



دانشگاه صنعتی امیرکبیر
(پلی تکنیک تهران)

دانشکده مهندسی کامپیوتر

گزارش سمینار داده کاوی

استاد راهنما:

دکتر امیرمزلقانی

دانشجویان:

علی نظری - ۹۶۳۱۰۷۵

عرفان عابدی - ۹۶۳۱۰۴۲۷

زمستان ۱۳۹۹

صفحه

فهرست مطالب

مقاله‌ی اول Categorical Data Sets in Data Mining A Fast Clustering Algorithm to Cluster Very Large ۱

۱-۱ - چکیده ۲

۱-۲ - معرفی ۲

۱-۳ - الگوریتم k-modes ۳

۱-۳-۱ - معیار عدم شباهت ۴

۱-۳-۲ - پیدا کردن mode دیتاست ۴

۱-۳-۳ - روند اجرای الگوریتم ۴

مقاله‌ی دوم LoOP: Local Outlier Probabilities ۶

۲-۱ - مقدمه ۷

۲-۲ - تعریف ریاضی ۷

مقاله‌ی اول

**A Fast Clustering Algorithm to Cluster Very Large
Categorical Data Sets in Data Mining**

۱-۱- چکیده

یکی از معروف‌ترین الگوریتم‌های خوشه‌بندی، k-means است که برای داده‌های عددی به نسبت پاسخ مناسبی می‌دهد ولی در داده‌هایی که مقادیرشان categorical بوده، پاسخ مناسبی را به همراه ندارد. در این مقاله الگوریتم k-modes را معرفی می‌کنیم تا بتوانیم این نوع داده‌ها را نیز خوشه‌بندی کنیم و از لحاظ مقیاس‌پذیری هم عملکرد مناسبی داشته باشیم.

۱-۲- معرفی

روش‌های خوشه‌بندی به صورت کلی بسیار کمک‌کننده هستند و ما را در جهت خوشه‌بندی داده‌های حجیم در خوشه‌هایی همگن بسیار یاری می‌کنند تا بتوانیم بر روی آن‌ها آنالیز بهتری انجام دهیم. همان‌طور که می‌دانیم، در داده‌کاوی با داده‌های حجیمی مواجه هستیم پس مقیاس‌پذیری الگوریتم‌های مورد استفاده نیز برای ما اهمیت دارد.

در این جا الگوریتمی ارائه می‌شود که روی داده‌های حجیم عملکرد مناسبی دارد و می‌تواند در خوشه‌بندی داده‌های categorical ما را یاری کند. یکی از روش‌هایی که برای خوشه‌بندی داده‌های categorical استفاده می‌شده به این صورت بوده که این داده‌ها را به مقادیر عددی تبدیل می‌کردیم و با آن‌ها مثل مقادیر عددی برخورد می‌کردیم ولی این نوع خوشه‌بندی نتایج معناداری به ما نمی‌داده مخصوصاً زمانی که این categoryها ترتیب خاصی نداشتند. الگوریتم k-modes که ارائه می‌شود، تلاش می‌کند که علاوه بر حفظ مزایای الگوریتم k-means، محدودیت‌های این روش در خوشه‌بندی داده‌های categorical را نیز مرتفع سازد.

پیش از این، الگوریتم دیگری تحت عنوان k-prototype ارائه شده بوده که در آن هم معیاری برای عدم شباهت ویژگی‌های عددی وجود داشته و هم معیاری برای عدم شباهت داده‌های categorical در آن وجود داشته است. در قسمت عددی، فاصله‌ی اقلیدسی محاسبه می‌شده و برای قسمت categorical هم به این صورت بوده که تعداد ویژگی‌های categorical که بین دو داده با هم متفاوت بودند رو محاسبه می‌کرده و در نهایت نیز این دو معیار را با هم جمع می‌کرده است. البته برای اینکه در

این جمع، توازن برقرار شود؛ یک وزن به معیار عدم شباهت در قسمت categorical داده می‌شده که مشکل این روش انتخاب همین وزن مناسب بوده است.

الگوریتم k-modes که در اینجا ارائه می‌شود تغییراتی در الگوریتم k-prototype داده‌است به این صورت که تمامی ویژگی‌های داده‌ها را categorical فرض کرده و اگر مقادیر عددی وجود داشت نیز آن‌ها را به صورت categorical در می‌آورد و مهم‌ترین ویژگی آن مقیاس‌پذیر بودن آن است.

البته یک روش دیگر نیز برای خوشه‌بندی داده‌های categorical با کمک k-means وجود داشته به این صورت که تمامی categoryها را به صورت یک ویژگی دودویی در می‌آورد و بعد به این صورت عمل می‌کرده که آیا این داده دارای این category می‌باشد یا خیر و بعد با این ویژگی‌ها به صورت عددی (۰ یا ۱) برخورد می‌کرده که البته مشکلات زیادی داشته‌است؛ به عنوان مثال وقتی که این categoryها زیاد باشند، ابعاد بسیار زیادی به دیتاست داده‌ها اضافه می‌شود. مورد دیگری که وجود داشته این بوده که محاسبه‌ی میانگین بین تعداد زیادی ۰ و ۱ می‌توانسته بی‌معنی باشد و نتوان آن را به عنوان یک مشخصه در نظر گرفت.

استفاده از خوشه‌بندی سلسله‌مراتبی نیز برای این دیتاست‌ها به علت حجم بالای داده‌ها، پیشنهاد نمی‌شده؛ پس استفاده از k-modes پیشنهاد می‌شود زیرا می‌توان از خود داده‌ها و همان categoryها استفاده کرد و در عین حال مقیاس‌پذیری را نیز حفظ کرد.

۳-۱- الگوریتم k-modes

رویکرد این الگوریتم به k-means شباهت زیادی دارد و البته نمونه‌ی ساده شده‌ی k-prototype نیز محسوب می‌شود. این الگوریتم سه تفاوت اصلی با الگوریتم k-means دارد. اول، معیار عدم شباهت است که در این الگوریتم فرق دارد، دوم اینکه mean را با mode جایگزین کردیم و سوم اینکه نحوه بروزرسانی modeها است که به کمک روشی بر پایه فرکانس categoryها انجام می‌شود.

۱-۳-۱- معیار عدم شباهت

این معیار به این صورت است که تمامی ویژگی‌های بین دو نمونه را نگاه می‌کند و آن تعداد ویژگی‌هایی که category یکسانی ندارند را به عنوان عدم شباهت، محسوب می‌کند. البته این روش به صورتی که تعریف شد، دارد به تمامی categoryها وزن یکسانی می‌دهد؛ پس حالا یک وزنی هم ایجاد شده و برای هر مقایسه‌ای این وزن را هم درگیر می‌کند، به این صورت که مجموع معکوس فرکانس این دو category را در معیار عدم شباهت قبلی‌ای که محاسبه کردیم، ضرب می‌کنیم. اگر دقت کنیم، می‌بینیم که ما با این کار داریم وزن بیشتری به داده‌هایی که فرکانس کمتری دارند می‌دهیم و مثلاً در جایی که می‌خواهیم کلاهبرداری را در داده‌های بیمه پیدا کنیم، به کار ما می‌آیند.

۱-۳-۲- پیدا کردن mode دیتاست

Mode دیتاست، وکتوری است که ما معیار عدم شباهت داده‌ها را نسبت به آن اندازه‌گیری می‌کنیم و تلاش می‌کنیم که این فاصله و معیار عدم شباهت را کمینه کنیم. این تابع وقتی کمینه می‌شود که تعداد تکرار ویژگی‌های مشابه با mode در داده‌ها از سایر categoryهای داده‌ها بیشتر باشد.

۱-۳-۳- روند اجرای الگوریتم

روند اجرای الگوریتم به این صورت است که دیتاست را به بخش‌هایی (خوشه) افراز می‌کنیم و mode هر کدام از این بخش را انتخاب می‌کنیم. تابع هزینه‌ای به وجود می‌آید که حاصل از ترکیب جمع عدم شباهت‌های هر داده با mode همان خوشه‌ای است که در آن قرار گرفته است.

اگر بخواهیم موارد بالا را بیشتر توضیح دهیم، می‌توانیم این‌گونه بگوییم که:

- K خوشه را که از افراز دیتاست به وجود آمده را در نظر می‌گیریم.
- برای هر خوشه یک وکتور به عنوان mode در نظر می‌گیریم (که در ادامه نحوه انتخاب آن را بیشتر توضیح می‌دهیم)

- به سراغ داده‌ها می‌رویم و اولین داده را به هر کدام از خوشه‌ها با محاسبه‌ی معیار عدم شباهت، نسبت می‌دهیم و mode را برابر همان داده‌ی اول در نظر می‌گیریم و بعد از آن به سراغ داده‌های دیگر می‌رویم و آن‌ها به خوشه‌ی نزدیک‌تر مربوط می‌کنیم و پس از هر نسبت‌دهی، mode را بروزرسانی می‌کنیم. این روند را تا زمانی که تمامی داده‌ها به یک خوشه نسبت داده شوند، ادامه می‌دهیم.

- زمانی که تمامی داده‌ها به یک خوشه نسبت داده شدند، دوباره تمامی داده‌ها را بررسی می‌کنیم که آیا با بروزرسانی mode لازم است که در خوشه‌ی دیگری قرار گیرند یا خیر و این روند را تا زمانی که دیگر تغییری در این نسبت بین داده‌ها و خوشه‌ها اتفاق نیفتد ادامه می‌دهیم.

این الگوریتم نیز به مانند k-means ممکن است در مینیمم محلی بیفتد و این بستگی به حالتی دارد که با آن، الگوریتم را آغاز می‌کنیم.

همان‌طور که اشاره شد، روش‌هایی برای انتخاب mode برای شروع الگوریتم وجود دارد که در اینجا به دو تا از آن‌ها اشاره شده است.

اول اینکه k داده‌ی اول را در نظر بگیریم و هر کدام از آن‌ها را به یکی از خوشه‌ها به عنوان mode نسبت دهیم.

دوم اینکه فرکانس هر کدام از category‌های هر ویژگی داده را محاسبه کنیم و آن‌ها را به صورت نزولی مرتب کنیم و بعد به صورت مساوی این‌ها را به عنوان mode به یک خوشه نسبت دهیم و بعد به سراغ داده‌ها می‌رویم و نزدیک‌ترین داده به هر خوشه را پیدا می‌کنیم و آن داده را با mode‌ی که قبلاً انتخاب کردیم، جایگزین می‌کنیم زیرا این mode‌ها به صورت مصنوعی انتخاب شده بودند و ممکن بود موجب ایجاد خوشه‌های خالی شوند.

این الگوریتم مقیاس‌پذیری بالایی هر در تعداد خوشه‌ی بالا و هم در تعداد داده‌ی زیاد در خوشه‌های زیاد دارد و با افزایش تعداد این پارامترها، مقیاس‌پذیری خود را از دست نمی‌دهد و به صورت خطی زمان اجرای الگوریتم با افزایش این مقادیر، افزایش می‌یابد.

مقاله‌ی دوم

LoOP: Local Outlier Probabilities

۱-۲- مقدمه

شیوه‌های متفاوتی برای مساله‌ی تشخیص داده‌های Outlier ارائه شده‌اند. این شیوه‌ها به دو دسته‌ی کلی (Global) و محلی (Local) تقسیم می‌شوند. شیوه‌های کلی برای تشخیص Outlierها متکی به بررسی کل داده‌های موجود هستند، در حالی که شیوه‌های محلی به بررسی داده‌های نزدیک داده‌ی مورد بررسی اکتفا می‌کنند. (روش‌هایی مشابه با k-nearest neighbors). خصیصه‌ی دیگری که در عمده‌ی این شیوه‌ها دیده می‌شود این است که شیوه‌های کلی صفر و یکی عمل کرده و داده‌های مورد بررسی را به دو دسته‌ی «Outlier» و «داده‌ی عادی» تقسیم‌بندی می‌کنند، در حالی که شیوه‌های محلی به داده‌ی مورد بررسی یک فاکتور «Outlierness» می‌دهند که نشان‌دهنده‌ی میزان حدودی Outlier بودن داده‌ی مورد بررسی است. هرچند که این امر همیشه صادق نیست.

۲-۲- تعریف ریاضی

در این مقاله، مدل تشخیص Outlier جدیدی معرفی شده که ایده‌های فاکتورهای محلی و مبتنی-بر-چگالی همانند LOF (و واریانت‌های آن، مثل LOCI) را با مفاهیم احتمالاتی ترکیب می‌کند تا میزان Outlier بودن یک داده را مدل کند. مزیت یک مدل احتمالاتی ایجاد نسبی یک تحمل (Tolerance) نسبت به داده‌های نویز است. مدل‌های سنتی حتی در مواردی تاثیر داده‌های نویز را زیاد می‌کنند. به عنوان مثال، متد LOF مبتنی بر محاسبه‌ی فاصله‌ی k تا داده‌ی نزدیک نسبت به داده‌ی فعلی است و انتخاب اشتباه عدد k می‌تواند منجر به نتایج ناپایدار شود. در ادامه، D به معنای مجموعه‌ای از n داده است و d به معنای تابع فاصله‌ای که برای تشخیص Outlierها محاسبه می‌شود. برای پایدارتر شدن نتایج محاسبات، مفهومی به نام «فاصله‌ی احتمالاتی» یا $\text{pdist}(o, S)$ معرفی می‌شود که پارامتر o در آن یکی از داده‌های موجود در مجموعه‌ی D است و S یک مجموعه‌ی محلی که زیرمجموعه‌ای از D است. این فاصله، باید دارای خصوصیت زیر باشد:

$$\forall s \in S: P[d(o, S) \leq pdist(o, S)] \geq \psi$$

این خصوصیت به این معناست که یک کره حول o با شعاع $pdist(o, S)$ هر داده‌ای در S را با احتمال ψ پوشش می‌دهد. فاصله‌ی $pdist(o, S)$ می‌تواند به عنوان «محدوده‌ی احتمالاتی» مجموعه‌ی محلی S در نظر گرفته شود. تفاوت اصلی این نوع محدوده با محدوده‌ی نرمال این است که محدوده‌ی احتمالاتی از عمد رخداد برخی خطاها را جایز می‌شمرد. معکوس فاصله‌ی احتمالاتی می‌تواند به عنوان تخمینی برای چگالی محدوده‌ی S در نظر گرفته شود:

$$pdens(S) = \frac{1}{pdist(o, S)}$$

اگر به جای ψ از $\lambda = \sqrt{2} \operatorname{erf}^{-1}(\psi)$ استفاده شود که در آن، erf نشان‌دهنده‌ی تابع ارور گاوسی است، در حین تخمین چگالی S ، می‌توانیم Outlier ها را به معنای داده‌هایی تفسیر کنیم که فاصله‌ی آن‌ها با میانگین S ، بیش‌تر از $\lambda \cdot \sigma$ است که σ در این جا به معنای انحراف معیار است. مقادیر λ توسط قانون تجربی ۳-سیگما تعیین می‌شوند، به این معنا که:

$$\lambda = 1 \Rightarrow \psi \simeq 68\%, \lambda = 2 \Rightarrow \psi \simeq 95\%, \lambda = 3 \Rightarrow \psi \simeq 99.7\%$$

با فرض این که o مرکز S باشد و توزیع فاصله‌های $s \in S$ حول o تقریباً نیم-گاوسی باشد، می‌توان «فاصله‌ی استاندارد» یک داده در S نسبت به o را مشابه با انحراف معیار محاسبه کرد:

$$\sigma(o, s) = \sqrt{\frac{\sum_{s \in S} d(o, s)^2}{|S|}}$$

فرق این فاصله با انحراف معیار این است که فرض می‌کند میانگین فاصله‌ها نسبت به o صفر است. (و نه برابر با $E[d(o, S)]$). به طور خاص، این تفاوت به این معناست که نمی‌توان توزیع فاصله‌ها را به صورت

نرمال فرض کرد، در عوض تصور می‌کنیم S حول o به صورت نرمال توزیع شده است. همچنین، «فاصله‌ی احتمالاتی مجموعه‌ای» از داده‌ی o به S با میزان اهمیت λ را به صورت زیر تعریف می‌کنیم:

$$pdist(\lambda, o, S) := \lambda \cdot \sigma(o, S)$$

این فاصله‌ی جدید، به ما قدرت کنترل تخمین چگالی S حول o را می‌دهد، اما بر روی ترتیب Outlierها تاثیری نخواهد گذاشت.

معمولاً، هر مدل‌سازی آماری بر اساس برخی مفروضات است. این مورد نیز مبتنی بر فرض‌های زیر است: یک. زیرمجموعه‌ی محلی S به مرکزیت داده‌ی مدنظر o است.

دو. توزیع فاصله‌ها مشابه با قسمت مثبت یک توزیع نرمال است.

در مورد فرض یک، در صورتی که o از C_S (نقطه‌ی مرکزی مجموعه‌ی S) فاصله‌ی زیادی داشته‌باشد، نشان‌دهنده‌ی این است که S نسبت به o نامتقارن است و اندیکاتور است برای Outlier بودن نقطه‌ی o . فرض دو نیز از تقریبی حول قضیه‌ی آماری مقدار میانی استخراج شده و استفاده از این قضیه نشان می‌دهد که ما توزیع واقعی نقاط S را محدود به توزیع خاصی نکرده‌ایم و هر توزیعی که داشته باشند، با توجه به قضیه‌ی مقدار میانی و فرض این که از فاصله‌های مدل L_p (مثل منهتن یا اقلیدسی استفاده می‌کنیم) می‌توان آن‌ها را با توزیع نرمال تقریب زد.

با توجه به این دلایل برای تقریب زدن چگالی یک مجموعه حول یک نقطه، «معیار احتمالاتی محلی Outlier» نسبت به یک داده‌ی $o \in D$ ، یک میزان اهمیت λ و یک مجموعه‌ی محلی $S(o) \subset D$ به صورت زیر تعریف می‌شود:

$$PLOF_{\lambda, S}(o) := \frac{pdist(\lambda, o, S(o))}{E_{a \in S(o)}[pdist(\lambda, s, S(s))]} - 1$$

مقدار PLOF یک داده‌ی o نسبت تقریب چگالی حول o که بر اساس $S(o)$ است به امیدریاضی تقریب چگالی حول تمامی داده‌ها در زیرمجموعه‌ی S را حساب می‌کند. این معیار که بسیار شبیه به معیار LOF است هنوز یک احتمال نیست و نیاز به نرمال‌سازی دارد. یک مقدار $0 \leq$ به معنای Outlier بودن نیست، اما مقادیر بالاتر میزان Outlier بودن بیش‌تری را نشان می‌دهند. برای نرمال‌سازی این مقدار فرای یک دیتاست خاص، مقدار متوسط $nPLOF$ بر اساس محاسبه‌ی $PLOF$ برای چند دیتاست به دست می‌آید:

$$nPLOF := \lambda \cdot \sqrt{E[(PLOF)^2]}$$

این مقدار می‌تواند به عنوان نوعی انحراف معیار نسبت به مقادیر $PLOF$ در نظر گرفته شود، یعنی $\sigma(PLOF) \cdot \lambda$ با فرض \bullet بودن میانگین. ما در نهایت دوباره تابع ارور گاوسی را بر محاسبات اعمال می‌کنیم تا یک مقدار احتمال به نام «احتمال محلی Outlier بودن» را حساب می‌کنیم که مستقیماً نشان دهنده‌ی احتمال Outlier بودن نقطه‌ی $o \in D$ است.

$$LoOP_{s(o)} := \max\{0, \text{erf}\left(\frac{PLOF_{\lambda, s(o)}}{nPLOF \cdot \sqrt{2}}\right)\}$$

این مقدار برای نقاطی که درون مناطق پرچگالی هستند نزدیک به صفر خواهد بود و برای نقاط دور افتاده نزدیک به یک. در نتیجه، بر خلاف مقادیری همانند LOF که بر روی دیتاست‌های مختلف نتایج متفاوتی ارائه خواهد کرد، این مقدار برای تمامی نقاط روی یک دیتاست و بر روی دیتاست‌های مختلف نتایج پایداری ارائه خواهد کرد. در جدول زیر نیز تفاوت این دو معیار برای یک دیتاست ساختگی را با مقادیر $k=20$ و $\lambda=3$ مشاهده می‌کنید:

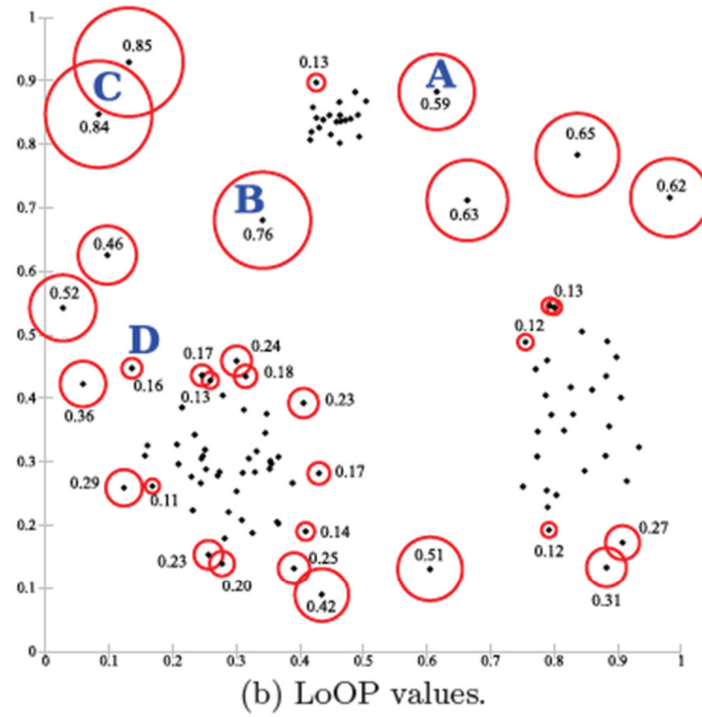
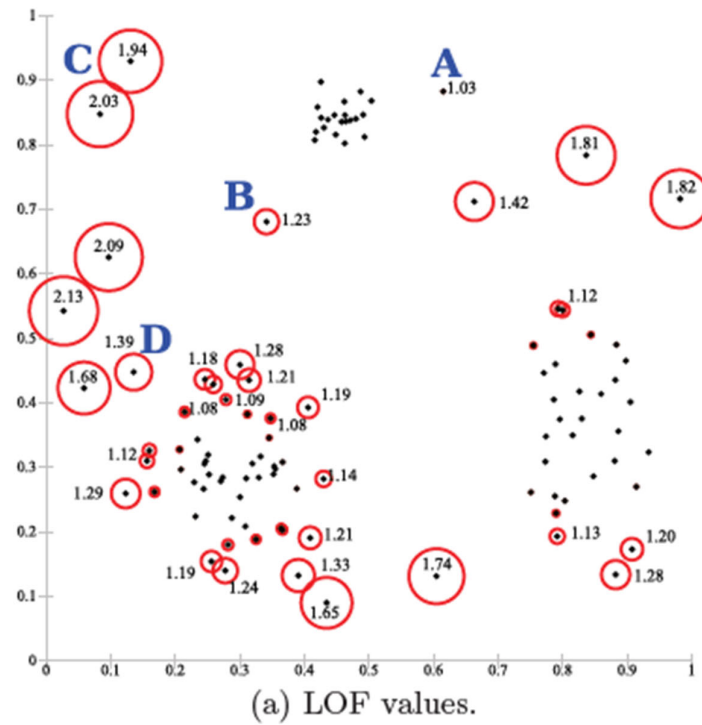


Figure 1: Comparison of the interpretability of local, density-based outlier scoring (here: LOF) values and LoOP values on 2D synthetic data. Both algorithms were run with $k = 20$, for LoOP $\lambda = 3$.