

یادگیری متریک نیمه نظارتی در فضای لایه‌ای با بهره‌گیری دقیق‌تر از دانش پیشین

زهرة کریمی، سعید شیرینی قیداری و روح‌اله رضانی

داده‌های بدون برچسب به میزان زیادی در دسترس هستند در سال‌های اخیر رویکردهای یادگیری متریک نیمه‌نظارتی مطرح شده است. تمرکز ما بر وظیفه دسته‌بندی و بر رویکردهایی است که فرض قرارگیری داده روی یک یا چند منیفلد را دارند. این رویکردها فرض همواربودن روی منیفلد را به مسأله بهینه‌سازی یادگیری متریک اعمال می‌کنند که بیانگر این است که داده‌های نزدیک به یکدیگر با احتمال بالایی برچسب یکسان دارند. هرچند این روش‌ها در مورد داده‌هایی که در فضای لایه‌ای قرار دارند در معرض چالش‌هایی قرار دارد. فضای لایه‌ای فضایی است که در آن داده‌ها روی چند منیفلد قرار می‌گیرند که ممکن است با یکدیگر اشتراک داشته و ابعاد ذاتی آنها متفاوت باشد. این چالش‌ها عبارتند از:

- شباهت داده‌های نزدیک به یکدیگر که در روش‌های موجود اعمال شده است در نواحی تقاطع منیفلدها لزوماً برقرار نیست.
- دسته‌بندی که از متریک یاد گرفته شده استفاده کرده و جهت تعیین برچسب داده‌ها به کار می‌رود به خاطر تعداد کم داده‌های برچسب‌دار دقت لازم را ندارد.

در این مقاله، یک رویکرد یادگیری متریک نیمه‌نظارتی برای دسته‌بندی داده‌هایی که در فضای لایه‌ای قرار دارد پیشنهاد شده که از دانش پیشین موجود که همان فرض همواربودن در فضای لایه‌ای است به صورت دقیق‌تر در موارد ذیل بهره‌برداری می‌کند:

- عدم اعمال فرض همواربودن در نواحی تقاطع منیفلدها
- انتشار داده‌های برچسب‌دار به نزدیک‌ترین نقاط به این داده‌ها در نواحی داخلی منیفلدها، بر اساس فرض همواربودن روی منیفلد
- موارد مذکور بر مبنای تفکیک نقاط داخلی منیفلد از سایر نقاط است که با توجه به تئوری موجود در این زمینه انجام می‌شود.
- در ادامه، ابتدا در بخش ۲ پژوهش‌های مرتبط بررسی گردیده و سپس در بخش ۳ جزئیات روش پیشنهادی شرح داده شده است. در بخش ۴ نتایج آزمایش‌های انجام‌شده و در بخش ۵ نتایج و کارهای پیشنهادی آمده است.

۲- پژوهش‌های مرتبط

در سال‌های اخیر روش‌های یادگیری متریک زیادی پیشنهاد شده [۱]، [۳] و [۴] که این روش‌ها در سه گروه بدون ناظر، نظارتی و نیمه‌نظارتی قرار می‌گیرد.

رویکردهای یادگیری متریک بدون ناظر سعی در حفظ توپولوژی و یا ساختار هندسی داده در حد امکان دارند. این رویکردها به فرم مسأله کاهش بعد و/یا یادگیری منیفلد مطرح می‌شوند. هدف اکثر رویکردهای یادگیری منیفلد، یادگیری بازنمایی بعد پایین داده است به نحوی که فاصله بین داده‌ها حفظ شود. معروف‌ترین رویکردهای اولیه در این زمینه

چکیده: یادگیری متریک نیمه‌نظارتی مبتنی بر منیفلد در سال‌های اخیر بسیار مورد توجه واقع شده است. این رویکردها، منظم‌سازی مبتنی بر فرض همواربودن داده‌ها روی منیفلد را اعمال می‌کنند، هرچند در معرض دو چالش قرار دارند: (۱) شباهت بین دسته‌های مختلف، تقاطع منیفلدها با یکدیگر را ایجاد می‌کند که با فرض همواربودن برچسب در این نواحی در تناقض است. (۲) دسته‌بند ۱NN که برای تعیین برچسب داده‌ها در مسایل یادگیری متریک اعمال می‌شود با وجود تعداد کم داده‌های برچسب‌دار دقت مناسب را ندارد. در این مقاله روشی برای یادگیری متریک نیمه‌نظارتی با فرض قرارگیری داده‌ها در فضای لایه‌ای ارائه شده که در آن از دانش پیشین موجود که همان فرض همواربودن داده‌ها روی هر منیفلد است به صورت دقیق‌تر بهره‌برداری شده است. در مرحله یادگیری متریک، فرض همواربودن در نواحی تقاطع اعمال نشده و در مرحله دسته‌بندی، داده‌های برچسب‌دار در نقاط داخلی منیفلدها بر اساس فرض همواربودن توسعه داده شده است. تفکیک نقاط تقاطع منیفلدها از سایر نقاط بر مبنای رفتار متمایز لاپلاسیان تابع هموار روی هر منیفلد در نقاط داخلی نسبت به سایر نقاط صورت می‌گیرد. آزمایش‌ها نشان‌دهنده دقت خوب روش پیشنهادی نسبت به روش‌های مشابه است.

کلیدواژه: یادگیری متریک نیمه‌نظارتی، فضای لایه‌ای، لاپلاسیان، فرض همواربودن.

۱- مقدمه

مطالعات سال‌های اخیر نشان داده که یادگیری متریک، نقشی اساسی در یادگیری ماشین از جمله دسته‌بندی، خوشه‌بندی و بازیابی اطلاعات دارد [۱] تا [۳]. از آنجا که هر وظیفه، فضای معنایی خود را جهت تعریف فاصله‌ها دارد، فاصله اقلیدسی اغلب نمی‌تواند این فضا را به درستی نشان دهد و انتخاب مناسب متریک تأثیر زیادی در نتیجه به دست آمده دارد. هدف یادگیری متریک، یادگیری یک تابع فاصله با مقدار حقیقی است که نقاط مشابه را به یکدیگر نزدیک نگه داشته و نقاط نامشابه را از یکدیگر دور می‌کند. اغلب روش‌های یادگیری متریک یک تابع فاصله متریک را از دانشی که به شکل محدودیت‌های دوتایی یا سه‌تایی در دسترس هستند یاد می‌گیرند. این محدودیت‌ها از داده‌های برچسب‌دار موجود استنتاج می‌شوند. از آنجا که برچسب‌گذاری داده‌ها هزینه زیادی داشته اما

این مقاله در تاریخ ۵ مهر ۱۳۹۶ دریافت و در تاریخ ۳۰ اردیبهشت ماه ۱۳۹۷ بازنگری شد.

زهرة کریمی، دانشکده مهندسی کامپیوتر، دانشگاه صنعتی امیرکبیر، تهران، ایران، (email: z_karimi@aut.ac.ir).

سعید شیرینی قیداری (نویسنده مسئول)، دانشکده علوم کامپیوتر، دانشگاه صنعتی امیرکبیر، تهران، ایران، (email: shiry@aut.ac.ir).

روح‌اله رضانی، دانشکده ریاضی و علوم کامپیوتر، دانشگاه دامغان، دامغان، ایران، (email: r_ramezani@du.ac.ir).

تنک^۶ نیمه‌نظارتی ماتریس مجاورت را با استفاده از انتشار مجاورت یاد می‌گیرد [۱۶]. یک ماتریس شاخص همسایگی در این روش تعریف می‌شود که فاصله بین زوج‌های مشابه را تا حد ممکن کم کرده و فاصله بین زوج‌های متفاوت را تا حد امکان زیاد می‌کند. باغشاه و همکاران، یک روش نیمه‌نظارتی یادگیری متریک مبتنی بر منیفلد پیشنهاد داده‌اند که بازنمایی خطی محلی در فضای ورودی را حفظ می‌کند [۱۷].

تمام روش‌های مذکور، نقاطی را که بر مبنای فاصله اقلیدسی نزدیک به یکدیگر قرار دارند مشابه یکدیگر در نظر گرفته و با تعریف محدودیت‌های محلی یا جملات منظم‌سازی سعی در حفظ این شباهت دارند. این رویکردها از دانش پیشین مربوط به شباهت داده‌های بدون برچسب بهره گرفته و دانش مربوط به عدم شباهت را در نظر نمی‌گیرند. Niu و همکاران، یک روش یادگیری متریک نیمه‌نظارتی مبتنی بر تئوری اطلاعات (SERAPH) از طریق منظم‌سازی انتروپی پیشنهاد داده که برخی محدودیت‌های رویکردهای مبتنی بر منیفلد را مرتفع نموده است [۱۸]. در این روش، ابتدا احتمال شرطی نسبت‌دادن برچسب‌ها با داشتن داده با فاصله ماله‌لونویس پارامتری شده و سپس اصل بیشینه انتروپی روی داده‌های برچسب‌دار و اصل کمینه انتروپی روی داده‌های بدون برچسب با استفاده از منظم‌سازی انتروپی اعمال شده است. SERAPH بر مبنای فرض جداسازی با چگالی پایین است.

Wang و همکاران، رویکرد یادگیری متریک نیمه‌نظارتی منظم‌شده^۷ (RSSML) را بر اساس هر دو فرض همواربودن روی منیفلد و فرض خوشه پیشنهاد داده‌اند [۱۹]. ماتریس مجاورت در این رویکرد بر اساس تابع هسته گوسی تعریف گردیده و چگالی محلی در آن با پنجره پارزن^۸ [۲۰] تخمین زده شده است. هدف اصلی این روش، کمینه‌کردن فاصله بین نقاط نزدیک در نواحی با چگالی زیاد جهت اعمال فرض جداسازی در چگالی پایین است.

رویکردهای مذکور در مواجهه با داده‌هایی که روی چند منیفلد متقاطع قرار دارد با چالش مواجه هستند. از آنجا که فرضیات منیفلد و جداسازی در چگالی کم در نواحی اشتراک منیفلدها می‌تواند نقض شود، در این مقاله، فرض همواربودن روی چند منیفلد در یادگیری متریک اعمال می‌شود. بدین معنا که تابع هدف به صورت قطعه‌ای روی هر منیفلد هموار است و دانش پیشین مرتبط با شباهت داده‌های برچسب‌دار و بدون برچسب در نواحی صحیح (و نه در مورد تمام داده‌ها) اعمال خواهد شد.

۳- روش پیشنهادی

در مسأله مورد نظر ما، مجموعه داده $\chi = \{x_i\}_{i=1}^{l+u} \in \mathbb{R}^{(l+u) \times d}$ شامل داده‌های برچسب‌دار $\chi_l = \{x_i\}_{i=1}^l$ و داده بدون برچسب $\chi_u = \{x_i\}_{i=l+1}^{l+u}$ است. $C_l = \{c_i\}_{i=1}^l$ مجموعه برچسب داده‌های برچسب‌دار و در یک دسته‌بندی دودویی $\{+1, -1\}$ است. داده‌ها از چند منیفلد $\Omega_i (1 \leq i \leq n_\Omega)$ که می‌توانند با یکدیگر اشتراک داشته باشند نمونه‌برداری شده است به طوری که برچسب نمونه‌ها روی هر منیفلد به صورت هموار تغییر می‌کند. n_Ω تعداد منیفلدها و نامشخص است. هدف، یافتن تابع دسته‌بندی f جهت برچسب‌گذاری داده‌های بدون برچسب است. منیفلدهای زیربنایی داده با گراف مجاورتی که از داده‌ها ساخته شده بازنمایی می‌شود. ماتریس لاپلاسی به صورت

شامل [۵] LLE، [۶] ISOMAP، [۷] LTSA^۱ و [۸] Eigenmap هستند. توسعه‌هایی از این روش‌ها در سال‌های اخیر پیشنهاد شده است. Wang و همکاران (۲۰۱۸) LTSA را با بهره‌گیری از منظم‌سازی نرم اتمی^۲ انجام می‌دهند تا همسایگی مناسب را در زمانی که برخی داده‌ها کامل نیستند بیابند [۹]. در [۱۰] نگاشت از فضای تانژانت با ابعاد پایین به فضای نمونه با ابعاد بالا و نگاشت معکوس آن به صورت هم‌زمان یاد گرفته می‌شود که عمل یادگیری منیفلد و نیز بازیابی داده‌های گمشده را به صورت هم‌زمان انجام می‌دهد.

تلاش اکثر روش‌های یادگیری منیفلد در راستای بازکردن^۳ منیفلد است و فرض می‌شود بعد ذاتی منیفلد به صورت دانش پیشین در دسترس است. Piterlis و همکاران یادگیری یک منیفلد را به صورت یافتن یک اطلس^۴ یا مجموعه‌ای از نمودارهای^۵ متقاطع بیان نموده است که هم‌زمان پارامترهای پیوسته تعریف نمودارها و نحوه تخصیص نقاط به نمودارها را می‌یابد [۱۱]. روش بیان‌شده امکان بازنمایی منیفلدهای بسته را نیز دارد. در سال‌های اخیر جهت بازنمایی، داده‌های مولتی‌منیفلد با ابعاد ذاتی متفاوت ارائه شده است. Wang و همکاران ابتدا بعد ذاتی هر داده را تخمین زده و سپس داده‌ها را به خوشه‌هایی تقسیم می‌کنند به نحوی که تعداد ابعاد ذاتی داده‌های هر خوشه یکسان باشد. سپس منیفلدهای هر خوشه را یاد گرفته و یک مدل متغیر نهان را برای هر منیفلد مقداردهی اولیه می‌کند و نهایتاً یک مدل آمیخته سلسله‌مراتبی برای منیفلدها می‌سازد [۱۲].

در روش‌های نظارتی یادگیری متریک تلاش بر آن است که نمونه‌های یک دسته تا حد امکان نزدیک به یکدیگر و نمونه‌های دسته‌های متفاوت دور از یکدیگر قرار گیرند. روش LMNN از معروف‌ترین روش‌ها است که از اطلاعات محلی مجموعه آموزش استفاده می‌کند [۱۳]. Der و سایرین یک مدل متغیر نهان را پیشنهاد داده‌اند و الگوریتم EM را جهت تخمین پارامتر به کار برده‌اند [۱۴]. ملاک‌های تئوری اطلاعات مانند دیورژانس لگاریتم دترمینان در برخی رویکردهای یادگیری متریک نظارتی نیز به کار رفته و به نتایج خوبی دست یافته‌اند.

رویکردهای نظارتی در صورت عدم وجود داده برچسب‌دار کافی در معرض بیش‌برازش قرار دارد و رویکردهای بدون ناظر از دانش موجود در داده‌های برچسب‌دار که می‌تواند متریک مناسب برای کاربرد خاص را ارائه کند بهره‌برداری نمی‌نماید. در ادامه، رویکردهای نیمه‌نظارتی که مرتبط با این مقاله هستند بررسی شده‌اند.

رویکردهای نیمه‌نظارتی بر اساس فرضی که در مورد داده‌های بدون برچسب دارند در دو دسته رویکردهای مبتنی بر فرض خوشه و مبتنی بر فرض منیفلد قرار می‌گیرند. در رویکردهای مبتنی بر منیفلد، ماتریس مجاورت، W ، که میزان نزدیکی تمام داده‌ها با یکدیگر را اندازه‌گیری می‌کند تعریف شده و یک جمله منظم‌سازی در جملات بهینه‌سازی یادگیری متریک سعی در حفظ این مجاورت دارد. یادگیری متریک منظم‌شده لاپلاسی (LRML) از داده‌های بدون برچسب بدین روش بهره گرفته و به خوشه‌بندی و بازیابی تصاویر اعمال شده است [۱۵]. در ماتریس مجاورت LRML، اگر x_i در بین k نزدیک‌ترین همسایه x_j باشد $W_{ij} = 1$ و در غیر این صورت برابر صفر است. یادگیری متریک

1. Local Tangent Space Alignment
2. Nuclear Norm
3. Unwrapping
4. Atlas
5. Charts

6. Sparse

7. Regularized Semi-Supervised Metric Learning

8. Parzen Window

جدول ۱: نمادهای استفاده‌شده در مقاله.

نماد	مفهوم
$\mathcal{X} = \{x_i\}_{i=1}^{l+u} \in \mathfrak{R}^{(l+u) \times d}$	مجموعه داده‌ها
n	برابر کل تعداد داده‌ها $n = l + u$
\mathcal{X}_l	مجموعه داده‌های برچسب‌دار
\mathcal{X}_u	مجموعه داده‌های بدون برچسب
$C_l = \{c_i\}_{i=1}^l$	مجموعه برچسب داده‌های برچسب‌دار
L	لاپلاسیان گراف، زمانی که تمرکز روی پارامتر آن است از L_ω استفاده می‌شود.
D	ماتریس درجه گراف
W	ماتریس وزن‌های گراف
ω	پارامتر هسته گاوسی
g	تابع قطعه‌ای هموار روی منیفلدهای متقاطع
g_i	تابع هموار روی i امین منیفلد
Γ_ω	عملگر لاپلاسیان با پارامتر ω
Δ	عملگر لاپلاس - بلترامی
K	تعداد نزدیک‌ترین نقاط در گراف nm
f	تابع دسته‌بندی
n_{nonint}	تعداد نقاط غیر داخلی
M	ماتریس مثبت نیمه‌معین
D_M	تابع فاصله ماهالونوبیس پارامتری شده با ماتریس M
$W^{(1)}$	ماتریس وزن مربوط به همسایه‌های داده‌های برچسب‌دار
$W^{(2)}$	ماتریس وزن مربوط به همسایه‌های تمام داده‌ها
X_b	مجموعه داده‌های برچسب‌گذاری شده شامل و داده‌های برچسب‌گذاری شده اولیه و داده‌های برچسب‌گذاری شده با عمل انتشار برچسب

اشتراک منیفلدها برقرار است و به بیان دیگر در نقاط داخلی منیفلدها این فرض برقرار می‌باشد. در واقع، اعمال فرض همواربودن در نواحی اشتراک منیفلدها دانش پیشین اشتباهی را در مورد داده‌هایی که در فضای لایه‌ای قرار دارد به مدل اعمال می‌کند. (۲) داده‌های بدون برچسبی که به اندازه کافی نزدیک به داده‌های برچسب‌دار هستند و در نواحی داخلی منیفلدها واقع شده‌اند با احتمال بالایی برچسب یکسانی با داده‌های برچسب‌دار همسایه خود دارند.

جهت اعمال موارد مذکور نیاز به تمایز بین نقاط داخلی منیفلد از سایر نقاط است. فرض کنید Ω مجموعه $\Omega_i (1 \leq i \leq n_\Omega)$ باشد. تابع قطعه‌ای هموار $\mathfrak{R} : \bar{\Omega} \rightarrow \mathfrak{R}$ بدین صورت تعریف می‌شود: اگر g_i تابع محدودشده به منیفلد i ام باشد، g_i در نقاط داخلی منیفلد i ام، $C_2 - C_1$ پیوسته است. اگر نقاط به صورت یکنواخت از یک منیفلد هموار نمونه‌برداری شده باشد در نقاط داخلی منیفلد هنگامی که تعداد نقاط، n ، به سمت بی‌نهایت میل کرده و ω با نرخ مناسب به سمت صفر میل کند،

$$\Gamma_\omega g(x) = \Delta g(x) + O(1) \quad (2)$$

که در آن Δ عملگر لاپلاس - بلترامی روی منیفلد است. لاپلاسیان گراف اعمال‌شده به تابع g به این صورت محاسبه می‌شود

$$L_\omega g(x) = \sum_{j=1}^n w_{ij}(\omega) [g(x) - g(x_j)] \quad (3)$$

$L_\omega g$ فرم ناپیوسته Γ_ω است. در بخش‌های بعد Lg و w_{ij} به ترتیب به جای $L_\omega g$ و $w_{ij}(\omega)$ استفاده خواهد شد مگر در صورتی که تمرکز ما روی پارامتر ω باشد.

نتایج به دست آمده از مطالعات نظری در سال‌های اخیر نشان داده که نقاط غیر داخلی شامل نقاط اشتراک منیفلدها، مرز و لبه‌ها جنبه‌های مهمی از داده هستند که در بسیاری از پژوهش‌ها در نظر گرفته نشده‌اند. برای نقطه x در همسایگی \in نقاط غیر داخلی، (۴) برقرار است [۲۲]

$$\Gamma_\omega g(x) = \frac{1}{\sqrt{\omega}} \frac{\pi}{2} \partial_{\bar{n}} g(x) + O\left(\frac{1}{\sqrt{\omega}}\right) \quad (4)$$

که در آن $\partial_{\bar{n}}$ نرمال واحد به سمت خارج^۱ در نقطه x ، $\partial_{\bar{n}} g(x)$ مشتق جهت‌دار g در جهت $\partial_{\bar{n}}$ و g تابع قطعه‌ای هموار شرح داده شده در بخش ۳-۱ است. از این ویژگی برای دسته‌بندی داده‌ها در فضای لایه‌ای استفاده شده است [۲۳].

بر طبق معادله بیان‌شده، $L_\omega g$ از مرتبه $O(1/\sqrt{\omega})$ است که بزرگ‌تر از $L_\omega g$ در نقاط داخلی است که برای مقادیر کوچک ω از مرتبه $O(1)$ است. بنابراین مقادیر بزرگ‌تر $L_\omega g$ ، بیانگر نقاط نزدیک نقاط غیر داخلی است.

تعریف تابع g به دلیل کمبود دانش قبلی در خصوص این که چه داده‌ای متعلق به کدام منیفلد است و نیز بعد بالای داده ورودی ساده نیست. هرچند با توجه به فرض همواربودن روی چند منیفلد، تابع دسته‌بندی روی منیفلد، f ، یک تابع قطعه‌ای هموار است، لذا پیشنهاد می‌گردد که با محاسبه لاپلاسیان تابع f ، نقاط داخلی از نقاط غیر داخلی تمیز داده شود.

جهت وضوح بیشتر، نمادهای استفاده‌شده مقاله در جدول ۱ آمده است.

$L = D - W$ تعریف شده که در آن W ماتریس مجاورت گراف و D ماتریسی قطری است که $D(i, i) = \sum_j w_{ij}$ می‌باشد. هر عنصر ماتریس مجاورت برای داده‌هایی که در k نزدیک‌ترین همسایگی یکدیگر هستند به صورت $w_{ij} = \exp(-\|x_i - x_j\|^2 / 2\omega)$ تعریف شده و برای سایر داده‌ها برابر صفر است.

هدف، یادگیری فاصله متریک ماهالونوبیس به فرم ذیل است

$$D_M(x_i, x_j) = \sqrt{(x_i - x_j)^T M (x_i - x_j)} \quad (1)$$

که در آن M ماتریس مثبت نیمه‌معین است. همانند LRML و RSSML مسأله تخمین پارامتر در مدل ما به عنوان مسأله بهینه‌سازی در فضای ماتریس‌های مثبت نیمه‌معین فرمول‌بندی شده است. همانند اکثر مطالعات مرتبط با یادگیری متریک با هدف دسته‌بندی، تلاش ما در بهبود دسته‌بندی نزدیک‌ترین همسایه است [۳] و [۲۱]. مؤلفه‌های اصلی مدل پیشنهادی عبارتند از (۱) اعمال محدودیت‌های محلی به صورت دقیق برای یادگیری متریک داده‌هایی که روی منیفلدهای متقاطع قرار دارند. (۲) انتشار داده‌های برچسب‌دار برای غلبه بر کمبود داده‌های برچسب‌دار. در ادامه، ابتدا شهود روش پیشنهادی و سپس جزئیات آن ارائه شده است.

۳-۱ شهود

مدل ما بر اساس شهودهای زیر است که از فرض همواربودن در فضای لایه‌ای استنتاج می‌شود: (۱) بر اساس فرض همواربودن روی چند منیفلد، برچسب هر داده با احتمال بالایی برابر با برچسب نزدیک‌ترین همسایه‌ها است. این فرض برای تمام داده‌ها به جز برای داده‌های ناحیه

$$\min tr(XL'XM) + \gamma \sum [\mathbb{1} + D_{ij}^y - D_{ik}^y]_+ \quad (12)$$

$$s.t. M \succ = 0.$$

که عملگر $\succ =$ بیانگر محدودیت مثبت نیمه معین بودن است. مسأله بهینه‌سازی مذکور به فرم زیر بازنویسی می‌شود:

$$\min tr(XL'XM) + \gamma \sum_{i,j,k=1}^l \xi_{ijk} \quad (13)$$

$$s.t.$$

$$D_{ij}^y - D_{ik}^y \geq \mathbb{1} - \xi_{ijk}$$

$$\xi_{ijk} \geq 0$$

$$M \succ = 0.$$

مسأله بهینه‌سازی مذکور از نوع برنامه‌ی نیمه‌معین^۱ است که می‌تواند با حل‌کننده‌های^۲ همه‌منظوره حل شود [۲۴]، هرچند این روش حل نسبت به تعداد محدودیت‌ها مقیاس‌پذیر نیست. Weinberger و همکاران یک روش حل کارا بر اساس نزول گرادیان پیشنهاد کرده‌اند که سریع‌تر از solver معمول است [۱۲]. ما از رویکردی مشابه جهت یادگیری متریک استفاده کرده‌ایم بدین صورت که در تکرار t ام، محدودیت‌هایی که ارضا نمی‌شوند مشخص شده $(\psi(t))$ و متریک کنونی $(M(t))$ جهت ارضای آنها بروز می‌شود. جهت گرادیان $(\nabla F(t))$ بر اساس این محدودیت‌ها محاسبه می‌شود. $F(t)$ به صورت زیر تعریف می‌شود

$$F(t) = \min tr(XL'XM(t)) + \gamma \sum_{i,j,k \in \psi(t)} [\mathbb{1} + D_{ij}^y(t) - D_{ik}^y(t)]_+ \quad (14)$$

رویه مذکور تا زمان ارضای تمام محدودیت‌ها یا رسیدن تعداد تکرارها به تعداد تکرار از قبل مشخص شده ادامه می‌یابد.

۳-۳ اعمال دسته‌بند نزدیک‌ترین همسایه با انتشار برچسب داده‌ها در سطح مکان

جهت دسته‌بندی، متریک یاد گرفته شده به یک دسته‌بند مبتنی بر فاصله که اغلب دسته‌بند نزدیک‌ترین همسایه (NN) است اعمال می‌شود [۳]. دسته‌بند نزدیک‌ترین همسایه یک روش ناپارامتری دسته‌بندی و بسیار توانمند است. هرچند تعداد کم داده‌های برچسب‌دار در دسته‌بند NN که جهت تعیین برچسب داده‌ها در رویکردهای مبتنی بر متریک به کار می‌رود مسأله‌ساز است، لذا مدلی برای انتشار داده‌های برچسب‌دار پیشنهاد می‌شود. انگیزه اعمال این مدل، استفاده بیشتر از دانش پیشین در دسترس است. از آنجا که فرض همواربودن داده‌ها روی هر منیفولد در فضای لایه‌ای برقرار است، می‌توان از این فرض جهت افزایش تعداد داده‌های برچسب‌دار با اطمینان بالا استفاده کرد. تلاش رویکردهای موجود، در اعمال این فرض فقط به داده‌های بدون برچسب و در مرحله یادگیری متریک بوده است. از فرض همواربودن برچسب‌ها روی هر منیفولد، چنین استنتاج می‌شود که نزدیک‌ترین داده بدون برچسب به هر داده برچسب‌دار در نقاط داخلی منیفولد با احتمال بالایی برچسب یکسان با داده برچسب‌دار مورد نظر دارد. لذا در این مرحله، k_l نزدیک‌ترین داده برچسب‌دار به داده‌های برچسب‌دار موجود در نواحی داخلی منیفولدها به نقاط برچسب‌دار

ورودی: $\mathcal{X} = \{x_i\}_{i=1}^{l+u}$ ، C_l ، n_{nonint} ، k_l و η ، stepsize

۲-۳ مسأله بهینه‌سازی

مسأله بهینه‌سازی ارائه‌شده شامل دو تابع هزینه است: تابع هزینه اول بر اساس محدودیت‌های دوتایی و تابع هزینه دوم بر اساس محدودیت‌های سه‌تایی است. با داشتن زوج‌های (x_i, x_j) یادگیری متریک سعی در یافتن ماتریس مثبت نیمه معین M دارد به نحوی که x_i و x_j را به یکدیگر نزدیک نگه دارد و بر مبنای محدودیت‌های سه‌تایی (x_i, x_j, x_k) تلاش می‌کند فاصله بین x_i و x_j را بزرگ‌تر از فاصله بین x_i و x_k داشته باشد. در توابع $\ell_{pull}(M)$ و $\ell_{push}(M)$ به ترتیب نقض موارد مذکور جریمه می‌شود. هدف، کمینه‌سازی جملات زیر است

$$\ell_{pull}(M) + \lambda \ell_{push}(M) \quad (5)$$

پارامتر وزن دهی λ اهمیت جملات مذکور نسبت به یکدیگر را مشخص می‌کند. بر طبق [۱۳] برای هر x_i ، x_j و x_k به ترتیب همسایه هدف و همسایه impostor گفته می‌شود و تمایز روش ما با این روش در تعریف این دو نوع همسایگی است. برای هر نقطه در داخل منیفولد، k_l همسایه با برچسب یکسان به عنوان همسایه‌های هدف و k_l نزدیک‌ترین همسایه با برچسب غیر یکسان به عنوان همسایه impostor در نظر گرفته شده است.

برای محاسبه $\ell_{pull}(M)$ ، دو ماتریس وزن $W^{(v)}$ و $W^{(s)}$ تعریف می‌شود: ماتریس وزن $W^{(s)}$ به صورت ماتریسی که عناصر درایه‌های همسایه‌های هدف داده‌های برچسب‌دار آن یک و سایر عناصر آن صفر است تعریف می‌شود. ماتریس وزن $W^{(v)}$ ماتریسی است که در آن k نزدیک‌ترین همسایه هر داده به عنوان همسایه‌های هدف تعریف شده و عناصر درایه‌های همسایه‌های هدف به صورت زیر تعریف می‌شود

$$\exp\left(-\frac{\|x_i - x_j\|^2}{2w}\right) \quad (6)$$

تابع هزینه $\ell_{pull}(M)$ به صورت زیر تعریف می‌شود

$$\ell_{pull}(M) = \frac{1}{\gamma} \sum_{i,j \in X_b} W_{ij}^{(s)} D_{ij}^y + \frac{1}{\gamma} \sum_{i,j \in X_b} W_{ij}^{(v)} D_{ij}^y = \frac{1}{\gamma} \sum \sum W_{ij}^{(v)} \|U^T(x_i - x_j)\|^2 = tr(U^T XL'XU) = tr(XL'XUU^T) = tr(XL'XM) \quad (7)$$

$$L' = L^{(s)} + L^{(v)} \quad (8)$$

$$L^{(m)} = D_w^{(m)} - W^{(m)} \quad (9)$$

$$D_w^{(m)}(j, j) = \sum_i w_{ij}^{(m)} \quad (9)$$

نماد tr نشان دهنده عملگر trace است.

جهت اعمال حاشیه بزرگ بین impostorها و دورترین همسایه هدف، باید نامساوی زیر برقرار باشد

$$D_{ij}^y + \mathbb{1} < D_{ik}^y \quad (10)$$

تابع هزینه $\ell_{push}(M)$ به صورت زیر تعریف می‌شود

$$\ell_{push}(M) = \sum_{i,j,k} [\mathbb{1} + D_{ij}^y - D_{ik}^y]_+ \quad (11)$$

نماد $[z]_+$ به صورت $\max\{z, 0\}$ تعریف می‌شود.

مسأله بهینه‌سازی کلی به صورت زیر است

1. Semi-Definite Program
2. Solvers

خروجی: M

گام ۱: مقداردهی اولیه

- تخمین برچسب‌ها با استفاده از منظم‌سازی منیفلد
- محاسبه لاپلاسیان برچسب‌های تخمین زده شده با استفاده از (۳) و مشخص کردن نقاط داخلی منیفلدها

گام ۲: یادگیری متریک

- $\psi(t)$ را مشخص کن.
 - اگر $\psi(t)$ تهی است خاتمه.
 - $\nabla F(t)$ را محاسبه کن.
 - متریک کنونی را به روز رسانی کن:
- $$M(t+1) = M(t) - \text{stepSize} \times \nabla F(t)$$
- $M(t+1)$ را به زیرفضای شامل ماتریس‌های مثبت نیمه‌معین تصویر کن.
 - $M(t)$ را برابر $M(t+1)$ قرار بده.

گام ۳: دسته‌بندی

- فرض همواربودن را به نزدیک‌ترین همسایه‌های داده‌های برچسب‌دار اعمال کن تا X_b برچسب‌گذاری شود.
- دسته‌بند $1NN$ را به باقیمانده داده‌های بدون برچسب اعمال کن.

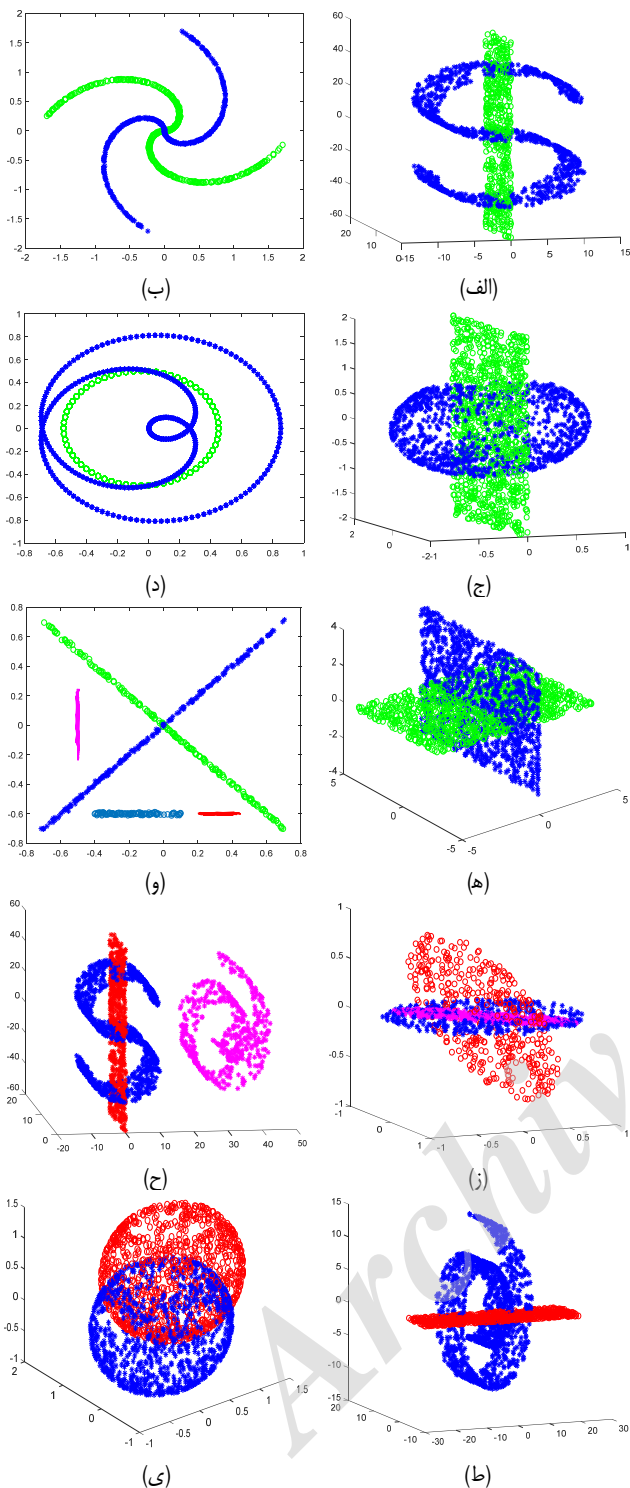
شکل ۱: گام‌های الگوریتم یادگیری متریک نیمه‌نظارتی در فضای لایه‌ای.

موجود اضافه می‌شوند و این مجموعه جدید داده‌های برچسب‌دار، X_b نام‌گذاری می‌شود. تعیین برچسب داده‌ها با دسته‌بند $1NN$ با استفاده از این مجموعه داده برچسب‌دار جدید انجام می‌شود. جزئیات الگوریتم در شکل ۱ ارائه می‌شود.

۴- آزمایش‌ها

الگوریتم پیشنهادی با فاصله اقلیدسی، دو روش یادگیری متریک نیمه‌نظارتی بروز شامل RSSML و Seraph و یک روش دسته‌بندی نیمه‌نظارتی، MR^1 [۲۵] مقایسه شده است. RSSML یک روش یادگیری متریک نیمه‌نظارتی بر اساس فرض‌های همواربودن روی منیفلد و جداسازی در چگالی کم است. Seraph یک روش یادگیری متریک نیمه‌نظارتی بر مبنای تئوری اطلاعات است که سعی کرده بر محدودیت‌های روش‌های مبتنی بر منیفلد غلبه کند. انواع متفاوتی از این الگوریتم ارائه شده که برخی از آنها برای بهره‌گیری از فرض منیفلد توسعه داده شده‌اند. MR فرمی از منظم‌سازی مبتنی بر تیخونوف^۲ است که در روش پیشنهادی برای جداسازی نقاط داخلی منیفلد از سایر نقاط به کار رفته است. تمام آزمایش‌ها ۳۰ مرتبه با داده‌های برچسب‌دار متفاوت انتخاب شده به صورت تصادفی تکرار گردیده و میانگین درصد خطا (نسبت تعداد داده‌های صحیح دسته‌بندی شده به تعداد کل داده‌ها) در این ۳۰ بار اجرا و نیز انحراف معیار درصد خطا گزارش شده است. کوچک بودن مقدار انحراف معیار بیانگر پایداربودن روش در برابر تغییر داده‌های برچسب‌دار است.

پارامترهای MR شامل λ_1 و λ_2 به ترتیب از مجموعه $\{0.001, 0.01, 0.02, 0.03, 0.04, 0.05, 0.06, 0.07, 0.08, 0.09, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1, 1.5, 2, 3, 4, 5, 6, 7, 8, 9, 10, 15, 20\}$ انتخاب گردیده است. تابع هزینه مربعی به MR اعمال شده و پارامتر K در RSSML، MR و الگوریتم پیشنهادی برابر ۱۰ است. تمام انواع Seraph شامل SeraphNone، SeraphPost، SeraphProj، SeraphHyper، SeraphMani، SeraphManiPost و SerpahManiHyper اجرا شده و پارامترهای آن بر اساس [۱۶] تنظیم شده است. چهار تنظیم اول بیان شده بر اساس فرض منیفلد بوده و سایر موارد بر اساس ترکیبی از روش‌های مبتنی بر منیفلد و مبتنی بر انتروپی است.



شکل ۲: مجموعه داده‌های مصنوعی، (الف) Dollar Sign (DS)، (ب) Two Spirals (Two Spirals)، (ج) Surface-Sphere (SS)، (د) Roll Curve and Circle (RCC)، (ه) Two (Two Spheres)، (و) Five Segments (FS)، (ز) Three Planes (TP2)، (ح) Dollar-Sign and Plane and Roll (DPR)، (ط) Roll and Plane (RP) و (ی) Two Spheres (TS2).

SeraphManiPost، SeraphMani، SeraphHyper، SerpahProj، SeraphPost و SerpahManiHyper اجرا شده و پارامترهای آن بر اساس [۱۶] تنظیم شده است. چهار تنظیم اول بیان شده بر اساس فرض منیفلد بوده و سایر موارد بر اساس ترکیبی از روش‌های مبتنی بر منیفلد و مبتنی بر انتروپی است.

1. Manifold Regularization
2. Tikhonov

جدول ۲: درصد خطا و انحراف معیار روی مجموعه داده‌های مصنوعی.

مجموعه داده											روش
DPR	TS۲	RP	FS	RCC	TS۱	TP۲	TP۱	SS	DS		
۲,۵۸	۱,۴۶	۰,۴۰	۱,۵۷	۳,۷۸	۰,۷۳	۲,۷۵	۵,۹۳	۱۱,۰۲	۱۴,۸۵		روش پیشنهادی
(۰,۰۳)	(۰,۰۱)	(۰,۰۰)	(۰,۰۲)	(۰,۰۱)	(۰,۰۱)	(۰,۰۱)	(۰,۰۲)	(۰,۰۴)	(۰,۰۹)		
۱۷,۱۶	۱۴,۴۸	۴,۱۰	۴,۰۱۶	۳۶,۶۴	۴۶,۱۱	۲۲,۵۵	۶,۴۵	۱۴,۵۸	۲۴,۸۵		MR
(۰,۰۲)	(۰,۰۲)	(۰,۰۰)	(۰,۰۵)	(۶,۵۰)	(۰,۰۶)	(۰,۰۳)	(۲,۴۶)	(۵,۶۹)	(۹,۰۲)		
۱۲,۹۵	۷,۶۱	۲,۸۶	۳,۸۶	۵,۸۶	۲,۰۵	۱۴,۸۰	۱۲,۷۷	۱۵,۸۲	۲۶,۰۵		RSSML
(۰,۰۱)	(۰,۰۱)	(۰,۰۰)	(۰,۰۲)	(۰,۰۱)	(۰,۰۱)	(۰,۰۲)	(۰,۴۷)	(۰,۰۵)	(۰,۰۷)		
۱۲,۸۱	۷,۵۰	۳,۱۷	۳,۶۸	۵,۳۹	۱,۹۷	۱۶,۷۲	۱۵,۳۰	۲۴,۰۷	۳۹,۸۳		فاصله اقلیدسی
(۰,۰۱)	(۰,۰۱)	(۰,۰۰)	(۰,۰۲)	(۰,۰۱)	(۰,۰۱)	(۰,۰۲)	(۰,۰۳)	(۰,۰۴)	(۰,۰۵)		
۳۶,۱۳	۲۰,۴۱	۱۹,۱۵	۴,۰۴	۵,۶۱	۱,۹۸	۱۹,۵۹	۲۱,۷۲	۲۲,۶۰	۲۴,۶۳		SERAPH _{None}
(۰,۰۴)	(۰,۰۱)	(۰,۰۳)	(۰,۰۲)	(۰,۰۱)	(۰,۰۱)	(۰,۰۵)	(۰,۱۱)	(۰,۱۰)	(۰,۰۹)		
۳۶,۲۱	۲۸,۶۹	۱۶,۴۵	۳,۹۸	۶,۱۶	۳,۰۶	۲۷,۲۵	۲۵,۵۹	۲۹,۵۱	۲۴,۷۵		SERAPH _{Post}
(۰,۰۴)	(۰,۰۲)	(۰,۰۴)	(۰,۰۲)	(۰,۰۱)	(۰,۰۵)	(۰,۱۰)	(۰,۱۴)	(۰,۱۳)	(۰,۰۹)		
۳۶,۱۳	۲۰,۳۱	۱۹,۰۶	۴,۰۴	۵,۶۱	۱,۹۸	۱۹,۶۲	۲۱,۹۴	۲۳,۵۳	۲۴,۶۳		SERAPH _{proj}
(۰,۰۴)	(۰,۰۱)	(۰,۰۳)	(۰,۰۲)	(۰,۰۱)	(۰,۰۱)	(۰,۰۵)	(۰,۱۱)	(۰,۱۱)	(۰,۰۹)		
۳۶,۲۱	۲۸,۷۹	۱۶,۶۶	۳,۹۸	۶,۱۶	۴,۱۲	۲۷,۴۹	۲۵,۵۹	۳۱,۲۴	۲۴,۷۱		SERAPH _{Hyper}
(۰,۰۴)	(۰,۰۱)	(۰,۰۴)	(۰,۰۲)	(۰,۰۱)	(۰,۰۸)	(۰,۱۱)	(۰,۱۴)	(۰,۱۳)	(۰,۰۹)		
۴۵,۶۴	۲۰,۳۲	۱۲,۸۱	۴,۱۸	۵,۶۱	۱,۹۹	۱۹,۵۸	۲۰,۵۱	۲۲,۲۳	۲۷,۱۴		SERAPH _{mani}
(۰,۱۲)	(۰,۰۱)	(۰,۰۱)	(۰,۰۲)	(۰,۰۱)	(۰,۰۱)	(۰,۰۵)	(۰,۱۱)	(۰,۱۰)	(۰,۱۲)		
۳۵,۵۲	۲۰,۹۳	۲۰,۶۶	۴,۱۴	۵,۵۰	۱,۹۷	۵۳,۸۵	۱۶,۱۹	۲۰,۱۲	۲۲,۴۰		SERAPH _{Post+mani}
(۰,۰۴)	(۰,۰۲)	(۰,۰۱)	(۰,۰۲)	(۰,۰۱)	(۰,۰۱)	(۰,۰۶)	(۰,۰۹)	(۰,۰۹)	(۰,۱۰)		
۴۵,۶۴	۲۰,۳۲	۱۲,۸۲	۴,۱۸	۵,۶۱	۱,۹۹	۱۹,۶۳	۲۰,۱۵	۲۳,۲۱	۲۷,۱۹		SERAPH _{proj+mani}
(۰,۱۲)	(۰,۰۱)	(۰,۰۱)	(۰,۰۲)	(۰,۰۱)	(۰,۰۱)	(۰,۰۵)	(۰,۱۱)	(۰,۱۱)	(۰,۱۲)		
۳۵,۵۲	۲۱,۰۵	۲۰,۶۶	۴,۱۵	۵,۴۸	۱,۹۷	۵۳,۹۶	۱۶,۲۱	۲۰,۰۳	۲۳,۹۶		SERAPH _{Hyper+mani}
(۰,۰۴)	(۰,۰۲)	(۰,۰۱)	(۰,۰۲)	(۰,۰۱)	(۰,۰۱)	(۰,۰۶)	(۰,۰۹)	(۰,۰۸)	(۰,۰۹)		

جدول ۳: مقایسه درصد خطا و تعداد داده‌های برچسب‌دار رویکرد پیشنهادی با رویکرد BOOTSTRAPPING استاندارد.

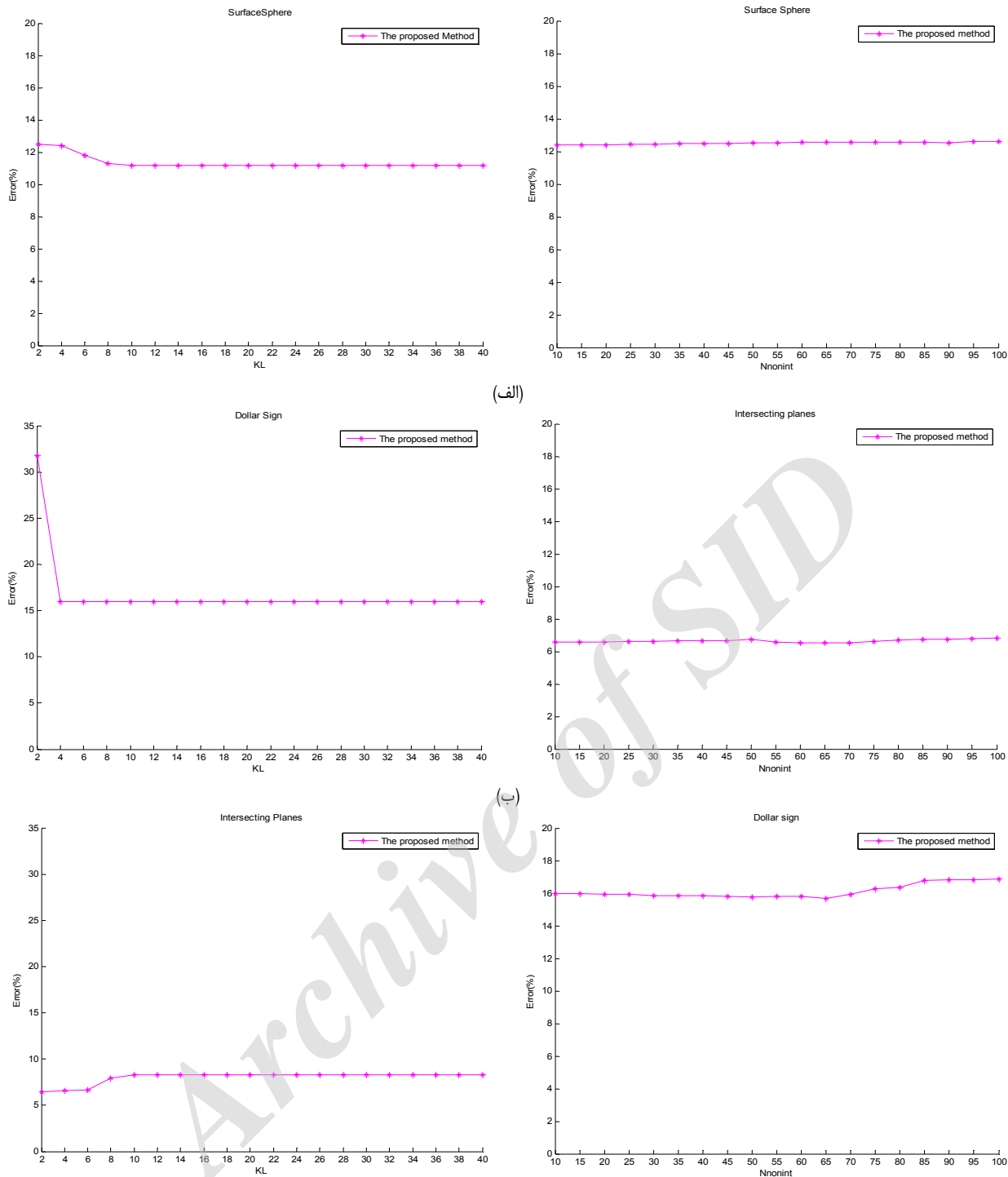
تعداد داده‌های برچسب‌دار				خطا				روش پیشنهادی	مجموعه داده
MRBS بعد از سومین تکرار	MRBS بعد از دومین تکرار	MRBS بعد از اولین تکرار	MR و روش پیشنهادی	MRBS بعد از سومین تکرار	MRBS بعد از دومین تکرار	MRBS بعد از اولین تکرار	MR		
۵۶۵	۴۹۱	۱۵۲	۱۰	۵۶,۲۶	۳۳,۳۳	۴۷,۶۰	۱۰,۴۲	۱۸,۴۶	DS
-	-	-	۲۰	-	-	-	۱۷,۹۵	۱۴,۷۵	SS
-	-	-	۲۰	-	-	-	۶,۲۰	۹,۱۴	TP۱
-	۴۶۹	۱۹۳	۱۲۰	-	۱۰۰,۰۰	۶۶,۶۶	۲۲,۲۵	۱۴,۵۴	TP۲
۳۱۵	۳۰۶	۲۸۴	۱۰۰	۵۶,۱۰	۶۰,۷۰	۶۶,۲۰	۲۴,۵۰	۲,۹۰	TS۱
۵۹۷	۵۲۲	۳۸۴	۲۰۰	۵۶,۷۰	۵۹,۰۰	۶۱,۳۰	۱۶,۸۰	۶,۸۳	TS۲
۵۰۱	۴۶۳	۳۶۶	۲۵۶	۷۲,۳۲	۷۲,۳۲	۷۵,۱۳	۲۴,۹۸	۶,۹۴	RCC
-	۴۹۵	۹۶	۷۰	-	۱۰۰,۰۰	۸۵,۷۱	۱۵,۴۳	۶,۸۳	FS
۱۳۰۷	۱۲۶۸	۱۱۴۲	۲۰۰	۳۳,۳۰	۳۳,۴۵	۳۳,۵۵	۳,۵۵	۲,۵۶	RP
۱۱۷۰	۴۷۶	۲۰۰	۲۰۰	-	۱۰۰,۰۰	۷۵,۰۰	۱۴,۹۵	۲,۵۶	DPR

تی-آستیوندت با سطح معناداری ۰/۰۵ نشان می‌دهد که روش پیشنهادی از سایر روش‌ها بهتر است.

۴-۱ بحث

همان‌طور که نتایج به دست آمده در جدول ۲ نشان می‌دهد با وجود موفقیت MR در بسیاری از داده‌ها، این روش روی داده‌هایی که روی مینیفلدهای متقاطع قرار دارند به واسطه انتشار برچسب اشتباه به نتایج

نتیجه آزمایش روی ۱۰ مجموعه داده مصنوعی استاندارد آمده است و شکل ۲ مجموعه داده‌ها را نشان می‌دهد. در این مجموعه داده‌ها، $w = 0.5$ ، $n = 2000$ و 10% داده‌ها برچسب‌دار است و تعداد داده‌های هر دو دسته با یکدیگر برابر هستند. عرض هسته در MR برابر 0.5 در نظر گرفته شده است. میانگین نرخ درصد خطا به ازای ۳۰ بار اجرای مستقل و انحراف معیار نتایج به دست آمده از اعمال یک دسته‌بند ۱NN با استفاده از معیارهای فاصله بیان‌شده در جدول ۲ آمده است. آزمون



(الف)

(ب)

(ج)

شکل ۳: نمودار درصد خطا به ازای تغییر مقدار پارامترهای n_{nonint} (ستون سمت راست) و k_l (ستون سمت چپ) روی مجموعه داده‌های مصنوعی، (الف) Surface-sphere، (ب) Intersecting planes و (ج) Dollar sign.

هرمرحله داده‌های دسته‌بندی شده با اطمینان بالا به داده‌های برچسب‌دار در دسترس اضافه می‌شود. MR به عنوان دسته‌بند اولیه انتخاب شده و دقت این دسته‌بند در چهار تکرار اول و نیز تعداد داده‌های برچسب‌دار در این تکرارها در جدول ۳ آمده است. همان گونه که مشخص است تکرار عمل اضافه‌نمودن داده‌هایی که با اطمینان بالایی درست دسته‌بندی شده‌اند در اکثر موارد موجب خطای بسیار زیادی شده و هر چند تعداد داده‌های برچسب‌دار در گام‌های دو به بعد زیاد می‌شود اما خطای دسته‌بندی افزایش می‌یابد.

خوبی نمی‌رسد. همچنین در برابر تغییر داده‌های برچسب‌دار نسبت به سایر روش‌ها ناپایدارتر (انحراف معیار بالاتر) است و فاصله اقلیدسی نیز به واسطه قرارگیری داده‌های روی منیفلد و نیز اشتراک منیفلدها کارایی لازم را ندارد. روش پیشنهادی بهتر از Seraph و RSSML که تقاطع منیفلدها را در نظر نمی‌گیرند عمل می‌کند.

در ادامه، جهت تحلیل دقیق‌تر، دقت روش پیشنهادی با روش معمول bootstrapping¹ (MRBS) مقایسه شده است. در روش معمول در

1. MR Bootstrapping

- [11] N. Pitelis, C. Russell, and L. Agapito, "Learning a manifold as an atlas," in *Proc. IEEE Conf. on Computer Vision and Pattern Recognition, CVPR'13*, pp. 1642-1649, Portland, OR, USA, 23-28 Jun. 2013.
- [12] X. Wang, T. Peter, and A. F. Mark, "Multiple manifolds learning framework based on hierarchical mixture density model," in *Proc. Joint European Conf. on Machine Learning and Knowledge Discovery in Databases*, pp. 566-581, Antwerp, Belgium, 15-19 Sept. 2008.
- [13] K. Q. Weinberger and L. K. Saul, "Distance metric learning for large margin nearest neighbor classification," *The J. of Machine Learning Research*, vol. 10, pp. 207-244, 12 Jan. 2009.
- [14] M. Der and L. K. Saul, "Latent coincidence analysis: a hidden variable model for distance metric learning," in *Proc. 25th Int. Conf. on Neural Information Processing Systems*, vol. 2, pp. 3230-3238, Lake Tahoe, NV, US, 3-6 Dec. 2012.
- [15] S. C. Hoi, W. Liu, and S. F. Chang, "Semi-supervised distance metric learning for collaborative image retrieval and clustering," *ACM Trans. on Multimedia Computing, Communications, and Applications*, vol. 6, no. 3, pp. 18-43, Aug. 2010.
- [16] W. Liu, S. Ma, D. Tao, J. Liu, and P. Liu, "Semi-supervised sparse metric learning using alternating linearization optimization," in *Proc. Int. 16th Conf. on Knowledge Discovery and Data Mining, ACM SIGKDD'10*, pp. 1139-1148, Washington, DC, US, 25-28 Jul. 2010.
- [17] M. Baghshah and S. Shouraki, "Semi-supervised metric learning using pairwise constraints," in *Proc. of the 21st Int. Joint Conf. on Artificial Intelligence*, pp. 1217-1222, Pasadena, CA, US, 11-17 Jul. 2009.
- [18] G. Niu, B. Dai, M. Yamada, and M. Sugiyama, "Information-theoretic semisupervised metric learning via entropy regularization," *Neural Computation*, vol. 26, no. 8, pp. 1717-1762, Aug. 2012.
- [19] Q. Wang, P. C. Yuen, and G. Feng, "Semi-supervised metric learning via topology preserving multiple semi-supervised assumptions," *Pattern Recognition*, vol. 46, no. 9, pp. 2576-2587, Sept. 2013.
- [20] R. O. Duda, P. E. Hart, and D. G. Stork, *Pattern Classification*, John Wiley and Sons, 2001.
- [21] L. Yang and R. Jin, *Distance Metric Learning: A Comprehensive Survey*, Michigan State University, 2006.
- [22] M. Belkin, Q. Que, Y. Wang, and X. Zhou, "Toward understanding complex spaces: graph laplacians on manifolds with singularities and boundaries," in *Proc. 25th Annual Conf. on Learning Theory, COLT'12*, pp. 1-26, Edinburgh, Scotland, 2012.
- [23] Z. Karimi and S. Shiry Ghidary, "Semi-Supervised Classification in Stratified Spaces by Considering Non-interior Points Using Laplacian Behavior," *Neurocomputing*, vol. 239, pp. 223-231, 24 May 2017.
- [24] L. Vandenberghe and S. Boyd, "Semidefinite programming," *SIAM Review*, vol. 38, no. 1, pp. 49-95, Mar. 1996.
- [25] M. Belkin, P. Niyogi, and V. Sindhwani, "Manifold regularization: a geometric framework for learning from labeled and unlabeled examples," *J. of Machine Learning Research*, vol. 7, pp. 2399-2434, Nov. 2006.

زهرة کریمی در رشته مهندسی کامپیوتر در گرایش نرم‌افزار در سال ۱۳۸۵ از دانشگاه شهید بهشتی فارغ‌التحصیل شد و مدرک کارشناسی ارشد خود را در سال ۱۳۸۸ از دانشگاه صنعتی شریف در همان رشته اخذ کرد. او دکتری خود را در رشته هوش مصنوعی در دانشگاه صنعتی امیرکبیر دریافت نموده است. زمینه‌های پژوهش وی شامل یادگیری ماشین، داده‌کاوی و یادگیری مبتنی بر منیفلد است.

سعید شیری قیداری استاد دانشکده علوم کامپیوتر دانشگاه صنعتی امیرکبیر هستند. ایشان مدرک کارشناسی خود را در رشته مهندسی الکترونیک و کارشناسی ارشد خود را در مهندسی کامپیوتر از دانشگاه صنعتی امیرکبیر اخذ نموده است. وی در سال ۱۳۸۱ دکتری خود را از دانشگاه کوبه ژاپن دریافت نموده و از سال ۱۳۸۳ استادیار دانشگاه صنعتی امیرکبیر است. زمینه‌های تحقیقاتی مورد علاقه وی رباتیک، بینایی ماشین، علوم ساختی و مدل‌سازی مغز است.

روح‌اله رضانی در رشته آمار در سال ۱۳۸۳ از دانشگاه بین‌المللی امام خمینی (ره) فارغ‌التحصیل شد و مدرک کارشناسی ارشد خود را در سال ۱۳۸۵ از دانشگاه صنعتی امیرکبیر در همان رشته اخذ کرد. زمینه‌های پژوهش وی شامل داده‌کاوی، پژوهش‌های آماری، کنترل کیفیت آماری و قابلیت اطمینان است. وی، هم اکنون عضو هیأت علمی دانشگاه دامغان است.

۴-۲ بررسی تأثیر تغییر مقدار پارامترها

در این بخش، تأثیر انتخاب پارامترها را روی دقت روش پیشنهادی بررسی می‌کنیم. دو پارامتر اصلی تعداد نقاط غیر داخلی و k_l در روش پیشنهادی وجود دارد.

میانگین درصد خطای روش پیشنهادی روی ۳۰ اجرای مستقل به ازای مقادیر مختلف پارامترها روی چند مجموعه داده بررسی شده است. در هر زمان، مقدار یک پارامتر تغییر داده شده و مقدار سایر پارامترها برابر با مقدار پیش‌فرض تنظیم شده است. ستون سمت راست در شکل ۳ نمودار تغییر درصد خطا به ازای تغییر مقدار پارامتر $n_{non-int}$ از ۱۰ تا ۱۰۰ را نشان می‌دهد. از آنجا که در برخی تنظیمات تعداد نقاط با مقادیر لاپلاسین تابع غیر صفر کمتر از مقدار مشخص شده $n_{non-int}$ است، بیشینه مقدار $n_{non-int}$ در نظر گرفته شده است. نمودار ارائه‌شده حاکی از آن است که در تمام مجموعه داده‌ها برای مقادیر زیادی از $n_{non-int} \ll n$ کارایی رویکرد پیشنهادی مستقل از مقدار $n_{non-int}$ است. نمودار تغییر مقدار درصد خطا به ازای تغییر مقدار k_l در ستون سمت چپ شکل ۳ آمده است. روش پیشنهادی در زمانی که مقدار k_l خیلی کوچک نباشد مستقل از این پارامتر است.

۵- نتیجه‌گیری

در این مقاله، یک رویکرد یادگیری متریک نیمه‌نظارتی با فرض قرارگیری داده در فضای لایه‌ای ارائه شد. با اعمال فرض همواربودن فقط در نواحی داخلی منیفلدها، رویکردهای مبتنی بر منیفلد موجود جهت یادگیری متریک در فضای لایه‌ای توسعه داده شده است. همچنین رویکرد جدیدی جهت اعمال فرض همواربودن به نزدیک‌ترین همسایه داده‌های برچسب‌دار تعریف شده است. موارد مذکور بر مبنای بهره‌گیری دقیق‌تر از دانش پیشین قرارگیری داده‌ها روی چند منیفلد بوده و منجر به کارایی بهتر رویکرد یادگیری متریک پیشنهادی در مقایسه با روش‌های مشابه شده است.

مراجع

- [1] J. Lu, X. Zhou, Y. P. Tan, Y. Shang, and J. Zhou, "Neighborhood repulsed metric learning for kinship verification," *IEEE Trans. on Pattern Analysis and Machine Intelligence*, vol. 36, no. 2, pp. 331-345, Feb. 2014.
- [2] J. Yu, M. Wang, and D. Tao, "Semisupervised multiview distance metric learning for cartoon synthesis," *IEEE Trans. on Image Processing*, vol. 21, no. 11, pp. 4636-4648, Nov. 2012.
- [3] B. Kulis, "Metric learning: a survey," *Foundations and Trends in Machine Learning*, vol. 5, no. 4, pp. 287-364, 2012.
- [4] Y. Ying and P. Li, "Distance metric learning with eigenvalue optimization," *J. of Machine Learning Research*, vol. 13, no. 1, pp. 1-26, Jan. 2012.
- [5] L. K. Saul and S. T. Roweis, "Think globally, fit locally: unsupervised learning of low dimensional manifolds," *the J. of Machine Learning Research*, vol. 4, pp. 119-155, 12 Jan. 2003.
- [6] J. B. Tenenbaum, V. De Silva, and J. C. Langford, "A global geometric framework for nonlinear dimensionality reduction," *Science*, vol. 290, no. 5500, pp. 2319-2323, 22 Dec. 2000.
- [7] Z. Zhang and H. Zha, "Principal manifolds and nonlinear dimensionality reduction via tangent space alignment," *SIAM J. on Scientific Computing*, vol. 26, no. 1, pp. 313-338, 2004.
- [8] M. Belkin and P. Niyogi, "Laplacian eigenmaps for dimensionality reduction and data representation," *Neural Computation*, vol. 15, no. 6, pp. 1373-1396, 1 Jun. 2003.
- [9] J. Wang, X. Sun, and J. Du, "Local tangent space alignment via nuclear norm regularization for incomplete data," *Neurocomputing*, vol. 273, pp. 141-151, 17 Jan. 2018.
- [10] M. A. Carreira-Perpin and Z. Lu, "Manifold learning and missing data recovery through unsupervised regression," in *Proc. IEEE 11th Int. Conf. on Data Mining, ICDM'11*, pp. 1014-1019, Vancouver, Canada, 11-14 Dec. 2011.