



روش های ایجاد پراکندگی و استفاده از خصیصه های متعدد در خوش بندی ترکیبی

سوسن حسینی

دانشجوی کارشناسی ارشد، نرم افزار کامپیوتر، دانشگاه آزاد اسلامی مشهد

Sosan.hosseini@gmail.com

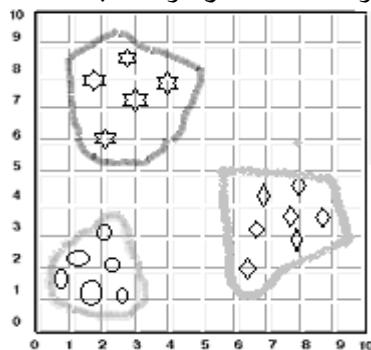
چکیده

خوش بندی یکی از تکنیک های داده کاوی می باشد و به عنوان مهم ترین مسئله در یادگیری بدون نظارت در نظر گرفته می شود. با ترکیب خوش بندی های توان به خوش بندی ترکیبی برای دست یافتن به نتیجه ای، پایدار و قدرتمند است یافت. و با انتخاب یک روش ایجاد پراکندگی مناسب به خوش بندی اصولی دست یافت.

کلید واژه- خوش بندی، داده کاوی، ایجاد پراکندگی، خصیصه های متعدد، خوش بندی ترکیبی.

۱- مقدمه

یکی از تکنیک های مهم داده کاوی، خوش بندی است که داده های مشابه و متفاوت را از هم جدا می کند. خوش بندی به معنای تقسیم بندی بدون نظارت است که در آن داده ها به دسته هایی که از نظر پارامترهای مورد علاقه، شباهت بیشتری به یکدیگر دارند، تقسیم می گردند. ایده اصلی [1] در خوش بندی اطلاعات جدا کردن نمونه ها از یکدیگر و قراردادن آنها در گروه های شبیه به هم می باشد. در شکل ۱ نمونه ای از خوش بندی دیده می شود.



شکل ۱: نمونه خوش بندی [2]

هر یک از الگوریتم های خوش بندی با توجه به اینکه بر روی جنبه های متفاوتی از داده ها تأکید می کنند، داده ها را به صورت های مختلفی طبقه بندی می کنند. به همین دلیل نیازمند شناخت و بررسی روش های ایجاد پراکندگی در خوش بندی ترکیبی هستیم. ادامه این مقاله اینگونه سازماندهی شده است: در بخش دوم مروری بر خوش بندی ترکیبی و مرحله آن و در بخش سوم به بررسی روش های ایجاد پراکندگی و در بخش چهارم استفاده از خصیصه های متعدد در خوش بندی و در بخش آخر به نتیجه گیری پرداخته خواهد شد.

۲- خوش بندی ترکیبی

هر یک از الگوریتم های خوش بندی با توجه به اینکه بر روی جنبه های متفاوتی از داده ها تأکید می کند، داده ها را به صورت های متفاوتی طبقه بندی می نماید. در واقع هدف اصلی [3] خوش بندی ترکیبی جستجوی بهترین خوش بندی است: از ترکیب نتایج الگوریتم های دیگر است. خوش بندی داده به صورت ترکیبی میتواند جواب های ایمن تر و موثر تری را از نظر سادگی، پایداری و انعطاف پذیری ارایه دهد. مطالعات اخیر [4] نشان می دهد که اجماع خوش بندی می تواند خارج از وضعیت های نوع رأی گیری با استفاده از روش های

مبتنی بر گراف و تئوری اطلاعات بدون حل دقیق مشکل تناظربرچسب ها انجام شود، روش خوشه بندی ترکیبی، با ترکیب افزارهای مختلف تولیدشده خوشه بندی پایه یک افزار مستحکم از داده را تولید می کند. در اکثر الگوریتم های پایه برای خوشه بندی ترکیبی از نمونه برداری داده ها استفاده می شود. مسئله اصلی در این روش ها چگونگی ارزیابی خوشه و افزار است. در چند سال اخیر [5] پایداری خوشه به عنوان یک معیار ارزیابی خوشه موردنوجه قرار گرفته است.

۱-۲ مراحل خوشه بندی ترکیبی

دوم رحله مهم در خوشه بندی ترکیبی وجود دارد، که عبارتند از [6]:

۱- ترکیب نتایج بدست آمده از خوشه بندی های متفاوت اولیه برای تولید خوشه بندی نهایی که این کارت و سط تابع توافقی (الگوریتم های ترکیب کننده) انجام می گردد.

۲- تولید نتایج مختلف از خوشه بندی های عنوان نتایج خوشه بندی اولیه براساس اعمال روش های مختلف که این مرحله ایجاد تنوع یا پراکندگی تعریف می شوند، که در ادامه بررسی می گردد.

۲-۲ تابع توافقی

ترکیب نتایج بدست آمده از خوشه بندی های متفاوت اولیه برای تولید خوشه نهایی برای پیدا کردن افزارهای ترکیب شده تعریف می شوند. به طور کلی این روش ها، روش مبتنی بر ماتریس-رأی گیری-تئوری اطلاعات و روش مبتنی بر، ابرگراف و مدل آمیزشی نامیده می شوند. اکثر مطالعات اخیر [7] در حوزه خوشه بندی ترکیبی سعی دارند، ابتدا خوشه بندی های اولیه تا حد ممکن پراکنده تولید شوند و سپس با اعمال یک تابع توافقی همه این نتایج با هم ترکیب شوند.

۲-۳ ایجاد تنوع

تنوع پراکندگی باعث داشتن نتیجه مطلوب ترد خوشه بندی می گردد. که برای، بدست آوردن نتایج متنوع و ویژگی های خاص از داده ها، می توان از دور روش کلی [8] استفاده نمود: ۱- استفاده از الگوریتم های خوشه بندی ۲۰. ۲- استفاده از روش های مختلف ایجاد پراکندگی بر روی تنها یک الگوریتم خوشه بندی پایه و تقسیم بندی داده به قسمت های متفاوت. شامل زیر مجموعه های مختلف از داده و ویژگی که در ادامه به آن خواهیم پرداخت.

۳- بررسی روش های ایجاد تنوع و پراکندگی

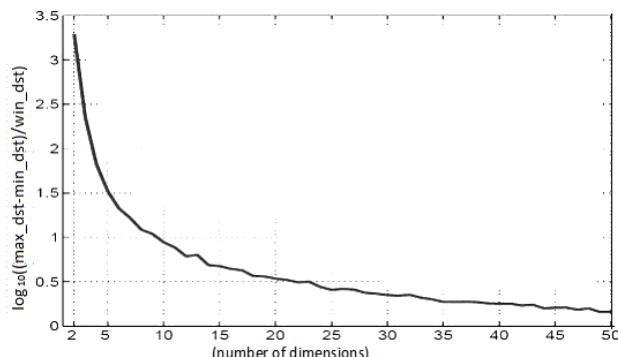
اکثر مطالعاتی [9] که تاکنون در زمینه خوشه بندی ترکیبی انجام گرفته است، جستجوی روش هایی است که بتواند پراکندگی نتایج مجمع رده بندی ها را بهتر کند. در ادامه به بررسی سه روش ایجاد تنوع و پراکندگی پرداخته خواهد شد.

۳-۱ تغییر مقادیر اولیه خوشه بندی:

یکی از راه های الگوریتم افزایش پراکندگی تغییر پارامتر های اولیه الگوریتم های خوشه بندی می باشد. در این روش میتوان تنها از یک الگوریتم خوشه بندی استفاده نمود. انتخاب اولیه مراکز خوشه ها می تواند به دو صورت تصادفی و هوشمندانه انجام گردد.

۳-۲ انتخاب زیر مجموعه مختلف از ویژگی ها:

این روش بر روی زیر مجموعه های مختلف از ویژگی ها اجرامی شود. و نشان داده شده است، که زیر مجموعه بهتر منجربه خوشه پذیری، بهتری می گردد. در شکل 2 می بیزید، وقتی ابعاد زیادی گردد و اسط بین نقاط می گردد و همچنین کاهش ابعاد منجر به تجسم بهتر داده و کاهش زمان و مقدار حافظه می گردد.



شکل ۲: مقایسه رابطه واسط نقاط و ابعاد [10].

۳-۳ انتخاب زیرمجموعه های مختلف ازداده:

یکی، از مهمترین عواملی که در بالابردن کیفیت نهایی خوش بندی تاثیردارد، ایجاد پراکنده بیشترمی باشد، که با انتخاب تعداد محدودی از نمونه ها [11] به جای کل نمونه هاباعث کاهش میزان محاسبات و افزایش پراکنده می گردد.

۴- استفاده از خصیصه های متعدد

در خوش بندی درجات متفاتی از شباهت میان داده ها و مراکز ثقل انتخاب شده وجوددارد. اگر هر الگودقیقاً با یک مرکز منطبق شود، تفکیک قطعی یا سخت امکان پذیر است و داده ها را می توان در خوش بندی مجزا دسته بندی کرد. اما در اغلب موارد مرزهای بین خوش بندی ها خوش تعریف نیستند؛ لذا الگوهار درجات تعلق مختلف به خوش بندی ها نسبت می دهیم. یکی از انواع خوش بندی فازی که کاربرد زیادی دارد خوش بندی عددی fuzzy c-) fcm (mean) که در اصل برای خوش بندی داده با خصیصه های عددی ارایه شده است. در معیار عدم تشابه بین داده ها یک فاصله اقلیدسی یا فاصله ای که روی داده های عددی دارای خصیصه های گروهی هستند تعریف می شود. خوش بندی داده های گروهی موضوع مهمی در داده کاوی و شناسایی الگو است. الگوریتم k-mode توسعه الگوی k-means برای خوش بندی داده های گروهی است که در آن از یک معیار تطابق عدم تشابه ساده برای اشیاء گروهی استفاده می شود. در این روش از mode بجای mean استفاده می شود. خصیصه های ترتیبی را می توان با استفاده از نگاشت های ترتیبی- عددی بصورت الگوریتم های خوش بندی عددی در نظر گرفت. یادگیری تطبیقی نگاشت ترتیبی- عددی، در [12] را می توان بعنوان بهبود کیفیت خوش بندی فازی برای اشیائی با خصیصه های فازی مورد استفاده قرارداد. به بود دیگری در معیار عدم تشابه بین اشیاء و الگو ها بطور مجزا در انواع مختلف الگوریتم KM در [13] برای اشیاء گروهی پیشنهاد شده است. در این الگوریتم، یک الگوی خوش بعنوان خصیصه گروهی خاص، لیستی از همه برچسبها در خصیصه به همراه رخدادهایشان در خوش می باشد. معیار عدم تشابه بین یک شئ والگو در سطح خصیصه، یک منهای تعداد رخدادهای مقدار خصیصه (برچسب) شئ در نظر گرفته می شود که الگوریتم FCM با الگو مستقل در [14] ارائه شده است. اگرچه همگرایی این دو الگوریتم در [15] به اثبات نرسیده است، با این وجود در کاربرد واقعی، می توان از معیار عدم تشابه مختلفی استفاده کرد. در جدول (۱) مجموعه های A,B,C به همراه مشخصات اشان درج شده است.

نام شیء	تعداد هر شیء	خصوصیت عددی یا غیر عددی	غیرگری گروهی	ماکزیمم برچسب خصوصیت گروهی	تعداد کلاس
A	47	0	35	7	4
B	101	1	15	2	7
C	653	6	9	14	2

جدول ۱ : مشخصات سه مجموعه از اشیاء واقعی . [16]

الگوریتم GFCM، ۲۰۰ بار روی داده های مقادیر مختلف p اجرا شده است و مفروضات زیر رنظر گرفته شده است: ۱- در هر اجرا، الگو های مختلفی بکار گرفته شده است. (مقدار $0.0001 =$ برای شرط خاتمه در نظر گرفته شده است). ۲- اگر تعداد ادبر چسب یک خصیصه گروهی خاص، کمتر از p باشد، تعداد دقیق برچسب به جای p استفاده می شود. با توجه به تعداد خوش بکار رفتہ، میانگین مقادیر آندیس برای همه مقادیر p اندازه و در جدول (۱-۱) نشان داده شده است.

C: (WITH M= 1,1)										
MPC	6	6	9 ,10 ,13	13	14	11	12	11 ,12	11 ,13	
CPI	6	6 ,7	9 ,10 ,13	9 ,11	14	11	12	11	11	
XB	11	13	12	11	11	13 ,14	10	11	11	
TANC	9	13	10	11	11	13	10	11	11	
COC	6	6	7 ,11 ,12	13	7	14	12 ,14	11	11 ,13	

جدول ۱-۱: بهترین مقادیر p با توجه به آندیس های صحت خوش فازی با تعداد متفاوت خوش . [16] .

مقادیر کوچک m بویژه ($1,1$) برای اشیائی با خصیصه های گروهی ، به خوش بندی دقیق تری منجر می گردد و بیشینه مقادیر p برای A,B,C به ترتیب 7 ، 2 و 14 می باشد. ضرایب فازی سازی متفاوتی برای m در نظر گرفته شده اند ، که این منجر به حصول مقادیر آندیس متفاوتی شده است.

۵-نتیجه گیری

در خوش بندی ترکیبی و ایجاد پراکندگی در خوش هاب رای داشتن یک خوش بندی ترکیبی مؤثر تر مطالعات نتایج بسیاری می توانند انجام پذیرد. مثلاً مطالعاتی در زمینه بهبود انتخاب هوشمندانه مقادیر اولیه الگوریتم ها و همچنین انتخاب زیرمجموعه ای بهینه و کامل ازداده ها و بویژگی ها برای ایجاد تنوع و پراکندگی مفید می باشد. وبکارگیری الگو های تک خصیصه ای و نیز متدهای تطبیق عدم تشابه ساده در الگوریتم های خوش بندی سنتی برای خصیصه های گروهی ، می تواند به نتایج خوش بندی غیر دقیق منجر شود. در [17] مفهوم جدید الگو های خوش p -mode فازی تعریف شده است که شامل چندین برچسب در سطح خصیصه گروهی می باشد، که همگرا یی این

الگوریتم جدید (GFCM) به اثبات رسیده است. بررسی ها روی مجموعه اشیاء واقعی نشان داده است که مقادیر p در الگوهای p -mode فازی و ضرایب فازی سازی در کارایی الگوریتم FCM موثرند در صورت بکارگیری متدهای تطبیق عدم تشابه ساده برای مقادیر گروهی، ضرایب فازی سازی کوچک، مثل (۱,۱) روی اشیائی با خصیصه های گروهی منجر به خوش بندی بهتری می شوند. همچنین اگر $p=\max$ اختیار شود، نتایج بهتری نسبت به ($p=1$) حاصل خواهد شد، البته این امر در همه موارد صادق نیست.

مراجع

- [1] F.kovas,c.legany,A.Babos,"cluster validity measurement techniques",department of automation and applied informatics,Budapest university of technology and economics,pp.2-3,2003.
- [2] Q.He,"A Review of clustering Algorithms as Applied in IR",Gr-aduate school of library and information science university of ilinois at urbana-champaign, ,pp97.-98,1999.
- [3] Dudoit and fridly "bagging to improve the accuracy of a clustering procedure ",bioinformatics,pp.1090-1091,2003.
- [4] Topchy,a.k and punch w.f."a mixture model for clustering ensemble",in proc.siam int1.conf.on data minig,icml,pp.389-390,2009.
- [5] Fischer,b and buhmann,j.m"path-based clustering for grouping of smooth curves and texture segmentation",ieee trans.on pami,vol.25,No.4,pp.513-514,2003.
- [6] Feyrn,x and brodle,c"random projection for high dimensional data clustering:a cluster ensemble approach ",in proc.20th int.conf.on machine learnimg,icml,pp.14-16, 2005.
- [7] M. Halkidi, Y. Batistakis and M. Vazirgiannis, "On Clustering Validation Techniques", Journal of Intelligent Systems, vol.23,No.7, pp 107-145, 2001.
- [8] Minaei-bidgoli,B.,jain,A.K and punch,W.F" a comparision of resampling method for clustering ensemble",in proc.int1.conf.machine learning methods technology and application ,Mlmtna 04,las vegas, vol.37,No.6,pp 150-154,2006.
- [9] Topchy,A.Minaei-Bidgoli,B.Jain,A.K.and punch,W.F"adaptive clustering ensemble",in proc.int1.conf on pattern recogation,ICPR,Cambridge,UK,pp.147-149,2004.
- [10] Levine.e and domany.E,"Resampling metgod for unsupervised estimation of cluster validaty".neural Computation, vol.17,No.13,pp.2583-2593,2001.
- [11] Javad azimi,Reza Davoodi,Morteza analoui , "fast Convergence Clustering ensemble,In Conference of Data mining and Data warehouse",pp.14-16,2006.
- [12] W.BOCK.Day,"consensus methods as tools for data analysis.in,editor,classification and related method for data analysis",Elsiver scince publishers B.V ,pp.317-320,1988.
- [13] Fayyad u,Piatetsky-shapiro , "Advances in know ledge discovery and data mining",pp.456-459,1996.
- [14] Brunk.c,Kelly.j,"an integrated system for data mining ",proc.3 ,int.1.conf.on knowledge discovery and data minig,pp.135-138,1997.
- [15] Becker.ra,chambers.jm,"an in teractive environment for data analysis and graphics ",pp.143-145,1984.

[16] M.Lee,w.pedrycz,”the fuzzy C-means algorithm with fuzzy p-mode prototypes for clustering object s having mixed features”,pp.15-17,2009.

[17] Stefano Benati,”categorical data fuzzy clustering:an analysis of local search heuristics,Computers &Operations research”,pp.766-775,2008.