

درس : کلان داده

نام و نام خانوادگی : فاطمه توکلی

۴۰۰۱۳۱۰۱۶

گزارش پایانی

ابتدا به صورت خلاصه به بیان مسئله و روش ارائه شده در آن میپردازیم:

در اینجا به بیان مسئله دسته بندی برای داده‌های غیر متوازن

خوشه بندی داده‌ها نامتوازن یک مشکل چالش برانگیز در یادگیری ماشین است. مشکل اصلی ناشی از عدم تعادل در اندازه خوشه و توزیع چگالی داده است. برای پرداختن به این مشکل، یک الگوریتم خوشه‌بندی جدید به نام LDPI را بر اساس پیک‌های چگالی محلی در این مقاله پیشنهاد شده است.

الگوریتم پیشنهادی LDPI شامل سه مرحله است:

(۱) یک طرح تولید زیر خوشه اولیه طراحی شده است که می‌تواند به طور خودکار نقاط نویز و مراکز زیر خوشه اولیه را شناسایی کند. خوشه‌های فرعی با طبقه بندی نقاط در مراکز زیر خوشه بر اساس اصل نزدیکترین همسایه ایجاد می‌شوند. (۲) یک روش به روز رسانی زیرخوشه معرفی شده است که می‌تواند مراکز زیر خوشه کاذب را از مراکز زیر خوشه اولیه شناسایی کند. متعاقباً، طرح آنها را حذف می‌کند تا مراکز فرعی به روز شده و زیرخوشه‌های به روز شده مربوطه را به دست آورد.

(۳) یک استراتژی ادغام زیرخوشه برای ادغام زیرخوشه‌های به روز شده و ایجاد خوشه‌های نهایی اعمال می‌شود.

(۱) تعیین آستانه فاصله مناسب  $d_c$  و چگالی

همانطور که قبلاً بحث شد، تخصیص دستی  $d_c$  در DPC معمولاً منجر به خوشه‌های نادرست، به ویژه با مجموعه داده‌های نامتعادل می‌شود. در راستای حل این مشکل برای یک مجموعه داده با  $n$  نقطه، اگر  $D_i$  فاصله بین نقطه  $i$  و نزدیکترین همسایه آن در نظر بگیریم، سپس مقدار حداکثر تمام فواصل بین هر نقطه و نزدیکترین همسایه آن به صورت زیر محاسبه می‌شود:

$$d_c = \max_{i=1}^n \{D_i\}$$

که این را به عنوان آستانه فاصله  $d_c$  در نظر می‌گیرد. موضوع دیگر، تعریف چگالی مناسب نقاط است.

در DPC چند نقطه ممکن است چگالی یکسانی داشته باشند و در نسخه بهبود یافته آن از معادله‌ای استفاده می‌کند که برای محاسبه چگالی هر نقطه بایستی تمامی نقاط در معادله گذاشته شوند و چگالی به دست آید.

در اینجا، برای غلبه بر اشکال روش بهبود یافته، به جای تمام نقاط موجود، چگالی محلی هر نقطه را با نقاط مجاور آن تعیین می‌کند، چگالی محلی هر نقطه را به صورت زیر محاسبه می‌کند

$$\rho_i = \sum_{j \in S/\{i\}, d_{ij} \leq d_c} e^{-\left(\frac{d_{ij}}{d_c}\right)^2}$$

(۲) تعیین نقاط نویز توسط یک گراف تصمیم جدید

نقاط نويز نمودار تصميم تاثير منفي قوي بر نتايج خوشه بندي دارند. بنابر اين بايد آنها را شناسايي کرد. در DPC، نقاطي با چگالي کمتر و معيار upward distance بيشتر به عنوان نقاط نويز در نظر گرفته مي شوند، که upward distance براي هر نقطه  $X_i$  صورت زير تعريف مي شود:

$$\delta_i = \begin{cases} \max_j \{d_{ij}\}, & \text{if } \rho_i > \rho_j \text{ for } \forall j \\ \min_j \{d_{ij} | \rho_j > \rho_i\}, & \text{otherwise.} \end{cases}$$

اما با مجموعه داده هاي نامتعادل، برخي از مراکز خوشه کوچک به اشتباه به عنوان نقاط نويز اختصاص داده مي شوند. براي شناسايي صحيح نقاط نويز در يک مجموعه داده نامتعادل، يک نمودار تصميم گيري جديد معرفي شده که از تعداد معکوس نزديکترين همسايگان هر نقطه استفاده مي کند. و در صورتي که سه شرط زير برقرار بود نقطه به عنوان نويز يا داده پرت در نظر گرفته مي شود:

- 1)  $\delta_i > \mu(\delta) + \sigma(\delta).$
- 2)  $\rho_i < \mu(\rho) - \sigma(\rho).$
- 3)  $RNN_i < \mu(RNN) - \sigma(RNN).$

که  $\mu$  و  $\delta$  در آن به ترتيب ميانيگين و انحراف معيار تمام نقاط براي  $RNN$ , densities, upward distance هستند.

### (۳) طرح اوليه توليد زير خوشه

براي جلوگيري از تخصيص تعداد خوشه ها به صورت دستي و فعال کردن الگوريتم براي خوشه بندي خودکار روي مجموعه داده هاي نامتعادل، ابتدا تعداد نسبتاً زيادي از مراکز زير خوشه را تنظيم کرده و سپس زير خوشه هاي اوليه را ساختيم. به طور دقيق تر پس از حذف نقاط نويز، نقاط باقيمانده که فاصله آنها بيشتر  $\mu(\delta) + \sigma(\delta)$  است به عنوان مراکز زير خوشه اوليه انتخاب مي شوند. از اين طريق اطمينان حاصل شود که فواصل بين مراکز مختلف به اندازه کافي بزرگ است. پس از آن، هر مرکز زير خوشه اوليه  $i$  به يک زير خوشه اوليه  $C_i$  مربوط مي شود. تمام نقاط باقيمانده به ترتيب نزولي چگالي آنها مرتب مي شوند. و به نزديک مرکز که چگالي آن از چگالي خود نقطه بيشتر است انتصاب مي يابند.

### (۴) به روزرساني زير خوشه

ابتدا، زير خوشه هاي کاذب را به صورت زير شناسايي مي کنيم: فرض کنيد که  $A_i$  مرکز زير خوشه اوليه  $C_i$  حاوي  $NC_i$  نقطه است، و همسايگي مرکز زير خوشه  $i$  با شعاع  $dc$  حاوي  $Ndc_i$  نقطه است. اگر  $C_i$  حاوي کمتر از نيمي از نقاط  $Ndc_i$ ، يعني  $NC_i < 0.5 Ndc_i$  باشد، به اين معنيست که چنين خوشه اي يک خوشه واقعي نيست و بايد حذف شود. سپس، با تخصيص نقاط در زير خوشه هاي نادرست به نزديکترين زير خوشه هاي همسايه، زير خوشه ها به روز مي شوند.

### (۵) استراتژي ادغام خوشه هاي فرعي

ابتدا نقاط مرزي هر زير خوشه مشخص مي شود. نقاط مرزي به عنوان نقاطي تعريف مي شوند که چگالي آنها کمتر از چگالي متوسط زير خوشه داده شده است.

برای ادغام زیرخوشه‌های به روز شده، اگر  $S$  مجموعه داده را نشان دهد،  $No(X)$  مجموعه‌ای از نقاط نويز را در  $S$  و  $In(X)$  نشانگر زیرمجموعه‌ای از  $X$  به استثنای تمام نقاط مرزی  $X$  باشد. شعاع ادغام  $r$  با معادله زیر تعریف می‌شود:

$$r = \begin{cases} d_c, & \text{if } No(S) = \emptyset \\ \max_{x_i \in S/No(S)} \{D_i\}, & \text{otherwise.} \end{cases}$$

و با استفاده از این شعاع و شروط زیر برای ادغام به خوشه بندی نهایی میرسد :

برای ادغام زیر خوشه های  $C_m$  و  $C_n$ ، باید نشان دهیم  $d(C_m, C_n) = \min\{d_{ij} \mid x_i \in C_m, x_j \in C_n\}$  و دو مورد باید در نظر گرفته شود:

(۱) اگر  $d(C_m, C_n) > r$  سپس  $C_m$  و  $C_n$  دور هستند و ادغام نمی‌شوند.

(۲) در غیر این صورت

(i) اگر  $d(C_m, C_n) = \min\{d_{ij} \mid x_i \in In(C_m), x_j \in In(C_n)\}$ ، سپس  $C_m$  و  $C_n$  بسیار نزدیک هستند و می‌توانند مستقیماً ادغام شوند. همانطور که در شکل ۲۸ نشان داده شده است، زیر خوشه های ۱ و ۲ شرایط فوق را برآورده می‌کنند و ادغام می‌شوند. به همین ترتیب، زیر خوشه های ۳ و ۴ نیز شرایط فوق را برآورده می‌کنند و ادغام می‌شوند.

(ii) در غیر این صورت،

الف)  $(x_t, x_s) = \operatorname{argmin}_{x_t \in C_m, x_s \in C_n} (d_{ij})$  نشان می‌دهد. اگر مجموع چگالی  $x_t$  و  $x_s$  از میانگین چگالی مراکز زیر خوشه آنها بزرگتر باشد، یعنی اگر  $\rho_t + \rho_s > \rho_{C_{mn}}$  سپس  $C_m$  و  $C_n$  ادغام می‌شوند.

ب) در غیر این صورت، دو خوشه فرعی با هم ادغام نمی‌شوند.

این فرآیند ادغام تا زمانی تکرار می‌شود که هیچ جفتی از زیر خوشه‌ها شرایط ادغام را برآورده نکنند. پس از فرآیند ادغام، اگر نقاط نويز باقی‌نماند، خوشه بندی کامل می‌شود. در غیر این صورت، هر نقطه نويز به نزدیکترین خوشه خود اختصاص داده می‌شود.

در نهایت با پرداختن به جزئیات بیان می‌کند که حل این مسئله با پیچیدگی زمانی  $O(n^2)$  ممکن است.

در نتیجه به عنوان مزایا میتوان به این موارد اشاره کرد که :

(۱) به هیچ پارامتر ورودی نیاز ندارد.

(۲) می‌تواند به طور خودکار مراکز خوشه و تعداد خوشه‌ها را تعیین کند.

۳) برای مجموعه داده ها و مجموعه داده های نامتعادل با اشکال و توزیع های دلخواه مناسب است