

## جداسازی داده های بدون برچسب به کمک برنامه ریزی خطی

حسین موسائی<sup>۱</sup> و نازیلا صباغی<sup>۲</sup>

<sup>۱</sup> استادیار، دانشگاه بجنورد، گروه ریاضی، بجنورد، ایران،

<sup>۲</sup> دانشجوی کارشناسی ارشد، دانشگاه بجنورد، گروه ریاضی، بجنورد، ایران

### چکیده:

در این مقاله سعی بر این است که داده هایی که بدون برچسب می باشند را به کمک نا مساوی قدرمطلق خوشه بندی کنیم. با استفاده از این روش داده های بدون برچسب به دو دسته تقسیم می شوند. به طوری که بیشتر داده ها در کلاسهای درست و جایگاه مناسب خود قرار گیرند. برای این منظور مدل بدست آمده را به مساله برنامه ریزی خطی تبدیل می کنیم و سپس مساله برنامه ریزی خطی را حل می کنیم. آزمایشات عددی نشان دهنده کارایی روش پیشنهادی می باشند.

### کلمات کلیدی:

خوشه بندی ; داده های بدون برچسب ; بردار پشتیبان ماشین ; نامساوی قدرمطلق.

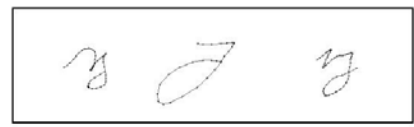
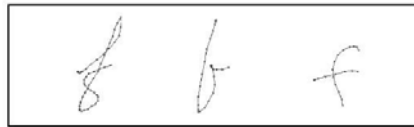
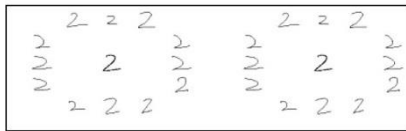
### ۱- مقدمه:

داده ها و الگوها یکی از شاخص های بسیار مهم در دنیای اطلاعات هستند. در مجموعه ای از داده هایی که طبقه بندی نشده اند، می توان ساختار مناسبی یافت که به خوشه بندی موسوم است. خوشه بندی یکی از بهترین روش هایی است که برای کار با داده ها حتی در حالتی که هیچگونه اطلاعاتی از ساختار داخلی داده ها نداریم و داده ها کاملاً بدون برچسب هستند، ارائه شده است. تجزیه و تحلیل خوشه ای روشی برای گروه بندی داده ها با توجه به شباهت یا درجه نزدیکی آنها به یکدیگر است. از این طریق می توان داده ها را به دسته های همگن و متمایز از هم تقسیم کرد. در روش خوشه بندی هیچ کلاسی از قبل وجود ندارد. بلکه ما به دنبال کلاس هایی از داده ها هستیم که به هم شباهت دارند و با کمک این شباهت ها می توان رفتارها را بهتر شناسایی کرد و بر مبنای آن طوری عمل کرد که نتیجه بهتری حاصل شود [۶].

به بیان دیگر می توان گفت که خوشه بندی قراردادن داده ها در گروه هایی است که اعضای هر گروه از زاویه خاصی شباهت دارند. در نتیجه اعضای یک خوشه به یکدیگر شباهت دارند و با اعضای خوشه های دیگر هیچ شباهتی ندارند.

برای مثال شکل های ۱ و ۲ به خوشه بندی دست خط های متفاوت در افراد اشاره می کند [۲].

<sup>1</sup>Corresponding author: E-mail address: hmoosaei@gmail.com



شکل ۱: دسته بندی مدل‌های مختلف

شکل ۲: سه کلاس مختلف برای نوشتن

شکل ۳: سه کلاس مختلف برای نوشتن

نوشتن عدد 2

حرف لاتین f

حرف لاتین y

معیار شباهت و هم‌گروهی داده‌ها در این مقاله فاصله می‌باشد. یعنی داده‌هایی که به یکدیگر نزدیک‌ترند در یک خوشه قرار می‌گیرند.

یکی از روش‌هایی که در حال حاضر به صورت گسترده برای جداسازی خوشه‌ها با محاسبه حداکثر فاصله بین داده‌های کلاس‌های مختلف مورد استفاده قرار می‌گیرد، روش ماشین بردار پشتیبان (SVM) است. ایده‌ی ماشین بردار پشتیبان می‌کوشد، ابرصفحاتی را پیدا کند که عمل تمایز نمونه‌های کلاس‌های مختلف داده‌ها را بطور بهینه انجام دهد. بدین گونه که ابرصفحه‌هایی با حداکثر حاشیه را بدست می‌آورد که دسته‌ها را از یکدیگر جدا می‌کند.

در این مقاله می‌خواهیم با استفاده از نامساوی قدر مطلق در ماشین بردار پشتیبان، داده‌های بدون برچسب را به دو خوشه تقسیم کنیم [۳].

نامساوی قدر مطلق محدب زیر را در نظر بگیرید:

$$|x^t w - \gamma| \leq 1, \quad (1)$$

که در آن  $\| \cdot \|$  بیان‌کننده قدر مطلق است. بردار ستونی  $x$  معرف داده در فضای  $n$ -بعدی،  $R^n$  می‌باشد.  $\gamma$  فاصله مبدا از صفحه را مشخص می‌نماید و "i" ترانهاده است.

نامساوی (۱) را به دو نیم فضای زیر می‌توان تقسیم کرد:

$$\begin{aligned} x^t w &\leq \gamma + 1, \\ x^t w &\geq \gamma - 1. \end{aligned} \quad (2)$$

هدف این است که بتوانیم به جای استفاده از دو نامساوی (۲) به کمک یک نامساوی داده‌ها را از یکدیگر جدا کنیم. بنابر این به

جای نامساوی اول از  $x^t w - \gamma \geq 0$  و به جای نامساوی دوم از  $x^t w - \gamma \leq 0$  استفاده می‌کنیم. در واقع هدف ما پیدا کردن ابر صفحه  $x^t w = \gamma$  می‌باشد [۴].

محتویات مقاله شامل موارد زیر است:

در بخش ۲ فرایند خوشه بندی داده‌های بدون برچسب و تشکیل تابع هدف را بررسی می‌کنیم. بخش ۳ شامل روش پیشنهادی

برای جداسازی داده‌های بدون برچسب می‌باشد. در بخش ۴ نتایج عددی را مشاهده می‌کنیم این نتایج مبین کارایی رویکرد

پیشنهادی است. در بخش ۵ نتیجه گیری را بیان می‌کنیم.

شیوه نمادگذاری مقاله حاضر بدین شرح است :

ضرب اسکالر دو بردار ستونی  $x, y$  در فضای  $n$ -بعدی به صورت  $x^t y$  می باشد. بردار ستونی یک با بعد دلخواه را با  $e$  و بردار ستونی صفر را با  $0$  نمایش می دهیم. s.t. مخفف "به طوری که" می باشد.

## ۲- خوشه بندی داده های بدون برچسب

فرض کنید داده های بدون برچسب شامل  $m$  نقطه در فضای  $n$ - بعدی  $R^n$  است که با ماتریس  $A, m \times n$  نشان می دهیم. هر عضو ماتریس  $A$  به یکی از کلاس های  $A^+$  یا  $A^-$  تعلق می گیرد، که در آن  $A^+$  در نیم فضای  $x^t w \leq \gamma + 1$  و  $A^-$  در نیم فضای  $x^t w \geq \gamma - 1$  واقع می باشد. لذا داریم:

$$|Aw - ey| \leq e. \quad (3)$$

در حقیقت ما به دنبال یافتن حداکثر فاصله دو ابر صفحه  $x^t w - \gamma = \pm 1$  هستیم.

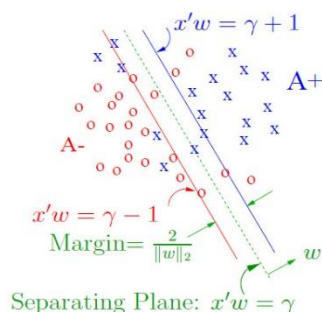
این حداکثر فاصله بین دو صفحه فوق که (حاشیه) نام دارد، یک مزیت برای دستگاه بردار پشتیبان ماشین محسوب می گردد.

$$\text{حاشیه} = \frac{2}{\|w\|_2} \quad \text{همچنین یک رابطه بین } w \text{ و حاشیه وجود دارد یعنی}$$

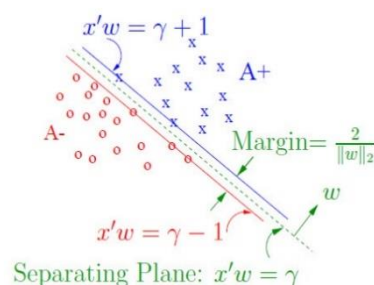
لذا به منظور افزایش حاشیه در تابع هدف باید  $w$  را کاهش دهیم.

از سویی وقتی دو کلاس به شدت خطی و جدا پذیر باشند، خطای مساله صفر است (شکل ۴). وقتی دو کلاس خطی و جداناپذیر باشند، مساله دارای خطا می باشد [۵] (شکل ۵).

بنابراین سعی ما بر آن است که در ضمن کاهش این خطا، حداکثر حاشیه را ایجاد کنیم. همچنین ابر صفحه  $x^t w = \gamma$  که وسط دو صفحه فوق می باشد را بیابیم.



شکل ۵: دو کلاس خطی و جدا ناپذیر



شکل ۴: دو کلاس خطی و جدا پذیر

به منظور خطی کردن رابطه (۳) می توان این رابطه قدر مطلق را با کران بالایش محدود کنیم. پس اگر  $r$  را کران بالای مورد نظر فرض کنیم، داریم:

$$-r \leq Aw - er \leq r. \quad (4)$$

با توضیحات فوق مساله را به صورت زیر مدل می کنیم:

$$\begin{aligned} \min \quad & -e^t |w| - |\gamma| + ve^t r \\ \text{s.t.} \quad & -r \leq Aw - e\gamma \leq r, \\ & r \leq e. \end{aligned} \quad (5)$$

که با یک خطی سازی پی در پی متناهی قابل حل است [۷]. در بخش بعد روشی جدید برای حل مساله (۵) ارائه می دهیم.

### ۳- یادگیری بدون نظارت به کمک برنامه ریزی خطی

در این قسمت روشی جدید برای حل مساله (۵) ارائه می دهیم. در روش پیشنهادی مدل را به صورت یک مساله برنامه ریزی خطی، تبدیل می کنیم و سپس به کمک نرم افزار MATLAB حل می نماییم.

برای این منظور ابتدا توجه داریم که:

$$\begin{aligned} |w| \leq k & \rightarrow -k \leq w \leq k, \\ |\gamma| \leq z & \rightarrow -z \leq \gamma \leq z. \end{aligned} \quad (6)$$

پس مساله جدید به صورت زیر مدل می شود:

$$\begin{aligned} \min \quad & -e^t w + \gamma - e^t k - z + ve^t r \\ \text{s.t.} \quad & -r \leq Aw - e\gamma \leq r, \\ & r \leq e, \\ & -k \leq w \leq k, \\ & -z \leq \gamma \leq z. \end{aligned} \quad (7)$$

لذا مساله زیر را به کمک MATLAB حل می نماییم و ابر صفحه  $x^t w = \gamma$  را تشکیل می دهیم.

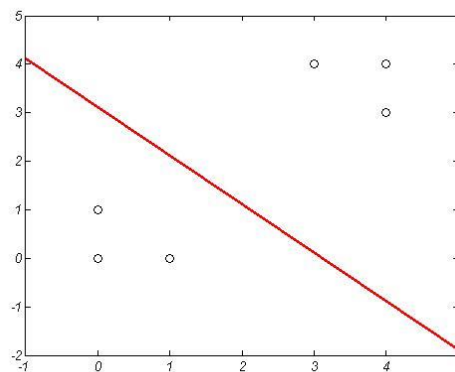
$$\begin{aligned}
 \min \quad & w + \gamma - e^t k - z + v e^t r \\
 \text{s.t.} \quad & Aw - e\gamma \leq r, \\
 & -Aw + e\gamma \leq r, \\
 & r \leq e, \\
 & w \leq k, \\
 & -w \leq k, \\
 & \gamma \leq z, \\
 & -\gamma \leq z.
 \end{aligned} \tag{۸}$$

در قسمت بعد به بیان نتایج عددی خواهیم پرداخت.

#### ۴- نتایج عددی

در این قسمت روش پیشنهادی را بر روی داده های بدون برچسب اجرا می کنیم. برای اجرای الگوریتم از نرم افزار مطلب ۲۰۱۳ بروی رایانه intel core i5-4210U و حافظه 6GB استفاده شده است.

ابتدا مثالی ۲-بعدی با شش داده را بررسی می کنیم. با اجرای روش پیشنهادی نقاط  $[0, 0; 0, 1; 1, 0; 4, 4; 3, 4; 4, 3]$  به دو دسته تقسیم می شود. همانطور که در شکل ۶ مشاهده می شود این جداسازی به طور مناسب انجام شده است.



شکل ۶: مثال دو بعدی و تقسیم به دو کلاس با روش پیشنهادی

مثال دوم بر روی داده های UCI، داده مرکز تشخیصی سرطان سینه دانشگاه ویسکانسین آمریکا WDBC و داده Cancer صورت گرفته است [۸ و ۱]. فرض می کنیم داده ها بدون برچسب هستند و آنها را به کمک روش پیشنهادی جدا می کنیم و سپس دقت را به دست می آوریم. جدول ۱ نتایج را نشان می دهد.

جدول ۱- نتایج روش پیشنهادی برای داده های WDBC , Cancer

نام داده	تعداد داده ها	تعداد داده های خوشه ها	دقت	زمان (ثانیه)
WDBC	۵۶۹	$A^+ = ۲۵۶$	۵۸,۱۳	۳,۹۰
		$A^- = ۳۱۳$		
Cancer	۳۰۳	$A^+ = ۲۱۶$	۶۰,۰۶	۱,۸۴
		$A^- = ۸۷$		

### ۵- نتیجه گیری

در این مقاله روشی به کمک نامساویهای قدر مطلق برای جداسازی داده های بدون برچسب پیشنهاد کردیم. برای حل مدل ایجاد شده ابتدا آن را به صورت برنامه ریزی خطی تبدیل کرده و سپس با استفاده از نرم افزار مطلب حل می‌نمائیم. نتایج عددی نشان دهنده کارایی این روش برای جداسازی داده های بدون برچسب می باشد.

### مراجع

- [1] D.W. Aha and P. M. Murphy, UCI machine learning repository, <http://archive.ics.uci.edu/ml>, 1992.
- [2] K. J. Anil, "Data clustering: 50 years beyond K-means" *Pattern recognition letters* Vol. 31, No.8, pp 651-666, 2010.
- [3] G. M. Fung and O.L. Mangasarian, "Data selection for support vector machine classifiers" *Proceedings of the sixth ACM SIGKDD international conference on Knowledge discovery and data mining*. USA, pp. 64-70, 2000.
- [4] G. M. Fung and O.L. Mangasarian, "Semi-supervised support vector machines for unlabeled data classification" *Optimization methods and software*, Vol. 15, No. 1, pp. 29-44, 2001.
- [5] G. M. Fung and O.L. Mangasarian, "Unsupervised and Semisupervised Classification via Absolute Value Inequalities", *Journal of Optimization Theory and Applications*, Vol. 168. No. 2, pp. 551-558, 2016.
- [6] C. Ma and G. Gan and J. Wu, *Data clustering: theory, algorithms, and applications*. SIAM, Vol. 20, 2007.
- [7] O. L. Mangasarian, "Unsupervised classification via convex absolute value inequalities" *Optimization*, Vol. 64, No. 1, pp. 81-86, 2015.
- [8] W. H. Wolberg and W. N. Street and O. L. mangasarian, WDBC: Wisconsin Diagnostic Breast Cancer Database, Computer Sciences Department University of Wisconsin Madison, <ftp://ftp.cs.wisc.edu/math-prog/cop-dataset/machine-learn/cancer/WDBC>, 1995.