۷
Olivetti face	۶,۸۲	۵,۳۷	۴,۳۷
MNIST	۳۳,۰۰	۲۲,۵۶	۲۱,۶۷

نتیجه‌گیری، پیشنهاد و توصیه‌های آتی

اکثر روش‌های یادگیری متریک عملکرد مناسبی در مقابل داده‌های شامل مقادیر پرت و برچسب نویز ندارند. از جمله این روش‌ها می‌توان به روش محبوب LMNN اشاره کرد. برای حل این مشکل، روش ارائه شده، متریک موردنظر را با قسمتی از داده‌ها که بدون تأثیرات ذکر شده است، آموزش می‌دهد. این روش با استفاده از نرم یک خطا داده‌های پرت را شناسایی می‌کند. برای حل مسئله بهینه‌سازی، از روش BCD استفاده شده است. کارایی روش معرفی شده، بر روی داده‌های UCI و تصویر در شرایط مختلف مورد بررسی قرار گرفت. بر اساس نتایج آزمایش‌ها، روش ارائه شده در حضور داده‌های پرت و برچسب نویزی نسبت به روش LMNN برتری قابل توجهی دارد. به‌طور کلی نتایج آزمایش‌ها روی داده‌های مختلف، برتری روش پیشنهادی را به خصوص در حضور برچسب نویزی نشان می‌دهد.

برخی توصیه‌های برای کارهای آینده در این زمینه عبارتند از:

- ۱- اعمال مدل ارائه شده بر روی سایر روش‌های یادگیری متریک.
- ۲- تبدیل مدل ارائه شده به صورت آنلاین برای اجرا بر روی پایگاه داده‌های بزرگ‌تر

مراجع:

- [۱۳] R. He, B. Hu, W.-S. Zheng, and Y. Guo, "Two-stage sparse representation for robust recognition on large-scale database".
- [۱۴] H. Wang, F. Nie, and H. Huang, "Robust Distance Metric Learning via Simultaneous L1-Norm Minimization and Maximization." pp. 1836-1844.
- [۱۵] S. Xiang, F. Nie, and C. Zhang, "Learning a Mahalanobis distance metric for data clustering and classification," *Pattern Recogn.*, vol. 41, no. 12, pp. 3600-3612, 2008.
- [۱۶] D. Wang, and X. Tan, "Robust distance metric learning in the presence of label noise".
- [۱۷] D. Wang, and X. Tan, "Robust Distance Metric Learning via Bayesian Inference," *IEEE Transactions on Image Processing*, vol. 27, no. 3, pp. 1542-1553, 2018.
- [۱۸] J. Kim, Y. He, and H. Park, "Algorithms for nonnegative matrix and tensor factorizations: A unified view based on block coordinate descent framework," *Journal of Global Optimization*, vol. 58, no. 2, pp. 285-319, 2014.
- [۱۹] Y. Xu, and W. Yin, "A block coordinate descent method for regularized multiconvex optimization with applications to nonnegative tensor factorization and completion," *SIAM Journal on imaging sciences*, vol. 6, no. 3, pp. 1758-1789, 2013.
- [۲۰] D. Tax, "Data description toolbox dd tools 1.9. 0," 2012.

- [۱] T. Cover, and P. Hart, "Nearest neighbor pattern classification," *IEEE transactions on information theory*, vol. 13, no. 1, pp. 21-27, 1967.
- [۲] F. Wang, and J. Sun, "Survey on distance metric learning and dimensionality reduction in data mining," *Data mining and knowledge discovery*, vol. 29, no. 2, pp. 534-564, 2015.
- [۳] A. Bellet, A. Habrard, and M. Sebban, "A survey on metric learning for feature vectors and structured data," *arXiv preprint arXiv:1306.6709*, 2013.
- [۴] S. Theodoridis, A. Pikrakis, K. Koutroumbas, and D. Cavouras, *Introduction to pattern recognition: a matlab approach*: Academic Press, 2010.
- [۵] E. P. Xing, M. I. Jordan, S. J. Russell, and A. Y. Ng, "Distance metric learning with application to clustering with side-information." pp. 521-528.
- [۶] D.-M. Tsai, and C.-C. Lin, "Fuzzy C-means based clustering for linearly and nonlinearly separable data," *Pattern recognition*, vol. 44, no. 8, pp. 1750-1760, 2011.
- [۷] D. Lim, G. Lanckriet, and B. McFee, "Robust structural metric learning." pp. 6, ۶۲۳-۱۵
- [۸] M. P. Kumar, P. H. Torr, and A. Zisserman, "An invariant large margin nearest neighbour classifier." pp. 1-8.
- [۹] E. Hasanbelliu, K. Kampa, J. C. Principe, and J. T. Cobb, "Online learning using a Bayesian surprise metric." pp. 1-8.
- [۱۰] K. Huang, R. Jin, Z. Xu, and C.-L. Liu, "Robust metric learning by smooth optimization," *arXiv preprint arXiv:1203.3461*, 2012.
- [۱۱] Z.-J. Zha, T. Mei, M. Wang, Z. Wang, and X.-S. Hua, "Robust distance metric learning with auxiliary knowledge".
- [۱۲] S. C. Hoi, W. Liu, and S.-F. Chang, "Semi-supervised distance metric learning for collaborative image retrieval and clustering," *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, vol. 6, no. 3, pp. 18, 2010.