

بررسی تاثیر بکارگیری توابع زیان مختلف بر عملکرد مدل خوشه‌بندی فازی برای داده‌های فازی در صورت وجود داده‌های دورافتاده

الهام اسکندری* و علیرضا خواستان

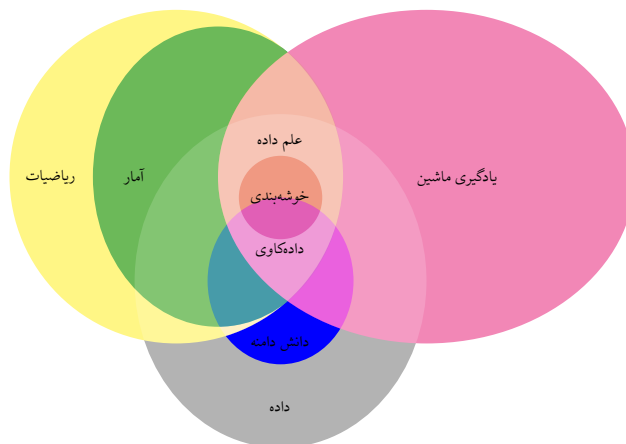
دانشکده ریاضی، دانشگاه تحصیلات تکمیلی علوم پایه، زنجان، ایران

تاریخ پذیرش: ۱۴۰۳/۰۵/۲۸

تاریخ دریافت: ۱۴۰۲/۰۷/۲۰

نوع مقاله: علمی-پژوهشی

چکیده. در این مقاله، استفاده از یک تابع زیان جدید در مدل خوشه‌بندی فازی برای داده‌های فازی را پیشنهاد و تاثیر استفاده از این تابع زیان و توابع زیان مربعی، هابر، خطی، سیگموئیدی و لگاریتمی بر عملکرد مدل را در صورت وجود داده‌های دورافتاده در مجموعه داده‌ها، در شبیه‌سازی صورت پذیرفته بررسی کرده‌ایم. مجموعه داده‌های مورد استفاده، از نظر تعداد ویژگی‌ها (۲ و ۳)، تعداد کلاس‌ها (۳ و ۴) و پخش و تعداد داده‌های دورافتاده (۲۰ و ۳۸ و ۴)، دارای تنوع مناسبی هستند. نتایج شبیه‌سازی موید آن است که تابع زیان هابر و تابع زیان جدید، نسبت به وجود داده‌های دورافتاده مقاوم هستند.



شکل ۱. نمودار ون نشان‌دهنده‌ی ارتباط میان علم داده، یادگیری ماشین، آمار و داده‌کاوی. چنانچه آمار در فرایند خوشه‌بندی به‌کار برده شود، خوشه‌بندی، زیرمجموعه‌ای از علم داده و داده‌کاوی خواهد بود [۱۴].

۱. سرآغاز

خوشه‌بندی یکی از الگوریتم‌های مورد استفاده در یادگیری بدون نظارت^۱ است. یادگیری بدون نظارت، نوعی یادگیری ماشینی^۲ است و یادگیری ماشین، ارتباط در هم تنیده‌ای با داده‌کاوی دارد. از طرفی، یادگیری ماشین (به‌جز روش شبکه عصبی مصنوعی^۳)، همان مباحث مربوط به آمار و یادگیری آماری^۴ است. آمار به‌عنوان شاخه‌ای از ریاضیات، در بسیاری از تحلیل‌های مربوط به علم داده^۵ به‌کار می‌رود. بنابراین، میان این حوزه‌ها تا حدودی هم‌پوشانی وجود دارد (شکل ۱) که اگر از آن‌ها به‌درستی در کنار یک‌دیگر استفاده کنیم، نتایج رضایت‌بخشی را استخراج می‌کنیم.

روش خوشه‌بندی فازی سی- میانگین^۶ که در سال ۱۹۷۴ به‌طور مستقل توسط دان^۷ و بزدک^۸ معرفی و در سال ۱۹۸۱ توسط بزدک [۱] تعمیم داده شد، شناخته شده‌ترین و

¹Unsupervised learning

²Machine learning

³Artificial neural network

⁴Statistical learning

⁵Data science

⁶Fuzzy C-Means (FCM)

⁷Dunn

⁸Bezdek

کاربردی‌ترین روش خوشه‌بندی فازی است. در این روش تابع زیان

$$J(\mathbf{U}, \mathbf{H}) = \sum_{i=1}^N \sum_{g=1}^c (u_{ig})^m d^{\gamma}(\mathbf{x}_i, \mathbf{h}_g),$$

باید با توجه به قیود

$$\sum_{g=1}^c u_{ig} = 1, \quad i = 1, \dots, N,$$

و

$$u_{ig} \geq 0, \quad i = 1, \dots, N, \quad g = 1, \dots, c,$$

به حداقل برسد که در آن نمونه‌ی i ام، \mathbf{h}_g مرکز خوشه‌ی g ام، u_{ig} درجه عضویت نمونه‌ی i ام در خوشه‌ی g ام، $m \in [1, \infty)$ پارامتر فازی کننده و

$$d^{\gamma}(\mathbf{x}_i, \mathbf{x}_{i'}) = \|\mathbf{x}_i - \mathbf{x}_{i'}\|^{\gamma},$$

مربع فاصله‌ی اقلیدسی است.

برخی پژوهشگران توجه خود را به توسعه‌ی روش‌های مقاوم^۱ در مقابل داده‌های دورافتاده^۲ معطوف کرده‌اند؛ بعضی از مهم‌ترین رویکردها عبارت‌اند از انتخاب واسط^۳ داده‌ها به عنوان سرخوشه^۴ به جای انتخاب میانگین آن‌ها (شکل ۲ را ببینید)، استفاده از یک معیار عدم تشابه مقاوم، در نظر گرفتن یک خوشه‌ی اضافی به نام خوشه‌ی نویز، حذف کسری از داده‌های دورافتاده، وزن‌دهی کم به تاثیر داده‌های دورافتاده و احتمالاتی^۵.

رویکردی دیگر، استفاده از توابع زیان مقاوم است که با مقاوم بودن معیار فاصله مرتبط است. دو تابع زیان بسیار معمول، توابع زیان مربعی^۶ یعنی

$$\mathcal{L}_{\text{SQR}}(e) = e^2,$$

و خطی^۷ یعنی

$$\mathcal{L}_{\text{LIN}}(e) = |e|,$$

¹Robust

²Outlier data

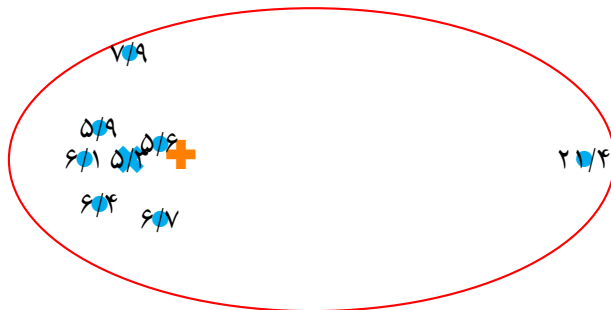
³Medoid

⁴Prototype

⁵Possibilistic

⁶Squared

⁷Linear



شکل ۲. تعداد ۸ نمونه که میانگین آن‌ها با + و واسط آن‌ها با × نشان داده شده است. عدد نوشته شده روی هر نمونه، مجموع فواصل آن نمونه از سایر نمونه‌ها است. واسط، یکی از اعضای مجموعه داده است که کمترین فاصله با سایر نمونه‌ها را دارد. میانگین، به شدت تحت تاثیر داده‌های دورافتاده قرار می‌گیرد و بنابراین نمی‌تواند نماینده خوبی برای اعضای مجموعه داده باشد [۱۱].

هستند. اگر مقدار خطا کمتر از ۱ باشد، زیان مربعی، کوچک و اگر بیشتر از ۱ باشد، زیان مربعی بزرگ می‌شود. در صورت وجود داده‌های دورافتاده که موجب افزایش خطا می‌شوند؛ تابع زیان مربعی، زیان زیادی را محاسبه می‌کند. اشکال اصلی تابع زیان خطی، بزرگ بودن مقدار گرادیان آن برای مقادیرهای کوچک خطاست.

یک روش متداول برای مدل‌سازی عدم قطعیت، استفاده از اعداد فازی است. اعداد فازی، مفهوم اعداد حقیقی کلاسیک را تعمیم می‌دهند. ساتو^۱ و ساتو^۲ [۱۲] نخستین کسانی بودند که در سال ۱۹۹۵ یک مدل کاملاً فازی یعنی یک مدل خوشه‌بندی فازی برای داده‌های فازی ارائه کردند.

هدف ما از مطالعه‌ی این موضوع، آزمودن توابع زیان مختلف در مدل خوشه‌بندی فازی برای داده‌های فازی و یافتن توابع زیان مقاوم نسبت به وجود داده‌های دورافتاده است. در این مقاله، استفاده از تابع زیان مقاوم لگاریتم کسینوس هذلولوی^۳ پیشنهاد شده است. این تابع زیان خاص به دلیل دارا بودن ویژگی‌های جالبی انتخاب شده است که در بخش ۴ به آن خواهیم پرداخت.

¹Sato

²Sato

³Logarithm of hyperbolic cosine

این مقاله شامل شش بخش است. در بخش ۲ توضیحاتی مقدماتی ایراد می‌شود. در بخش ۳ به تاریخچه‌ی موضوع مورد مطالعه اشاره می‌شود. ویژگی‌های تابع زیان پیشنهادی، در بخش ۴ تبیین می‌شود. در بخش ۵ شبیه‌سازی صورت گرفته معرفی و نتایج ارائه می‌شود. بخش ۶ نیز به تفسیر نتایج دست‌یافته اختصاص دارد.

۲. پیش‌نیازها

مجموعه داده‌ی فازی را می‌توان در یک ماتریس $N \times J$ ذخیره کرد، به طوری که N تعداد نمونه و J تعداد ویژگی است.

تعریف ۱.۰۲. فرض کنید $L, R : [0, 1] \rightarrow [0, 1]$ دو تابع پیوسته و نزولی هستند که

$$L(0) = R(0) = 1,$$

و

$$L(1) = R(1) = 0.$$

مجموعه‌ی فازی N ، $x_i : \mathbb{R} \rightarrow [0, 1]$ ، $i = 1, 2, \dots, N$ که

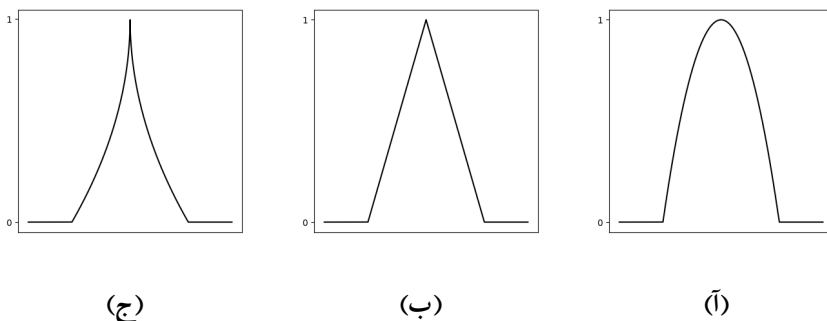
$$x_i(u) = \begin{cases} 0, & u \leq c_{1i} - l_i, \\ L\left(\frac{c_{1i} - u}{l_i}\right), & c_{1i} - l_i \leq u \leq c_{1i}, \\ 1, & c_{1i} \leq u \leq c_{2i}, \\ R\left(\frac{u - c_{2i}}{r_i}\right), & c_{2i} \leq u \leq c_{2i} + r_i, \\ 0, & u \geq c_{2i} + r_i, \end{cases}$$

یک عدد فازی $L - R$ است که در این مقاله، آن را با $x_i = (c_{1i}, c_{2i}, l_i, r_i)$ نشان می‌دهیم به طوری که c_{1i} و c_{2i} به ترتیب، مرکزهای چپ و راست x_i و l_i و r_i به ترتیب، گستره‌های چپ و راست x_{ij} نامیده می‌شوند. وقتی

$$L(x) = R(x) = 1 - x,$$

¹Center

²Spread



شکل ۳. تابع عضویت از نوع (آ) سهموی، (ب) مثلثی و (ج) ریشه‌ی دوم.

تابع عضویت^۱ از نوع دوزنقه‌ای^۲ است. چنان‌چه توابع L و R خطی باشند و داشته باشیم $c_{1i} = c_{2i} = c_i$ ، آن‌گاه x_i را عدد فازی مثلثی^۳ می‌نامیم.

۳. مروری بر کارهای پیشین

در سال ۱۹۹۶، یانگ^۴ و کو^۵ [۱۵] فاصله‌ی

$$d^2(x_i, x_{i'}) = (c_i - c_{i'})^2 + ((c_i - \lambda l_i) - (c_{i'} - \lambda l_{i'}))^2 + ((c_i + \rho r_i) - (c_{i'} + \rho r_{i'}))^2, \quad (۱.۳)$$

را بین دو عدد فازی $L - R$ تک‌متغیره پیشنهاد کردند که در آن $\lambda = \int_0^1 L^{-1}(\omega) d\omega$ و $\rho = \int_0^1 R^{-1}(\omega) d\omega$ پارامترهایی هستند که شکل دُم^۶ چپ و راست تابع عضویت را بیان می‌کنند؛ برای مثال وقتی $\lambda = \rho = \frac{2}{3}$ ، تابع عضویت از نوع سهموی^۷ (شکل ۳ج)، وقتی $\lambda = \rho = \frac{1}{3}$ ، از نوع مثلثی^۸ (شکل ۳ب) و وقتی $\lambda = \rho = \frac{1}{4}$ ، از نوع ریشه‌ی دوم^۸ (شکل ۳ج) است [۴]. مدل پیشنهادی یانگ و کو به شدت تحت تاثیر داده‌های دورافتاده است.

¹Membership function

²Trapezoidal

³Triangular

⁴Yang

⁵Ko

⁶Tail

⁷Parabolic

⁸Square root

طی سه دهه‌ی اخیر، مدل‌های خوشه‌بندی فازی مقاومی از جمله مدل‌های یانگ و لیو^۱ [۱۶]، بوتکیویچ^۲ [۲]، هونگ^۳ و یانگ [۱۰]، زرندی^۴ و رضایی^۵ [۱۷] و دورسو^۶ و جیوانی^۷ [۵] برای داده‌های فازی ارائه شدند.

در سال ۲۰۲۰، دورسو و لسکی^۸ [۶] با در نظر گرفتن نسخه‌ی چندمتغیره فاصله‌ی پیشنهاد شده توسط یانگ و کو (۱.۳)، فاصله‌ی

$$\begin{aligned} \mathcal{D}(\mathbf{x}_i, \mathbf{h}_g) = & \mathcal{L}(\mathbf{c}_{1i} - \mathbf{c}_{1g}) + \mathcal{L}(\mathbf{c}_{2i} - \mathbf{c}_{2g}) + \\ & \mathcal{L}[(\mathbf{c}_{1i} - \lambda \mathbf{l}_i) - (\mathbf{c}_{1g} - \lambda \mathbf{l}_g)] + \\ & \mathcal{L}[(\mathbf{c}_{2i} + \rho \mathbf{r}_i) - (\mathbf{c}_{2g} + \rho \mathbf{r}_g)], \end{aligned}$$

را ارائه کردند که در آن \mathcal{L} یکی از توابع زیان مربعی، خطی، هابر^۹ یعنی

$$\mathcal{L}_{\text{HUB}}(e) = \begin{cases} \frac{e^2}{\delta^2}, & |e| \leq \delta, \\ \frac{|e|}{\delta}, & |e| > \delta, \end{cases}$$

که $\delta > 0$ ، سیگموئیدی^{۱۰} یعنی

$$\mathcal{L}_{\text{SIG}}(e) = \frac{1}{1 + \exp(-\alpha(|e| - \beta))},$$

که $\alpha, \beta > 0$ و لگاریتمی^{۱۱} یعنی

$$\mathcal{L}_{\text{LOG}}(e) = \log(1 + e^2),$$

است. آن‌ها مدل

$$(۲.۳) \quad J(\mathbf{U}, \mathbf{H}) = \sum_{i=1}^N \sum_{g=1}^c \beta_i (u_{ig})^m \mathcal{D}(\mathbf{x}_i, \mathbf{h}_g),$$

¹Liu

²Butkiewicz

³Hung

⁴Zarandi

⁵Razaei

⁶Durso

⁷Giovanni

⁸Leski

⁹Huber

¹⁰Sigmoidal

¹¹Logarithmic

را ارایه کردند که در آن $\beta_i \in [0, 1]$ یک پارامتر مناسب است که تاثیر داده‌های دورافتاده را کاهش می‌دهد.

چنانچه \mathcal{L} تابع زیان مربعی باشد، فاصله‌ی پیشنهاد شده توسط دورسو و لسکی همان فاصله‌ی پیشنهاد شده توسط یانگ و کو است. تابع زیان هابر، برای مقدارهای کوچک $|e|$ ، مربعی و برای مقدارهای بزرگ آن، خطی است. در نتیجه نسبت به تابع زیان مربعی، کمتر تحت تاثیر داده‌های دورافتاده است. همچنین، بر خلاف تابع زیان خطی، مشتق‌پذیر بوده و کمینه‌سازی آن شدنی است. بنابراین تابع زیان هابر از مزایای هر دو تابع زیان مربعی و خطی بهره‌مند و فاقد معایب آن‌هاست.

وجه تمایز نتایج به دست آمده در این مقاله با نتایج دورسو و لسکی آن است که ما تاثیر استفاده از توابع زیان مختلف را با وزندهی یکنواخت به داده‌ها ($\beta_i = 1$) مورد مطالعه قرار می‌دهیم.

۴. تابع زیان پیشنهادی

در این بخش، ویژگی‌های تابع زیان پیشنهاد شده برای بکارگیری در مدل خوشه‌بندی فازی برای داده‌های فازی را تبیین می‌کنیم.

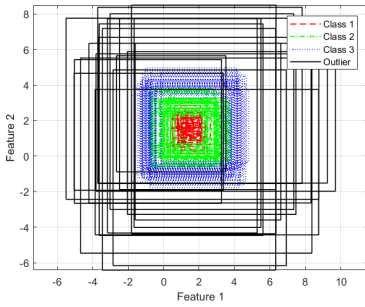
۱.۴. تابع زیان لگاریتم کسینوس هذلولوی. تابع زیان لگاریتم کسینوس هذلولوی یعنی

$$\mathcal{L}_{\text{Log-cosh}}(e) = \log(\cosh(e)),$$

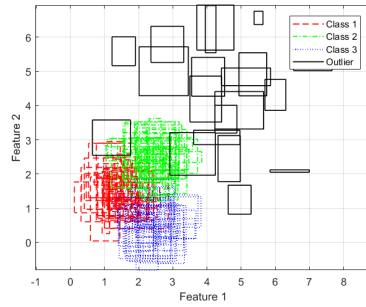
برای مقدارهای کوچک e ، تقریباً برابر $\frac{e^2}{4}$ و برای مقدارهای بزرگ آن، تقریباً برابر $|e| - \log(2)$ است. در نتیجه تحت تاثیر داده‌های دورافتاده قرار نمی‌گیرد. وجود مشتق دوم، مزیت دیگر این تابع زیان است که در تابع زیان هابر از آن بی‌بهره بودیم. بنابراین با استفاده از تابع زیان لگاریتم کسینوس هذلولوی، ضمن برخورداری از مزایای تابع زیان هابر، استفاده از روش‌هایی که برای بهینه‌سازی نیاز به مشتق دوم دارند نیز میسر است.

۵. مطالعه‌ی شبیه‌سازی

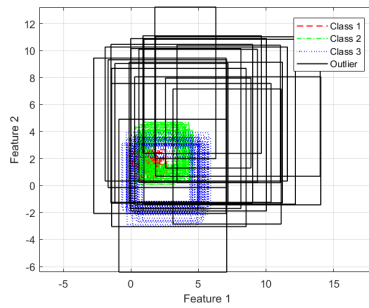
بزدک بر اساس تجربه $m \in [1/1, 5]$ را پیشنهاد کرد [۱]. مدل (۲.۳) با $\beta_i = 1$ و $m = 2$ و در نظر گرفتن توابع زیان مربعی، هابر با $\delta = 5$ ، خطی، سیگموییدی با $\alpha = \beta = 2$ ، لگاریتمی و لگاریتم کسینوس هذلولوی، در MATLAB R2021a پیاده‌سازی شده است.



(ب)



(آ)



(ج)

شکل ۴. مجموعه داده‌های فازی مثلثی تولید شده طبق سناریوهای (آ) مرکز دورافتاده، (ب) گستره‌های دورافتاده و (ج) مرکز و گستره‌های دورافتاده.

همه‌ی آزمایش‌ها بر روی یک ماشین (سیستم عامل: macOS، ظرفیت حافظه RAM: 8GB، پردازنده: 1.1GHz dual-core Intel Core i3) انجام شده است.

۱.۵. مجموعه داده‌های تولید شده. برای تولید مجموعه داده‌های فازی مثلثی، سه سناریوی مرکز دورافتاده مثل شکل ۴آ، گستره‌های دورافتاده مثل شکل ۴ب و مراکز و گستره‌های دورافتاده مثل شکل ۴ج در نظر گرفته شده است [۵، ۷، ۸]. هر مجموعه داده در \mathbb{R}^2 شامل ۱۲۰ نمونه در ۳ کلاس دارای هم‌پوشانی با تعداد اعضای برابر و یک گروه ۲۰ عضوی از داده‌های دورافتاده است. جدول ۱، سازوکار تولید مجموعه داده‌ها را بیان می‌کند.

۲.۵. مجموعه داده‌های معیار. مجموعه داده‌ی مصنوعی Zelnik4 [۱۸] در \mathbb{R}^2 که در شکل ۵ نشان داده شده است، شامل ۴۸۴ نمونه در ۴ کلاس و یک گروه ۱۳۸ عضوی از

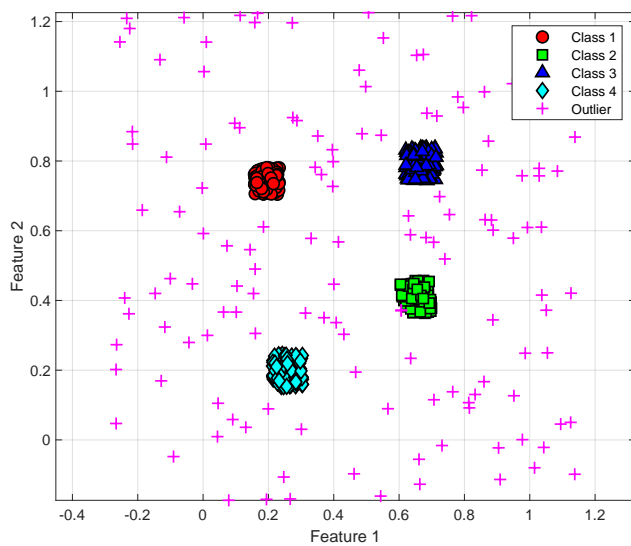
جدول ۱. سازوکار تولید مجموعه داده‌های فازی مثلثی طبق سناریوهای مرکز دورافتاده، گستره‌های دورافتاده و مرکز و گستره‌های دورافتاده.

سناریو	بُعد	مرکز	گستره‌ها
مرکز	کلاس ۱	۱ و ۲	$U[0,1]$
	کلاس ۲	۱ و ۲	$U[0,1]$
	کلاس ۳	۱	$U[2,3]$
		۲	$U[0,1]$
داده‌های دورافتاده	۱ و ۲	$\mathcal{N}(5,2)$	$U[0,1]$
گستره‌ها	کلاس ۱	۱ و ۲	$U[0,1]$
	کلاس ۲	۱ و ۲	$U[1,2]$
	کلاس ۳	۱ و ۲	$U[2,3]$
	داده‌های دورافتاده	۱ و ۲	$U[1,2]$
مرکز و گستره‌ها	کلاس ۱	۱ و ۲	$U[0,1]$
	کلاس ۲	۱ و ۲	$U[1,2]$
	کلاس ۳	۱	$U[2,3]$
		۲	$U[0,1]$
داده‌های دورافتاده	۱ و ۲	$\mathcal{N}(5,2)$	$\mathcal{N}(5,2)$

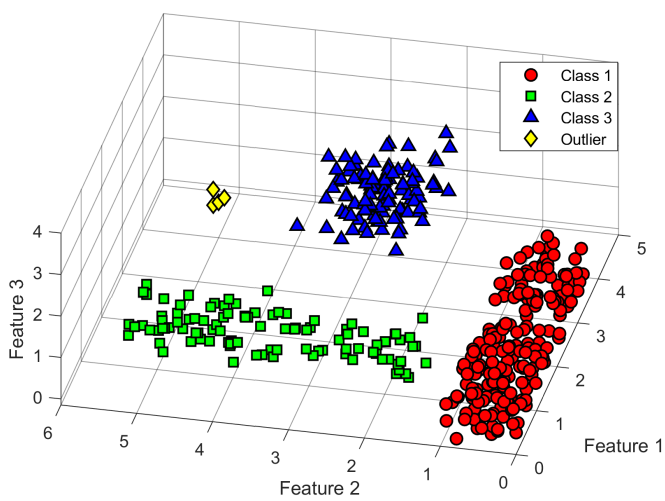
داده‌های دورافتاده است؛ تعداد ۱۱۱ نمونه با برچسب یک، ۱۱۴ نمونه با برچسب دو، ۱۵۰ نمونه با برچسب سه و ۱۰۹ نمونه با برچسب چهار.

مجموعه داده‌ی مصنوعی Lsun3D [۱۳] در \mathbb{R}^3 که در شکل ۶ نشان داده شده است، شامل ۴۰۰ نمونه در ۳ کلاس جدا از هم با اشکال هندسی متفاوت و یک گروه کوچک ۴ عضوی از داده‌های دورافتاده است؛ تعداد ۲۰۰ نمونه با برچسب یک، ۱۰۰ نمونه با برچسب دو و ۱۰۰ نمونه با برچسب سه.

برای فازی‌سازی مجموعه داده‌های Lsun3D و Zelnik4، هر دو مجموعه داده را به عنوان مرکز داده‌های فازی مثلثی در نظر گرفته‌ایم و فرض کرده‌ایم که گستره‌های این داده‌ها نیز از توزیع آماری $\mathcal{N}(1, 0/3)$ پیروی می‌کنند.



شکل ۵. مجموعه داده‌ی Zelnik4.



شکل ۶. مجموعه داده‌ی Lsun3D.

جدول ۲. نتایج بکارگیری توابع زیان مختلف در مدل خوشه‌بندی فازی برای داده‌های فازی.

میانگین	Lsun3D			Zelnik4			مرکز و گستره‌ها			مرکز و گستره‌ها			مجموعه داده‌ها		
	95% CI	میانگین	95% CI	میانگین	95% CI	میانگین	95% CI	میانگین	95% CI	میانگین	95% CI	میانگین	95% CI	میانگین	شاخص
۰/۶۸۶	[۰/۶۸, ۰/۷۲]	۰/۷۰	[۰/۷۵, ۰/۷۹]	۰/۷۷	[۰/۶۳, ۰/۶۸]	۰/۶۶	[۰/۵۸, ۰/۵۹]	۰/۵۸	[۰/۷۰, ۰/۷۳]	۰/۷۲	Frand				
۰/۴۷۶	[۰/۵۱, ۰/۵۵]	۰/۵۳	[۰/۵۰, ۰/۵۶]	۰/۵۳	[۰/۴۴, ۰/۵۱]	۰/۴۷	[۰/۲۸, ۰/۳۲]	۰/۳۰	[۰/۵۳, ۰/۵۸]	۰/۵۵	HUL				
۱/۷۵۳		۲۰/۷۶		۶۲/۲۶		۱/۴۴		۱/۱۴		۲/۰۷	Runtime				
۰/۶۹۸	[۰/۶۸, ۰/۷۱]	۰/۷۰	[۰/۷۵, ۰/۷۹]	۰/۷۷	[۰/۶۶, ۰/۷۱]	۰/۶۹	[۰/۵۸, ۰/۶۰]	۰/۵۹	[۰/۷۲, ۰/۷۶]	۰/۷۴	Frand				
۰/۴۹۴	[۰/۵۰, ۰/۵۵]	۰/۵۲	[۰/۵۲, ۰/۵۷]	۰/۵۴	[۰/۴۹, ۰/۵۵]	۰/۵۲	[۰/۲۹, ۰/۳۳]	۰/۳۱	[۰/۵۶, ۰/۶۱]	۰/۵۸	HUL				
۱/۸۳۵		۲۲/۷۵		۶۳/۹۴		۱/۵۳		۱/۴۱		۲/۱۲	Runtime				
۰/۶۱۴	[۰/۶۱, ۰/۶۴]	۰/۶۲	[۰/۶۳, ۰/۶۶]	۰/۶۵	[۰/۶۱, ۰/۶۴]	۰/۶۳	[۰/۵۳, ۰/۵۴]	۰/۵۳	[۰/۶۲, ۰/۶۵]	۰/۶۴	Frand				
۰/۳۰۸	[۰/۳۵, ۰/۳۹]	۰/۳۷	[۰/۲۷, ۰/۳۰]	۰/۲۹	[۰/۳۳, ۰/۳۷]	۰/۳۵	[۰/۱۶, ۰/۱۸]	۰/۱۷	[۰, ۰/۳۴, ۰/۳۸]	۰/۳۶	HUL				
۱/۱۸۷		۱۳/۸۳		۴۱/۱۸		۱/۴۵		۱/۱۴		۱/۷۵	Runtime				
۰/۶۰۶	[۰/۶۳, ۰/۶۶]	۰/۶۵	[۰/۵۳, ۰/۵۴]	۰/۵۳	[۰/۶۴, ۰/۶۸]	۰/۶۶	[۰/۵۲, ۰/۵۳]	۰/۵۳	[۰/۶۵, ۰/۶۷]	۰/۶۶	Frand				
۰/۲۹۲	[۰/۴۱, ۰/۴۵]	۰/۴۳	[۰/۰۳, ۰/۰۳]	۰/۰۳	[۰/۳۹, ۰/۴۴]	۰/۴۱	[۰/۱۴, ۰/۱۶]	۰/۱۵	[۰/۴۲, ۰/۴۵]	۰/۴۴	HUL				
۱۳/۶۹		۱۸/۰۸		۴۵/۰۸		۱/۵۰		۱/۴۲		۲/۳۶	Runtime				
۰/۶۹۴	[۰/۶۲, ۰/۶۵]	۰/۶۴	[۰/۷۵, ۰/۷۹]	۰/۷۷	[۰/۶۹, ۰/۷۳]	۰/۷۱	[۰/۵۸, ۰/۶۰]	۰/۵۹	[۰/۷۴, ۰/۷۸]	۰/۷۶	Frand				
۰/۴۶۸	[۰/۴۰, ۰/۴۴]	۰/۴۲	[۰/۵۱, ۰/۵۶]	۰/۵۳	[۰/۴۹, ۰/۵۴]	۰/۵۱	[۰/۲۸, ۰/۳۱]	۰/۳۰	[۰/۵۵, ۰/۶۰]	۰/۵۸	HUL				
۱۹/۴۶		۱۸/۰۸		۷۳/۰۵		۱/۹۰		۱/۳۴		۲/۹۴	Runtime				
۰/۶۹۶	[۰/۶۵, ۰/۶۸]	۰/۶۷	[۰/۷۴, ۰/۷۷]	۰/۷۶	[۰/۶۷, ۰/۷۱]	۰/۶۹	[۰/۵۸, ۰/۶۰]	۰/۵۹	[۰/۷۵, ۰/۷۹]	۰/۷۷	Frand				
۰/۴۸۰	[۰/۴۵, ۰/۴۹]	۰/۴۷	[۰/۴۹, ۰/۵۴]	۰/۵۲	[۰/۴۸, ۰/۵۳]	۰/۵۱	[۰/۲۹, ۰/۳۲]	۰/۳۰	[۰/۵۸, ۰/۶۳]	۰/۶۰	HUL				
۱۹/۶۵		۲۲/۲۱		۶۹/۸۴		۱/۸۱		۱/۴۵		۲/۹۴	Runtime				

۳.۵. نتایج. میانگین شاخص رند فازی^۱ Frand [۳] و میانگین شاخص هولرمایر^۲ HUL [۹] و بازه اطمینان آن‌ها و همچنین میانگین زمان اجرا برای ۱۰۰ بار تکرار آزمایش در جدول ۲ گزارش شده است. هرچه مقدار به دست آمده برای شاخص‌های Frand و HUL به یک نزدیک‌تر و زمان اجرا کمتر باشد، نتیجه‌ی خوشه‌بندی، رضایت‌بخش‌تر خواهد بود.

۶. نتیجه‌گیری

بکارگیری تابع زیان لگاریتم کسینوس هذلولوی در مدل خوشه‌بندی فازی برای داده‌های فازی را پیشنهاد و تاثیر استفاده از این تابع و نیز توابع زیان مربعی، هابر، خطی، سیگموئیدی و لگاریتمی بر عملکرد مدل را در صورت وجود داده‌های دورافتاده در مجموعه داده‌ها، بررسی کرده‌ایم.

شکل ۷ نمودارهای توابع زیان مربعی، هابر، خطی، سیگموئیدی، لگاریتمی و لگاریتم کسینوس هذلولوی را مقایسه می‌کند. به طور کلی، آن دسته از توابع زیان که

(۱) در نزدیکی $e = 0$ شیب کمی داشته باشند؛

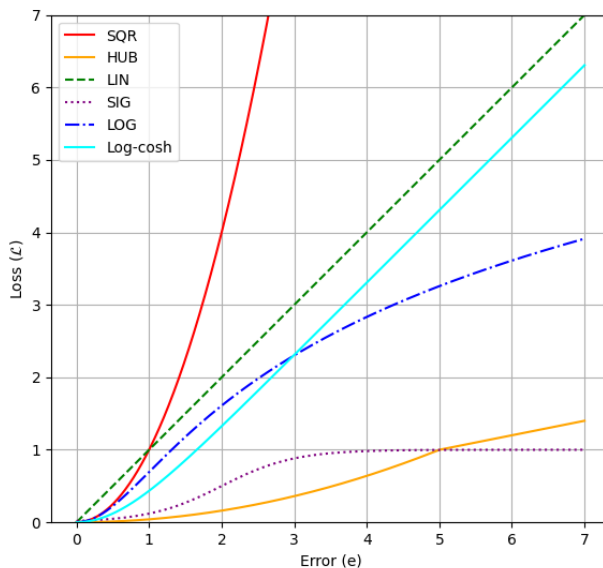
(۲) با افزایش خطا و دور شدن از نقطه‌ی $e = 0$ ، شیب تابع با آهنگ ملایمی افزایش یابد؛

(۳) در صورت وجود داده‌های دورافتاده و افزایش چشمگیر مقدار خطا، زیان مربوطه را با بزرگ‌نمایی محاسبه نکنند؛

کارایی بهتری خواهند داشت. نتایج شبیه‌سازی نیز موید آن است که بهترین عملکرد روش خوشه‌بندی، برای توابع زیان هابر و لگاریتم کسینوس هذلولوی به دست آمد. اشکال تابع زیان خطی، شیب زیاد آن در نزدیکی $e = 0$ است. تابع زیان مربعی نیز برای $e > 1$ ، زیان زیادی را محاسبه می‌کند و به همین دلیل، نسبت به وجود داده‌های دورافتاده مقاوم نیست. کمترین زمان اجرای روش خوشه‌بندی، برای توابع زیان خطی و سیگموئیدی به دست آمد.

¹Fuzzy rand

²Hullermeier



شکل ۷. مقایسه‌ی نمودارهای توابع زیان مربعی، هابر با $\delta = 5$ ، خطی، سیگموییدی با $\alpha = \beta = 2$ ، لگاریتمی و لگاریتم کسینوس هذلولوی.

مراجع

- [1] Bezdek, J. (1981) Pattern recognition with fuzzy objective function algorithms. *Plenum Press, New York*.
- [2] Butkiewicz, B. (2005) Robust fuzzy clustering with fuzzy data. *Advances in Web Intelligence, Third International Atlantic Web Intelligence Conference, AWIC 2005, Lecture Notes in Computer Science, Springer, Berlin, Heidelberg*, 3528, 76–82.
- [3] Campello, R. (2007) A fuzzy extension of the Rand index and other related indexes for clustering and classification assessment. *Pattern Recognition Letters*, 28 (7), 833–841.
- [4] D’Urso, P. (2007) Fuzzy clustering of fuzzy data. *Advances in Fuzzy Clustering and Its Applications*, 155–192.
- [5] D’Urso, P., De Giovanni, L. (2014) Robust clustering of imprecise data. *Chemometrics and Intelligent Laboratory Systems*, 136, 58–80.

- [6] D'Urso, P., Leski, J. (2020) Fuzzy clustering of fuzzy data based on robust loss functions and ordered weighted averaging. *Fuzzy Sets and Systems*, **389**, 1–28.
- [7] Eskandari, E., Khastan, A. (2023) A robust fuzzy clustering model for fuzzy data based on an adaptive weighted L^1 norm. *Iranian Journal of Fuzzy Systems*, **20** (6), 1-20.
- [8] Eskandari, E., Khastan, A., Tomasiello, S. (2022) Improved determination of the weights in a clustering approach based on a weighted dissimilarity measure between fuzzy data. *2022 IEEE Int. Conf. Fuzzy Syst. (FUZZ-IEEE), Padua, Italy*, 1-6.
- [9] Hullermeier, E., Rifqi, M., Henzgen, S., Senge, R. (2012) Comparing fuzzy partitions: a generalization of the Rand index and related measures. *IEEE Transactions on Fuzzy Systems*, **20** (3), 546–556.
- [10] Hung, W., Yang, M. (2005) Fuzzy clustering on LR-type fuzzy numbers with an application in Taiwanese tea evaluation. *Fuzzy Sets and Systems*, **150** (3), 561–577.
- [11] Jin, X., Han, J. (2010) K-medoids clustering. *Encyclopedia of Machine Learning*, 697-700.
- [12] Sato, M., Sato, Y. (1995) Fuzzy clustering model for fuzzy data. *Proceedings of 1995 IEEE International Conference on Fuzzy Systems*, **4**, 2123–2128.
- [13] Thrun, M. (2018) Projection-based clustering through self organization and swarm intelligence. *Springer Vieweg, Wiesbaden*.
- [14] Urbanowicz, R.J. [@DocUrbs]. (2018, Jun 15). *New proposed field/term Venn diagram for an upcoming talk. My take on illustrating the relationship between #Data-Science, #MachineLearning, #ArtificialIntelligence, #Statistics* [Tweet]. Twitter. <https://twitter.com/DocUrbs/status/1007375834347376642>
- [15] Yang, M., Ko, C. (1996) On a class of fuzzy c-numbers clustering procedures for fuzzy data. *Fuzzy Sets and Systems*, **84** (1), 49–60.
- [16] Yang, M., Liu, H. (1999) Fuzzy clustering procedures for conical fuzzy vector data. *Fuzzy Sets and Systems*, **106** (2), 189–200.
- [17] Zarandi, M., Razaee, Z. (2010) A fuzzy clustering model for fuzzy data with outliers. *International Journal of Fuzzy System Applications (IJFSA)*, **1** (2), 29–42.
- [18] Zelnik-Manor, L., Perona, P. (2004) Self-tuning spectral clustering. *Proceedings of the 17th International Conference on Neural Information Processing Systems*, **17**, 1601–1608.