

تاریخ دریافت مقاله: ۹۷/۱۲/۱۰

تاریخ پذیرش مقاله: ۹۸/۰۸/۰۳

انتخاب نمونه‌های آموزشی بهینه بر اساس معیارهای فاصله برای آموزش رده‌بندی احساسات

شیوا نوری سرای

دانشجوی کارشناسی ارشد، دانشکده مهندسی فناوری اطلاعات و کامپیوتر- دانشگاه صنعتی ارومیه - ارومیه - ایران
پست الکترونیکی: shivanoorisaray@it.uut.ac.ir

جعفر طهمورث نژاد*

استادیار دانشکده مهندسی فناوری اطلاعات و کامپیوتر- دانشگاه صنعتی ارومیه - ارومیه - ایران
پست الکترونیکی: j.tahmores@it.uut.ac.ir

چکیده

در این مقاله یک ترکیب خطی از معیارهای فاصله بین توزیع دامنه‌های منبع و هدف است که بهترین دامنه منبع را برای یادگیری رده‌بندی انتخاب می‌کند. روش پیشنهادی بر روی مجموعه داده‌های همگن^۱ و ناهمگن^۲ ارزیابی شده است. همان‌طور که نتایج نشان می‌دهد، مدل پیشنهادی، در مجموعه داده همگن با احتمال ۴۷،۱ درصد (۵،۹ درصد در مدل تصادفی) و در مجموعه داده ناهمگن با احتمال ۲۳،۱ درصد (۸،۳ درصد در مدل تصادفی) می‌تواند دامنه منبع صحیح را انتخاب کند که حاکی از بهبود چشمگیر عملکرد مدل پیشنهادی نسبت به مدل تصادفی در انتخاب دامنه منبع صحیح است.

واژه‌های کلیدی: رده‌بندی احساسات، معیار فاصله، یادگیری انتقالی، انتقال دانش

۱- مقدمه

در رده‌بندی احساسات، مدلی بر روی یک دامنه منبع

3- homogeneous
4- heterogeneous

افزایش چشمگیر دسترس‌پذیری به نظرها و توصیه‌های برخط باعث می‌شود رده‌بندی احساسات در متون کوتاه یکی از موضوع‌های جالب توجه در تحقیقات علمی و صنعتی باشد. در زمینه رده‌بندی احساسات، اصطلاحات به کار برده شده در دامنه‌های مختلف ممکن است متفاوت باشند. در نتیجه مدلی که با داده‌های برچسب‌دار آموزشی (دامنه منبع) یادگیری می‌شود ممکن است عملکرد خوبی در برچسب‌گذاری داده‌های آزمایشی (دامنه هدف) نداشته باشد. یادگیری انتقالی^۱ و انطباق دامنه^۲ دو راه حل مفید برای مواجهه با این مشکل هستند. یادگیری انتقالی و انطباق دامنه، توزیع‌های دامنه منبع و هدف را به هم نزدیک می‌کنند تا عملکرد رده‌بندی در دامنه هدف بهبود یابد اما نکته حائز اهمیت این است که کدام دامنه از مجموعه دامنه‌های نامزد به عنوان دامنه منبع انتخاب شود. روش پیشنهادی

* sentiment classification

1- transfer learning
2- domain adaptation

* نویسنده مسئول

با داده‌های برچسب‌دار آموزش داده می‌شود، سپس مدل برای برچسب‌گذاری داده‌های دامنه هدف به کار می‌رود. در رده‌بندی احساسات، داده‌های دامنه، متن‌های کوتاه هستند که با توجه به محتوای متن‌ها، رده‌بند، برچسب مثبت یا منفی به متن‌ها می‌دهد. رده‌بندی احساسات در زمینه‌های مختلف مانند تجارت الکترونیکی، بازیابی اطلاعات، سیستم‌های پیش‌بینی بازار بورس، پردازش گفتار، سیستم‌های توصیه‌گر به کار می‌رود. یکی از مسائل بحث برانگیز در رده‌بندی احساسات، اختلاف توزیع بین داده‌های دامنه منبع و هدف، یعنی نظراتی که مدل با آن‌ها ساخته می‌شود و نظراتی که مدل با آن‌ها آزموده می‌شود است که عملکرد مدل روی دامنه هدف را تحت تأثیر قرار می‌دهد.

روش‌های یادگیری انتقالی زیادی از جمله انطباق دامنه بصری از طریق یادگیری ویژگی انتقالی [۱] و استخراج ویژگی‌های مشترک برای دامنه‌های چندمنبعی در یادگیری انتقالی [۲] ارائه شده است تا افت عملکرد رده‌بند که ناشی از اختلاف توزیع دامنه‌ها است را کاهش دهد.

اختلاف توزیع دامنه‌های منبع و هدف، ناشی از اختلاف توزیع حاشیه‌ای^۵ و توزیع شرطی^۶ بین دامنه‌های منبع و هدف است. اختلاف توزیع حاشیه‌ای، به صورت اختلاف در احتمال وقوع هر یک از ویژگی‌ها در دامنه‌های منبع و هدف با ویژگی‌های یکسان و اختلاف توزیع شرطی، به صورت اختلاف در احتمال وقوع برچسب‌های متفاوت به ازای داده‌های یکسان از دو دامنه منبع و هدف تعریف می‌شود. روش‌های یادگیری انتقالی، به دو گروه نیمه‌نظارت‌شده^۷ و بدون نظارت^۸ تقسیم می‌شوند. روش‌های یادگیری انتقالی، در صورتی که تمام داده‌های دامنه منبع و قسمت کوچکی از داده‌های دامنه هدف برچسب داشته باشند، به صورت نیمه‌نظارت‌شده و اگر فقط داده‌های دامنه منبع برچسب داشته باشند به صورت بدون نظارت تعریف می‌شوند.

5- marginal distribution
6- conditional distribution
7- semi-supervised
8- unsupervised

روش‌های یادگیری انتقالی، بر کاهش اختلاف توزیع دامنه‌های منبع و هدف و انتقال دانش از دامنه منبع به دامنه هدف تمرکز دارند در حالی که به انتخاب دامنه مناسب به عنوان دامنه منبع، توجهی ندارند. برای انتخاب دامنه منبع می‌توان عواملی مانند موضوع دامنه، هدف از تشکیل دامنه و یا چارچوبی که دامنه در آن تعریف شده است را در نظر گرفت ولی توجه به این عوامل الزاماً بهترین دامنه منبع را به دست نمی‌دهد. در اکثر دیدگاه‌های انتخاب دامنه منبع از معیارهای شباهت دامنه استفاده می‌شود. از جمله معیارهای شباهت، تابع فاصله حداکثر اختلاف میانگین (MMD)^۹ است [۳]. این تابع از کرنل‌های مختلفی برای محاسبه فاصله بین دو دامنه منبع و هدف استفاده می‌کند ولی از معایب این معیار این است که محاسبه ماتریس کرنل برای مجموعه داده‌های بزرگ بسیار هزینه‌بر است. روش پیشنهادی این مقاله، یک روش هوشمند انتخاب دامنه بدون نظارت است که یک یا چند دامنه منبع مناسب را به صورت بدون نظارت از مجموع دامنه‌های در دسترس برای یادگیری رده‌بند انتخاب می‌کند در حالی که از برچسب‌های دامنه هدف استفاده نمی‌کند.

در ادامه، مقاله به صورت زیر سازماندهی شده است. در بخش دوم مقاله، مروری بر روی کارهای مربوطه انجام شده است. در بخش سوم مقاله، مسئله و روش پیشنهادی به تفصیل شرح داده شده است. در بخش چهارم، نتایج ارزیابی عملکرد الگوریتم پیشنهادی در مقایسه با روش‌های دیگر انتخاب دامنه منبع گزارش شده است. در انتها، مقاله با نتیجه‌گیری و ارائه پیشنهادهایی برای ادامه کار در آینده به اتمام رسیده است.

۲. مروری بر کارهای مربوطه

یادگیری انتقالی به دو دسته یادگیری انتقالی همگن و یادگیری انتقالی ناهمگن تقسیم می‌شود [۴]. در یادگیری انتقالی همگن فضای ویژگی دامنه‌ها با هم برابر است و به چهار زیرگروه که به شرح زیر می‌باشند، تقسیم می‌شود:

9- Maximum Mean Discrepancy

● یادگیری انتقالی مبتنی بر نمونه: در این دیدگاه به منظور کاهش اختلاف توزیع‌های بین دامنه‌های منبع و هدف، نمونه‌هایی از دامنه منبع که در کاهش اختلاف توزیع داده‌های منبع و هدف موثر هستند، انتخاب شده و یا وزندهی مجدد می‌شوند. روش انتخاب مرزما^{۱۰} [۵]، از جمله روش‌های مبتنی بر نمونه است. به نمونه‌هایی از دامنه منبع که کمترین اختلاف توزیع را با نمونه‌های دامنه هدف دارند، مرزما گفته می‌شود. در این روش، از معیار حداکثر اختلاف میانگین (MMD) برای محاسبه اختلاف توزیع نمونه‌های دامنه منبع با نمونه‌های دامنه هدف بهره گرفته می‌شود و نمونه‌هایی از دامنه منبع که دارای کمترین اختلاف توزیع با نمونه‌های دامنه هدف باشند، وزن بیشتری می‌گیرند. ایراد اصلی در روش انتخاب مرزما این است که به ویژگی‌های مختص هر دامنه توجهی نمی‌شود.

● یادگیری انتقالی مبتنی بر ویژگی: در این دیدگاه به منظور کاهش اختلاف توزیع بین دامنه‌های منبع و هدف، داده‌ها به فضای ویژگی مشترک با اختلاف توزیع کمتر نگاشت می‌شوند. JDA [۶] و GFK [۷] را از جمله روش‌های این دیدگاه می‌توان نام برد. در روش JDA، هدف یافتن ماتریس نگاشتی است که نمونه‌های منبع و هدف را به زیرفضای مشترکی با حداقل اختلاف توزیع شرطی و حاشیه‌ای بین دامنه‌ها نگاشت کند. مشکل اصلی روش JDA عدم حفظ ساختار داده‌ها می‌باشد. روش GFK با استخراج ویژگی‌های پنهان مشترک بین دامنه‌های منبع و هدف، اختلاف توزیع حاشیه‌ای بین دو دامنه را کاهش می‌دهد. روش GFK، برای کاهش اختلاف توزیع حاشیه‌ای بین دامنه‌های منبع و هدف، ابعاد فضای اصلی داده‌ها را کاهش می‌دهد که به علت کاهش بیش از حد ابعاد، ساختار اصلی داده‌ها حفظ نمی‌شود.

● یادگیری انتقالی مبتنی بر پارامتر / مدل: در این دیدگاه، پارامترهای مدل آموزش دیده روی نمونه‌های منبع با ساختار نمونه‌های دامنه هدف تطبیق داده می‌شود.

- 10- landmark
11- Joint Domain Adaptation
12- Geodesic Flow Kernel

تنظیم‌پذیری تطبیقی برای یادگیری انتقالی^{۱۳} [۸] یک روش تطبیق دامنه بدون نظارت است که از طریق کاهش اختلاف توزیع هندسی و آماری بین دامنه‌های منبع و هدف و همچنین کاهش خطای رده‌بند در دامنه منبع، یک رده‌بند تطبیق یافته بین دامنه‌های منبع و هدف ایجاد می‌کند.

به طور کلی، مدل‌های مختلفی برای یادگیری انتقالی استفاده می‌شود، از جمله این مدل‌ها می‌توان به ماشین بردار پشتیبان (SVM)^{۱۴} [۹]، شبکه‌های عصبی احتمالی (PNN)^{۱۵} [۱۰]، مدل مخلوطی گوسی (GMM)^{۱۶} [۱۱] و مدل حداکثر آنتروپی (ME)^{۱۷} [۱۲] اشاره کرد. مدل ماشین بردار پشتیبان، نتایج خیلی خوبی در رده‌بندی داده‌ها به دست می‌دهد. در مدل ماشین بردار پشتیبان سعی بر پیدا کردن بهترین ابرصفحه بین رده‌ها برای رده‌بندی بهتر داده‌ها است. این مدل بیشتر برای رده‌بندی احساسات درباره فیلم‌ها استفاده می‌شود. مدل شبکه عصبی احتمالی برای رده‌بندی داده‌های با بعد زیاد مفید است. در مدل‌های مخروطی گوسی، مولفه‌های گوسی برخی از ویژگی‌هایی که در برچسب‌گذاری مفیدتر هستند را نشان می‌دهند و می‌تواند انواع مختلف توزیع‌ها را مدل کند. مدل حداکثر آنتروپی، مدل مبتنی بر ویژگی است که توزیع هر رده را پیدا می‌کند و سبب کلمات^{۱۸} دوتایی و عبارات را می‌توان به عنوان ویژگی به مدل حداکثر آنتروپی داد بدون این‌که نگران همپوشانی ویژگی‌ها باشیم. با این حال، برای رده‌بندی می‌توان ترکیبی از رده‌بندهای فوق را نیز استفاده کرد.

● یادگیری انتقالی مبتنی بر رابطه: این دیدگاه شامل روش‌های یادگیری الگوهای ساختار گرامری و جمله‌ای دامنه‌ها و یافتن رابطه بین دامنه‌های منبع و هدف است. RAP^{۱۹} [۱۳] یک روش انطباق دامنه تحت نظارت است که برای استخراج دانش از تحلیل احساسات استفاده می‌کند

- 13- Adaptation Regularization for Transfer Learning
14- Support Vector Machine
15- Probability Neural Network
16- Gaussian Mixture Model
17- Maximum Entropy
18- Bag of Words
19- Relational Adaptive bootstraPping

و به این صورت عمل می‌کند که ابتدا کلمات احساسی و مبحث دامنه هدف را مشخص می‌کند سپس آن‌ها را توسعه می‌دهد و رابطه بین مبحث و احساسات را تعیین می‌کند.

در یادگیری انتقالی ناهمگن فضای ویژگی دامنه‌های منبع و هدف متفاوت هستند. دو دیدگاه اصلی برای حل مسئله فضای ویژگی ناهمگن وجود دارد که به شرح زیر است:

یادگیری انتقالی مبتنی بر ویژگی متقارن: در این دیدگاه، دامنه‌های منبع و هدف به صورت جداگانه به فضای ویژگی پنهان انتقال داده می‌شوند. روش DAMA [۱۴] از این دیدگاه استفاده می‌کند. در این روش به منظور ایجاد فضای پنهان و اتصال فضای ویژگی دامنه‌ها از تابع نگاشت متفاوت برای هر دامنه استفاده می‌شود. همچنین در روش DAMA، از برچسب نمونه‌ها به جای ویژگی‌های معادل نمونه‌ها برای انطباق هندسی استفاده می‌شود.

● یادگیری انتقالی مبتنی بر ویژگی نامتقارن: در این دیدگاه فضای ویژگی دامنه منبع به فضای ویژگی دامنه هدف انتقال داده می‌شود. روش HDP [۱۵]، یک روش تشخیصی ناقص بر اساس تطبیق معیار با استفاده از تحلیل آماری است. انطباق معیار برای HDP به توزیع داده‌های به اندازه کافی بزرگ نیاز دارد. مشکل عمده این روش این است که حداقل ۵۰ نمونه برای انتقال دانش و تشخیص عیب مورد نیاز است.

انتخاب دامنه منبع بهینه برای آموزش رده‌بند، می‌تواند عملکرد را در روش‌های انطباق دامنه و یادگیری انتقالی به طور چشمگیری بهبود دهد. در واقع دیدگاه‌های یادگیری انتقالی با روش‌های مختلف، توزیع دامنه‌های منبع و هدف را به هم نزدیک می‌کنند تا رده‌بند روی دامنه‌ای نزدیک به دامنه هدف آموزش داده شود. بدین ترتیب اگر از ابتدا یک دامنه آموزشی بهینه برای یادگیری رده‌بند انتخاب کنیم و سپس با استفاده از روش‌های یادگیری انتقالی و انطباق دامنه اختلاف توزیع دامنه منبع بهینه و دامنه هدف را کم

کنیم بهبود چشمگیری در عملکرد رده‌بند در رده‌بندی داده‌های دامنه هدف به وجود خواهد آمد. بنابراین انتخاب دامنه منبع مناسب به عنوان پیش مرحله روش‌های یادگیری انتقالی، می‌تواند اتلاف تطبیق^{۲۲} را بهبود دهد.

در این مقاله، روش CCEM^{۲۳} ایده جدیدی برای انتخاب دامنه منبع پیشنهاد می‌کند که نیاز به برچسب داده‌های دامنه هدف ندارد. در واقع این روش، ترکیب خطی از معیارهای فاصله است که وزن‌ها در ترکیب خطی، با یادگیری روی دامنه‌های منبع در دسترس به دست می‌آیند و به این صورت مدل انتخاب دامنه منبع یاد گرفته می‌شود.

۳. بیان مسئله

در رده‌بندی احساسات، داده‌های دامنه، اسنادی هستند که هر سند با یک متغیر تصادفی X و برچسب آن سند با متغیر تصادفی Y نشان داده می‌شود که P و \bar{P} دو توزیع روی $X \times Y$ هستند و توزیع داده‌های دامنه‌های منبع و هدف را نشان می‌دهند و توزیع حاشیه‌ای دامنه‌های منبع و هدف روی X به ترتیب با P_X و \bar{P}_X مشخص می‌شود. یک رده‌بند برای هر کدام از متغیرهای تصادفی یک برچسب تعیین می‌کند و بهترین رده‌بند آن است که کمترین خطا را در تعیین برچسب داشته باشد که با بیان ریاضی به صورت رابطه (۱) تعریف می‌شود:

$$h_p^* := \arg \min_h \text{Prob}(h(x) \neq y | (x, y) \sim P) \quad (1)$$

در رابطه (۱)، h_p^* بهترین رده‌بند برای رده‌بندی داده‌های دامنه هدف است که از بین طیف وسیعی از رده‌بندهای $h(x)$ انتخاب می‌شود.

دوم مفهوم خطای رده‌بندی داخلی دامنه $\xi(P)$ و خطای رده‌بندی بین دامنه‌ای $\xi(P, \bar{P})$ به این صورت تعریف می‌شود که خطای رده‌بندی داخلی دامنه، احتمال رده‌بندی اشتباه در دامنه‌ای است که رده‌بند روی آن دامنه آموزش دیده است در حالی که خطای رده‌بندی بین دامنه‌ای به صورت احتمال رده‌بندی اشتباه در دامنه هدف با استفاده

22- adaptation loss

23- linear combination of (Chi2),(CMD),(EMD),(MKL)

20- Domain Adaptation Manifold Alignment

21- Heterogeneous Defect Prediction

از رده‌بند یاد گرفته شده روی دامنه منبع تعیین می‌شود. رابطه‌های (۲) و (۳) به ترتیب خطای رده‌بندی داخلی دامنه و خطای رده‌بندی بین دامنه‌ای را نشان می‌دهند:

$$\xi(P) := \text{Prob}(h_p^* \neq y | x, y \sim P) \quad (۲)$$

$$\xi(P, \bar{P}) := \text{Prob}(h_p^* \neq y | x, y \sim \bar{P}) \quad (۳)$$

با توجه به اختلاف توزیع‌های دامنه منبع و هدف همواره $\xi(P, \bar{P})$ بزرگتر از $\xi(P)$ است [۱۶]. به اختلاف بین خطای رده‌بندی داخلی دامنه و خطای رده‌بندی بین دامنه‌ای، اتلاف تطبیق گفته می‌شود. بهترین رده‌بند آن است که خطای رده‌بندی بین دامنه‌ای کمی داشته باشد و در واقع هدف، انتخاب دامنه منبعی (P^*) است که خطای رده‌بندی بین دامنه‌ای را کم کند. بدین ترتیب اگر \mathcal{P} را مجموعه دامنه‌های منبع در دسترس تعریف کنیم هدف مقاله به صورت رابطه (۴) تعریف می‌شود:

$$P^* := \arg \min_{P \in \mathcal{P}} \xi(P, \bar{P}) \quad (۴)$$

مسئله‌ای که در پیدا کردن بهترین دامنه منبع با استفاده از رابطه (۴) به وجود می‌آید این است که برای پیدا کردن خطای رده‌بندی بین دامنه‌ای به داده‌های برچسب‌دار دامنه هدف نیاز است در حالی که دامنه هدف، داده‌های بدون برچسب داشته و فقط توزیع حاشیه‌ای دامنه هدف در دسترس است.

۱.۳ روش پیشنهادی CCEM

در این مقاله برای حل مشکل عدم وجود برچسب در دامنه هدف روشی به نام CCEM معرفی می‌شود که ترکیب خطی از چهار معیار فاصله و خطای رده‌بندی داخلی دامنه $\xi(P)$ است که معیارهای فاصله، فاصله آماری بین توزیع حاشیه‌ای دامنه‌های منبع و هدف را محاسبه می‌کند. برای یافتن دامنه منبع بهینه فرض می‌شود دامنه منبعی که کمترین خطای رده‌بندی داخلی دامنه و کمترین فاصله را با دامنه هدف دارد $\hat{d}(P, \bar{P}_x)$ ، کمترین خطای رده‌بندی بین دامنه‌ای $\xi(P, \bar{P})$ را نیز دارد و در واقع خطای رده‌بندی بین دامنه‌ای را می‌توان با خطای رده‌بندی داخلی دامنه منبع و فاصله بین دامنه منبع و هدف تقریب زد که به صورت

رابطه (۵) نشان داده می‌شود و دامنه، با خطای رده‌بندی داخلی و فاصله تا دامنه هدف کمتر، برای آموزش رده‌بند انتخاب می‌شود.

$$\xi(P, \bar{P}) \approx \hat{d}(P, \bar{P}_x) + \xi(P) \quad (۵)$$

فرض می‌کنیم \mathcal{D} مجموعه‌ای از توابع فاصله است و \bar{P} مجموعه‌ای از دامنه‌ها است و به دلیل این که در آزمایش‌های دامنه هدف جداگانه نداریم، در هر تکرار از آزمایش یکی از دامنه‌ها به عنوان دامنه هدف و بقیه دامنه‌ها به عنوان دامنه‌های منبع در مجموعه نامزد در نظر گرفته می‌شوند. برای انتخاب بهترین معیار فاصله از رابطه (۶)

$$\hat{d} := \arg \min_{d \in \mathcal{D}} \sum_{P \in \bar{\mathcal{P}}} |\xi(P, \bar{P}) - \hat{d}(P, \bar{P}_x)| \quad (۶)$$

برای انتخاب مجموعه‌ای از معیارهای فاصله، از بردار M با n معیار فاصله استفاده می‌کنیم که هر کدام از M_i ها یک معیار فاصله با ضریب $\alpha_i \in \mathbb{R}^+$ است. خطای رده‌بندی دامنه منبع نشان‌دهنده مناسب بودن منبع برای یادگیری رده‌بند است. اگر خطای دامنه زیاد باشد دامنه ارزش کمتری برای انتخاب دارد به همین دلیل خطای رده‌بندی داخلی هر دامنه را محاسبه می‌کنیم و به ترکیب خطی معیارهای فاصله اضافه می‌کنیم. یافتن ضرایب بهینه برای معیارهای فاصله به صورت رابطه (۷) تعریف می‌شود.

$$\hat{d} := \arg \min_{\alpha \geq 0} \sum_{P \in \bar{\mathcal{P}}} |\xi(P, \bar{P}) - \alpha_M(P, \bar{P}_x)| \quad (۷)$$

بنابراین بدون نیاز به داده‌های برچسب‌دار دامنه هدف می‌توان بهترین دامنه منبع را برای آموزش رده‌بند انتخاب کرد زیرا برای محاسبه معیارهای فاصله فقط توزیع حاشیه‌ای دامنه‌های منبع و هدف استفاده می‌شود بنابراین تعریف رابطه (۴) به صورت رابطه (۸) تغییر می‌کند که هدف مقاله است:

$$P^* := \arg \min_{P \in \mathcal{P}} \hat{d}(P, \bar{P}_x) \quad (۸)$$

۲.۳ معیارهای فاصله

از آنجایی که توزیع دامنه‌های منبع و هدف اختلاف دارند در بیشتر رویکردها برای انتخاب دامنه منبع بهینه، که

P_x و \bar{P}_x توزیع‌های حاشیه‌ای دامنه‌های منبع و هدف در فاصله فشرده $[a, b]^N$ هستند و $E(p_x)$ و $c_k(p_x)$ به ترتیب امید ریاضی و بردار گشتاور مرکزی از درجه k هستند. $\| \cdot \|_2$ نرم مربع است.

(۱۲)

$$CMD(p_x, \bar{p}_x) = \frac{1}{|b-a|} \|E(p_x) - E(\bar{p}_x)\|_2 + \sum_{k=2}^{\infty} \frac{1}{|b-a|} \|c_k(p_x) - c_k(\bar{p}_x)\|_2$$

$$c_k(p_x) = (E(\prod_{i=1}^N ((p_x)_i - E(p_x)_i)^{r_i})) \quad (13)$$

$$\text{Subject to: } \begin{aligned} r_1 + \dots + r_N &= k, \\ r_1, \dots, r_N &\geq 0 \end{aligned}$$

مدل پیشنهادی از ترکیب خطی چهار معیار ذکر شده، خطای انحراف و خطای رده‌بندی داخلی دامنه منبع تشکیل شده است که به صورت رابطه (۱۴) تعریف می‌شود.

$$d(P_x, \bar{P}_x) = \alpha_1 \text{chi2}(P_x, \bar{P}_x) + \alpha_2 \text{CMD}(P_x, \bar{P}_x) + (14)$$

$$+ \alpha_3 \text{EMD}(P_x, \bar{P}_x) + \alpha_4 \text{MKL}(P_x, \bar{P}_x) + \alpha_5 \xi(P) + \alpha_0$$

ضرایب α از طریق حداقل کردن تابع اتلاف، که چهار معیار فاصله، برای دامنه‌های منبع و خطای رده‌بندی داخلی دامنه‌های منبع را به‌عنوان ورودی می‌گیرد، به دست می‌آید. با جایگذاری این ضرایب در رابطه (۱۴)، برای هر یک از دامنه‌های منبع و دامنه هدف، فاصله هر یک از دامنه‌های منبع و دامنه هدف به دست می‌آید. دامنه منبعی که کمترین فاصله با دامنه هدف $d(P_x, \bar{P}_x)$ را داشته باشد، به‌عنوان دامنه منبع بهینه انتخاب می‌شود.

۴. ارزیابی عملکرد

در این مقاله برای کاهش پیچیدگی دامنه‌ها، اطلاعات اضافی با حذف علامت‌گذاری‌ها و تبدیل همه حروف به حروف کوچک تقلیل شده است و برای نمایش دامنه‌ها از سبک کلمات یکتایی و دوتایی استفاده شده است که کلمه‌ها فضای ویژگی‌ها را تشکیل می‌دهند و تعداد تکرار هر کلمه به‌عنوان مقدار آن ویژگی تلقی می‌شود. هر سند با مجموع بردارهای ویژگی‌های نشان داده می‌شود. در این مقاله از مدل وزن‌دهی TF-IDF^{۲۸} برای وزن‌دهی به ویژگی‌ها

اختلاف توزیع کمتری با دامنه هدف دارد، شباهت دامنه‌ها اندازه‌گیری می‌شود. در این مقاله از چهار معیار فاصله برای محاسبه میزان شباهت توزیع حاشیه‌ای دامنه‌ها در مدل‌های مختلف استفاده می‌شود. چهار معیار فاصله عبارتند از آزمون آماری پیرسون^{۲۴} [۱۷]، فاصله حرکت زمین^{۲۵} [۱۸]، رگرسیون کولبک-لیبلر متوسط^{۲۶} [۱۹]، اختلاف گشتاور مرکزی^{۲۷} [۲۰] که در ادامه به‌صورت اجمالی به آن‌ها پرداخته می‌شود.

• اولین معیار، فاصله آزمون آماری پیرسون (Chi2) است که شباهت بین دو توزیع P_x و \bar{P}_x را به‌صورت رابطه (۹) محاسبه می‌کند که $p(w)$ و $\bar{p}(w)$ احتمال کلمه w به ترتیب در توزیع‌های P_x و \bar{P}_x هستند و N تعداد کلمه‌های با بیشترین تکرار است.

$$\text{chi2}(P_x, \bar{P}_x) := \sum_{i=1}^N \frac{(P(w_i) - \bar{P}(w_i))^2}{P(w_i) + \bar{P}(w_i)} \quad (9)$$

دومین معیار، تابع فاصله حرکت زمین (EMD) است که مقدار جرم احتمالی و مسافتی که جرم احتمالی باید جابه‌جا شود تا دو توزیع دامنه منبع و هدف به هم نزدیک‌تر شوند را محاسبه می‌کند و به‌صورت رابطه (۱۰) تعریف می‌شود که f_{ij} جرم احتمالی و $d_{i,j}$ مسافت بین خوشه‌های i و j از نمایش خوشه‌ای b و \bar{b} است که هر کدام از b و \bar{b} شامل m خوشه هستند.

$$\text{EMD}(P_x, \bar{P}_x) := \min_F \sum_{i=1}^m \sum_{j=1}^m d_{i,j} f_{i,j} \quad (10)$$

$$\text{Subject to: } \begin{aligned} f_{i,j} &\geq 0, \sum_{j=1}^m f_{i,j} = \bar{p}_j \\ \sum_{i=1}^m f_{i,j} &= p_i \end{aligned}$$

سومین معیار نسخه اصلاح شده رگرسیون کولبک-لیبلر (MKLD) است که با نام آنتروپی رابطه‌ای نیز شناخته می‌شود و به‌صورت (۱۱) تعریف می‌شود. m تعداد خوشه‌های موجود در نمایش خوشه‌ای است.

$$\text{MKL}(p_x, \bar{p}_x) = \sum_{i=1}^m E(p_x)_i \log \frac{E(p_x)_i}{E(\bar{p}_x)_i} + E(\bar{p}_x)_i \log \frac{E(\bar{p}_x)_i}{E(p_x)_i} \quad (11)$$

• چهارمین معیار، اختلاف گشتاور مرکزی (CMD) است و به‌صورت رابطه‌های (۱۲) و (۱۳) تعریف می‌شود.

24- statistic of Pearson's χ^2 test

25- Earth Mover's Distance

26- Mean Kullback-Leibler

27- Central Moment Discrepancy

28- Term Frequency-Inverse Document Frequency

$$l(x, y)_{LR} = \sum_{i=1}^{n_d} \log(1 + e^{-(2y_i - 1)x_i \alpha + \alpha_0}) + \mu \frac{1}{2} \|\alpha\|^2 \quad (15)$$

که در مدل CCEM برای مجموعه داده‌های همگن و ناهمگن ضرایب بهینه α به صورت زیر به دست می‌آید:

$$\alpha_{\text{همگن}} = [0.01, 0.002, 0.04, 0.04, 1.61, 0.14]$$

$$\alpha_{\text{ناهمگن}} = [0.021, 0.015, 0.58, 0.007, 1.28, 0.24]$$

پس از یافتن ضرایب α در رابطه (15)، از رابطه (16) برای رده‌بندی دودویی داده‌های دامنه هدف استفاده می‌شود.

$$h(x) = \begin{cases} 1 & \text{if } ax + \alpha_0 > 0 \\ 0 & \text{else} \end{cases} \quad (16)$$

برای پیدا کردن بهترین دامنه منبع با کمترین خطای رده‌بندی بین دامنه‌ای $\xi(P, \bar{P})$ جهت آموزش رده‌بند، از خطای رده‌بندی بین دامنه رابطه‌ای (17) استفاده می‌شود.

$$\xi_{\text{related}}(\hat{p}, \bar{p}) = \xi(\hat{p}, \bar{p}) - \xi(p^*, \bar{p}) \quad (17)$$

در رابطه (17) \hat{p} ، \bar{p} و p^* به ترتیب توزیع دامنه انتخابی منبع توسط روش پیشنهادی، توزیع دامنه هدف و توزیع دامنه منبع بهینه است. می‌دانیم p^* بهترین دامنه منبع با کمترین خطای رده‌بندی بین دامنه‌ای است و اگر مدل پیشنهادی بهترین دامنه منبع را انتخاب کند خطای رابطه‌ای صفر می‌شود.

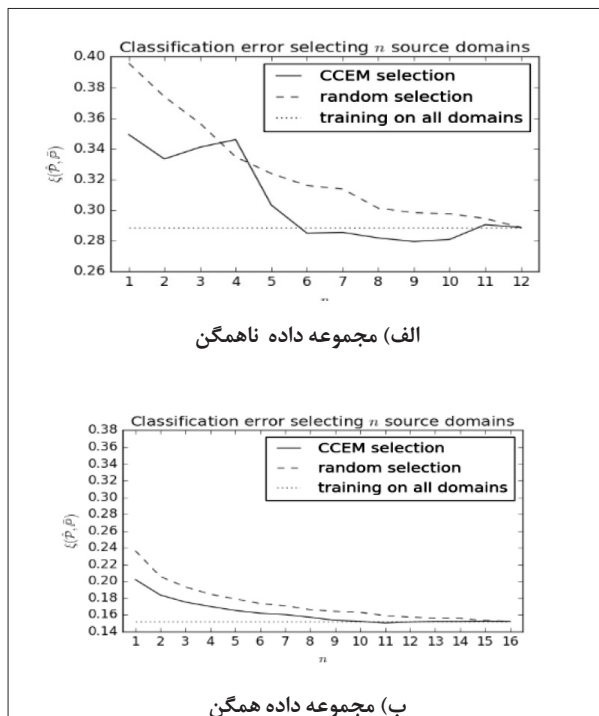
در شکل (1) توزیع خطای رابطه‌ای مدل پیشنهادی با مدل تصادفی برای انتخاب یک دامنه منبع مقایسه شده است. محور افقی خطای رده‌بند بین دامنه‌ای رابطه‌ای است که از رابطه (17) به دست می‌آید که نشان‌دهنده خطای انتخاب دامنه صحیح به عنوان دامنه منبع است و محور عمودی، احتمال انتخاب دامنه با خطای رده‌بندی بین دامنه‌ای رابطه‌ای مشخص به عنوان دامنه منبع است.

همان‌طور که در شکل (1) دیده می‌شود مدل CCEM، در مجموعه داده ناهمگن با احتمال 0.53 و در مجموعه داده همگن با احتمال 1 دامنه منبعی با کمترین خطای رابطه‌ای انتخاب می‌کند که نشان از قدرتمند بودن روش پیشنهادی در انتخاب دامنه منبع بهینه، برای آموزش رده‌بند است.

در شکل (2) خطای رده‌بندی بین دامنه‌ای مطلق را

استفاده شده است. برای ارزیابی روش پیشنهادی از دو مجموعه داده همگن و ناهمگن استفاده شده است این دو مجموعه داده در چارچوب، موضوع، تعداد سندهای موجود در دامنه و طول سندها با همدیگر متفاوت هستند. مجموعه داده همگن، شامل 17 دسته پیکره^{۲۹} از DRANZIERA است [21] که از هر پیکره 5000 سند مثبت و 5000 سند منفی انتخاب شده است و طول سندها بین 64 تا 123 کلمه متغیر است. هر پیکره شامل احساسات کاربران آمازون روی محصولات از دسته‌های مختلف است. در هر مرحله از آزمایش 17 مرحله‌ای مجموعه داده همگن، یک دسته برای دامنه هدف و 16 دسته دیگر به عنوان دامنه منبع در مجموعه نامزد، در نظر گرفته می‌شود تا برای هر دسته از 17 دسته پیکره به عنوان دامنه هدف، یک دامنه منبع بهینه از مجموعه نامزد انتخاب شود. مجموعه داده ناهمگن نیز از 13 دسته پیکره تشکیل شده است که متوسط تعداد سند در پیکره‌ها 6326 است و طول سندها بین 10 تا 15 کلمه متغیر است. در مجموعه داده ناهمگن، احساسات درباره فیلم‌ها [22]، محصولات مصرفی [23]، مجموعه توییت‌ها با موضوع‌های مختلف [24]، پیکره‌ها را تشکیل می‌دهند. در هر مرحله از آزمایش 13 مرحله‌ای مجموعه داده ناهمگن، یک دسته برای دامنه هدف و 12 دسته دیگر به عنوان دامنه منبع در مجموعه نامزد، در نظر گرفته می‌شود تا برای هر دسته از 13 دسته پیکره به عنوان دامنه هدف، یک دامنه منبع بهینه از مجموعه نامزد انتخاب شود. داده‌های آموزشی، از ترکیب چهار معیار فاصله بین دامنه‌های مجموعه نامزد و خطای داخلی دامنه‌های مجموعه نامزد و داده‌های آزمون، از ترکیب چهار معیار فاصله بین دامنه هدف و دامنه‌های مجموعه نامزد و خطای داخلی دامنه‌های مجموعه نامزد تشکیل می‌شود.

برای تابع اتلاف از رگرسیون خطی استفاده شده است که به صورت رابطه (15) تعریف می‌شود که با یافتن حداقل مقدار تابع اتلاف، بهترین ضرایب α برای ترکیب خطی CCEM در رابطه (14) به دست می‌آید.

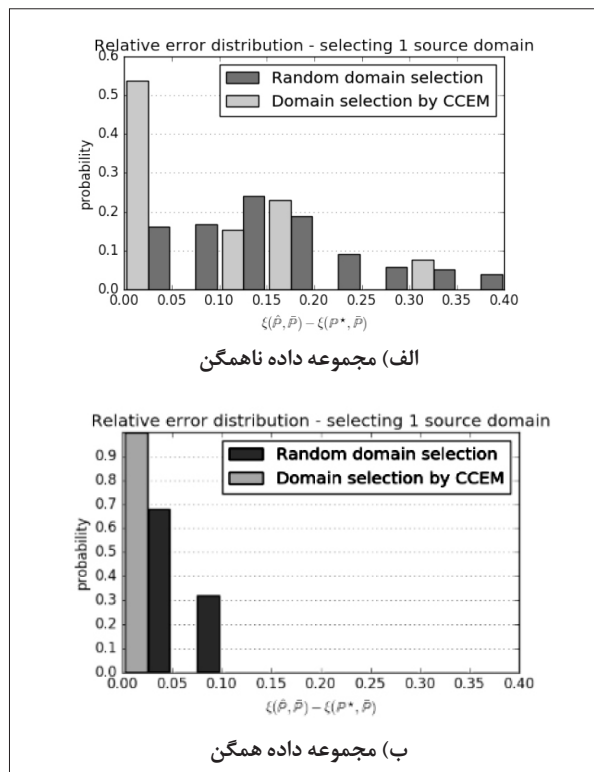


شکل ۲: مقایسه خطای رده‌بندی بین دامنه‌های مطلق برای مدل‌های CCEM و تصادفی با انتخاب n دامنه منبع و مدلی که رو همه مجموعه داده‌ها آموزش دیده، (الف) برای مجموعه داده ناهمگن، (ب) برای مجموعه داده همگن

در جدول (۱) پارامتر احتمال انتخاب دامنه منبع صحیح، نشان‌دهنده این است که دامنه صحیح با چه احتمالی به‌عنوان دامنه منبع انتخاب می‌شود که مدل پیشنهادی در مجموعه داده همگن با احتمال نزدیک به ۵۰ درصد (۱، ۰.۴۷ درصد) و در مجموعه داده ناهمگن با احتمال ۲۳،۱ درصد می‌تواند دامنه منبع صحیح را انتخاب کند که نسبت به مدل تصادفی بهبود چشمگیری دارد.

پارامتر احتمال انتخاب یکی از سه دامنه منبع با بدترین خطای بین دامنه‌ای، در جدول (۱) نشان می‌دهد که روش انتخابی با چه احتمالی یکی از سه دامنه‌ای که بیشترین خطای رده‌بندی بین دامنه‌ای دارند را به‌عنوان دامنه منبع انتخاب می‌کند که مدل CCEM با احتمال صفر درصد در مجموعه داده همگن و با احتمال ۷،۷ درصد در مجموعه داده ناهمگن عملکرد بهتری نسبت به مدل تصادفی دارد و نشان‌دهنده این است که روش ارائه شده، دامنه منبع با خطای رده‌بندی بین دامنه‌ای بالا را انتخاب نمی‌کند.

پارامتر خطای رده‌بندی بین دامنه‌ای مطلق در جدول



شکل ۳: توزیع خطای رابطه‌ای مدل CCEM و مدل تصادفی برای انتخاب یک دامنه منبع، (الف) برای مجموعه داده ناهمگن، (ب) برای مجموعه داده همگن

زمانی که مدل بیش از یک دامنه منبع را انتخاب می‌کند نشان می‌دهد که در این شکل مدل‌های CCEM، تصادفی و مدلی که روی همه مجموعه داده‌ها آموزش دیده، با هم مقایسه شده‌اند. محور افقی، تعداد دامنه منبع انتخابی و محور عمودی، خطای رده‌بندی بین دامنه‌ای است که از رابطه (۳) به دست می‌آید.

در شکل (۲) همان‌طور که دیده می‌شود مدل CCEM نسبت به مدل تصادفی خوب عمل می‌کند و حتی عملکرد مدل پیشنهادی در مجموعه داده ناهمگن زمانی که بین ۶ تا ۱۱ دامنه منبع انتخاب می‌شود نسبت به مدل آموزش دیده روی همه مجموعه داده‌ها بهتر است.

در جدول (۱) احتمال انتخاب بهترین دامنه منبع درست، احتمال انتخاب یکی از ۳ دامنه منبع با بدترین خطای رده‌بندی بین دامنه‌ای و خطای رده‌بندی بین دامنه‌ای مطلق برای سه مدل CCEM، تصادفی و حالت بهینه نشان داده شده است.

جدول ۱: عملکرد سه مدل انتخابی حالت بهینه، CCEM و تصادفی

مدل تصادفی (%)	مدل CCEM (%)	مدل بهینه (%)	مدت انتخابی
۵.۹	۴۷.۱	۱۰۰	احتمال انتخاب دامنه منبع صحیح
۷۰.۶	۰	۰	احتمال انتخاب یکی از سه دامنه منبع با بدترین خطای بین دامنه‌های
۲۴.۷	۲۰.۲	۱۹.۹	خطای رده‌بندی بین دامنه‌های مطلق
۸.۳	۲۳.۱	۱۰۰	احتمال انتخاب دامنه منبع صحیح
۴۱.۷	۷.۷	۰	احتمال انتخاب یکی از سه دامنه منبع با بدترین خطای بین دامنه‌های
۴۰.۳	۳۳.۲	۲۵.۲	خطای رده‌بندی بین دامنه‌های مطلق

و با استفاده از توزیع دامنه منبع و توزیع حاشیه‌ای دامنه هدف، می‌تواند با احتمال خوبی بهترین دامنه منبع را به دست دهد. با تعریف n به‌عنوان تعداد دامنه‌های مجموعه نامزد، مدل بر روی $n(n-1)$ جفت دامنه منبع و هدف مورد آزمایش قرار می‌گیرد تا با یافتن بهترین بردار ضرایب α برای ترکیب خطی رابطه (۱۴)، بتواند با حداقل کردن خطای مطلق تخمین، دامنه منبع بهینه‌ای انتخاب کند. مدل پیشنهادی بر روی داده‌های همگن و ناهمگن با مدل تصادفی و مدل آموزش دیده روی همه دامنه‌ها مقایسه شد.

برای بهبود این روش می‌توان روی روش‌های انتخاب ویژگی تمرکز کرد که ویژگی‌هایی با اطلاعات مفیدتر درباره برچسب داده‌ها انتخاب شوند. در این مقاله رده‌بندی روی یک دامنه آموزش می‌بیند در حالی که اگر از روش‌هایی استفاده شود که رده‌بندی روی چند دامنه یادگیری کند نتایج بهتری به دست می‌آید. برای تعیین تعداد بهینه دامنه‌ها برای یادگیری نیز می‌توان از روش‌های متفاوت استفاده کرد. همچنین اگر بتوان معیارهای اندازه‌گیری فاصله‌ای را پیدا کرد که روی همه مجموعه داده‌ها خوب کار کنند یکی از معایب این روش، که عملکرد ضعیف بعضی معیارها روی مجموعه داده‌های مختلف است، رفع خواهد شد.

مراجع

1. Tahmoresnezhad, J., and Hashemi, S., "Visual domain adaptation via transfer feature learning", Knowledge and Information Systems, vol.50, no. 2, pp. 585-605, 2016.
2. Tahmoresnezhad, J., and Hashemi, S., "Common feature extraction in multi-source domains for transfer learning", information and knowledge technology (IKT), vol. 7, IEEE, 2015.
3. Gretton, A., Borgwardt, K. M., Rasch, M. J., Schölkopf, B., and Smola, A., "A kernel two-sample test", Journal of Machine Learning Research", vol. 13, pp. 723-773, Mar 2012.
4. Weiss, K., Khoshgoftaar, T.M. and Wang, D., "A survey of transfer learning", Journal of Big data, vol.3, no.1, pp.9, Dec 2016.
5. Gong, B., Grauman, K. and Sha, F., "Connecting the dots with landmarks: Discriminatively learning domain-invariant features for unsupervised domain adaptation", In International Conference on Machine Learning, vol.28, no.1 pp. 222-230, February 2013.

(۱) نشان‌دهنده میانگین خطای رده‌بندی بین دامنه‌های همه دامنه‌های هدف است. به عبارت دیگر هر یک از دامنه‌های مجموعه نامزد یک بار به‌عنوان دامنه هدف انتخاب می‌شود و خطای رده‌بندی بین دامنه‌ای برای هر کدام از دامنه‌های هدف محاسبه شده و میانگین آن‌ها را به‌عنوان خطای رده‌بندی بین دامنه‌ای مطلق در نظر می‌گیریم. مدل CCEM با خطای ۳۳.۲ درصد نسبت به مدل تصادفی با خطای ۴۰.۳ درصد عملکرد نسبتاً خوبی دارد.

۵. نتیجه‌گیری

در این مقاله با توجه به بار محاسباتی آموزش رده‌بندی روی همه مجموعه داده‌ها، روش CCEM پیشنهاد شد که از ترکیب خطی چهار معیار فاصله، به همراه خطای رده‌بندی داخلی دامنه و خطای انحراف^۳ تشکیل می‌شود

30- offset error

- empirical comparison”, in 2015 IEEE International Symposium on Multimedia (ISM), pp. 233–236, Dec. 2016.
19. Zhuang, F., Cheng, X., Luo, P., Pan, S. J., and Qing He., “Supervised representation learning: Transfer learning with deep auto encoders”, In International Joint Conference on Artificial Intelligence, Jun. 2015.
 20. Zellinger, W., Grubinger, T., Lughofer, E., Natschläger, T., Saminger-Platz, S., “Central moment discrepancy (cmd) for domain-invariant representation learning”, In International Conference on Learning Representations, 2017.
 21. Dragoni, M., Tettamanzi, A. G., Pereira, C. D. C., Dranziera: an evaluation protocol for multi-domain opinion mining, in Tenth International Conference on Language Resources and Evaluation (LREC 2016), European Language Resources Association (ELRA), pp. 267–272, May. 2016.
 22. University of Michigan, UMich SI650 - Sentiment Classification, 2011.
 23. Kotzias, D., Denil, M., De Freitas, N., and Smyth, P., “From group to individual labels using deep features”, in Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pp. 597–606, Aug. 2015.
 24. ‘CrowdFlower, Data for everyone’, [Online]. Available: <https://www.figure-eight.com/data-for-everyone/>. [Accessed: 2017].
 6. Long, M., Wang, J., Ding, G., Sun, J., and Yu, PS., “Transfer feature learning with joint distribution adaptation”, IEEE international conference on computer vision, pp. 2200–07, Dec. 2013.
 7. Gong, B., Shi, Y., Sha, F., and Grauman, K., “Geodesic flow kernel for unsupervised domain adaptation”, IEEE conference on computer vision and pattern recognition, pp. 2066–73, Jun. 2012.
 8. M. Long, J. Wang, G. Ding, S. J. Pan and P. Yu, “Adaptation regularization: a general framework for transfer learning”, IEEE Trans. Knowl. Data Eng, vol. 26, pp. 1076–1089, 2013.
 9. Maulik, U., and Chakraborty, D., “Remote sensing image classification: a survey of support-vector-machine-based advanced techniques”, IEEE Geoscience and Remote Sensing Magazine, vol. 5, no. 1, pp. 33-52. Mar. 2017.
 10. Iounousse, J., Er-Raki, S., El Motassadeq, A. and Chehouani, H., “Using an unsupervised approach of Probabilistic Neural Network (PNN) for land use classification from multitemporal satellite images”, Applied Soft Computing, vol. 30, pp. 1-13. May. 2015.
 11. Gui, C., Li, W., Pan, Z., Zhang, J., Zhu, J., Cui, D., “A classifier for diagnosis of manic psychosis state based on SVM-GMM”, Sydney: The 10th International Conference on Information Technology and Applications (ICITA2015), Jul. 2015.
 12. Wetzel, D., Lopez, A. and Webber, B., “A maximum entropy classifier for cross-lingual pronoun prediction”, In Proceedings of the Second Workshop on Discourse in Machine Translation, pp. 115-121, 2015.
 13. Li, F., Pan, S. J., Jin, O., Yang, Q., and Zhu, X., “Cross-domain co-extraction of sentiment and topic lexicons”, In: Proceedings of the 50th annual meeting of the association for computational linguistics long papers, vol. 1, pp. 410–19, July. 2012.
 14. Wang, C., and Mahadevan, S., “Heterogeneous domain adaptation using manifold alignment”, In: Proceedings of the 22nd international joint conference on artificial intelligence, vol. 2, pp. 541–46, Jun. 2011.
 15. Nam, J., Fu, W., Kim, S., Menzies, T. and Tan, L., “Heterogeneous defect prediction”, IEEE Transactions on Software Engineering, vol. 44, no.9, pp. 874-896, Sep. 2018.
 16. Blitzer, J., Dredze, M., Pereira, F., and et al., “Biographies, bollywood, boom-boxes and blenders: Domain adaptation for sentiment classification”, in ACL, vol. 7, pp. 440–447, 2007.
 17. K. Pearson, “X. on the criterion that a given system of deviations from the probable in the case of a correlated system of variables is such that it can be reasonably supposed to have arisen from random sampling”, The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science, vol. 50, no. 302, pp. 157–175, Jul. 1900.
 18. Beecks, C., Uysal, M., and Seidl, T., “Earth mover’s distance vs. quadratic form distance: An analytical and