



بررسی داده کاوی توزیع شده با الگوریتم k-means

نجمه تقی زاده^۱، لاله مدادح علی^۲، دکتر محبوبه شمسی^۳، علیرضا آراسته^۴

1. دانشجوی کارشناسی ارشد فناوری اطلاعات دانشگاه قم
2. دانشجوی کارشناسی ارشد فناوری اطلاعات دانشگاه قم
3. عضوهیئت علمی گروه مهندسی برق و کامپیوتر دانشگاه صنعتی قم
4. فارغ التحصیل مهندسی نرم افزار دانشگاه پیام نور مرکز قم

نویسنده مسئول : نجمه تقی زاده

چکیده

اکثر الگوریتم های خوشه بندی نیاز به داده های متتمرکز دارند ، اما این الگوریتم ها با توسعه اینترنت و در برخورد با داده های توزیع شده ، با دو چالش روبرو شدند. اول، حجم داده های تولید شده حتی برای ابر کامپیوتراها هم خیلی زیاد شده است. دوم، داده ها در چندین مکان ذخیره شده اند و متتمرکز کردن آنها در یک جا بسیار پرهزینه خواهد بود، هم چنین محدودیت پهنای باند و حریم شخصی نیز از نگرانی ها و موانع متتمرکز سازی داده ها می باشد. به همین دلیل برای حل این مشکلات، داده کاوی توزیع شده یک حوزه تحقیقاتی پر طرفدار شده است. یکی از الگوریتم های خوشه بندی، الگوریتم کامیز است که به عنوان یکی از ابا نفوذترین الگوریتم های داده کاوی مورد استفاده قرار می گیرد و بسیار ساده و مقیاس پذیر است. در سال های اخیر نسخه هایی از این الگوریتم انتشار یافته است که می تواند در برخورد با داده های توزیع شده ، به خوبی عمل کرده و نتایج خوبی را ارائه دهد. در این الگوریتم ها ، نیازی به جمع آوری کردن اطلاعات و داده ها در یک مجموعه متتمرکز نیست. در این مقاله قصد داریم که این الگوریتم ها را معرفی و بررسی کنیم.

واژه های کلیدی: داده کاوی توزیع شده، خوشه بندی، الگوریتم کامیز، الگوریتم کامیز توزیع شده، نرم افزار

مقدمه

حجم داده های تولید شده در سال های اخیر در حال رشد است. مجموعه اسناد، عکس ها، بیانفورماتیک و دیگر داده ها با تکنولوژی های جدید، در حال افزایش هستند. در بسیاری از موارد، داده ها در سایت های مختلف داده ، به طور طبیعی توزیع ، تولید و ذخیره شده اند. در این زمینه، الگوریتم ها باید قادر به استخراج اطلاعات مناسب از داده های توزیع شده با عملکرد خوب محاسباتی باشند. با این حال، اکثر تکنیک های داده کاوی که پیشنهاد می شوند داده های متتمرکز را مورد توجه قرار می دهند ، اما متتمرکز کردن آنها در یک جا بسیار پرهزینه خواهد بود ، هم چنین محدودیت پهنای باند و حریم شخصی نیز از نگرانی ها و موانع متتمرکز سازی داده ها می باشد.

خوشه بندی یکی از تکنیک های پر کاربرد داده کاوی است که در آن ، مجموعه ای از اشیا که معمولاً به صورت چند بعدی هستند به گروه ها یا کلاس هایی تقسیم بندی می شوند که در هر گروه، اشیا بر اساس معیارهای از پیش تعیین شده ، به هم شبیه هستند. به طور کلی، فاصله اقلیدسی از مراکز خوشه ها به عنوان معیار شbahat یا تفاوت در نظر گرفته می شود. الگوریتم های خوشه بندی در زمینه های مختلفی مثل داده کاوی، تشخیص الگو ، تئوری یادگیری و دیگر زمینه ها کاربرد دارد.

یکی از متدائل ترین الگوریتم های خوشه بندی ، الگوریتم کامیز است . این الگوریتم کامیز به خاطر سادگی، مقیاس پذیری و کاربرد در نرم افزارهای مختلف، به عنوان یکی از 10 الگوریتم موثر و بانفوذ داده کاوی انتخاب شده است. با این حال، الگوریتم کامیز، به انتخاب نمونه خوشه های اولیه حساس است، پس اگر خوشه های اولیه بدست انتخاب نشوند نتایج خوب و بهینه بدست نمی آید. هم چنین در این الگوریتم باید تعداد خوشه های مورد نظر (K) را مشخص کرد و این کار (مشخص کردن تعداد خوشه ها) می تواند کاملاً محدود کننده باشد. در سال های اخیر نسخه هایی از این الگوریتم انتشار یافته است که می تواند در برخورد با داده های توزیع شده به خوبی عمل کند.

در قسمت دوم مقاله، به بررسی داده کاوی توزیع شده (DDM) می پردازیم. در قسمت سوم الگوریتم استاندارد کامیز و نسخه بهبود یافته آن را بررسی می کنیم. در قسمت چهارم، مقاله، الگوریتم های کامیز توزیع شده را مرور می کنیم.

داده کاوی توزیع شده

داده کاوی توزیع شده ، تکنیک کشف الگو ها یا تولید مدل هایی از داده های توزیع شده است که به صورت متتمرکز نه امکان پذیر است و نه مطلوب. داده کاوی توزیع شده می تواند در ابر کامپیوتراهای موازی، شبکه های نظیر و شبکه های حسگر مورد استفاده قرار گیرد. اما، محیط های مختلف مسائل و نگرانی های مختلفی دارند. DDM به دسته های زیر تقسیم می شود.[1]



خوشه های متمرکز در مقابل نظیر به نظر

خوشه متمرکز یک هماهنگ کننده دارد. هماهنگ کننده کار را به کامپیوترهای مختلف تقسیم می کند. خوشه متمرکز به کارگیری و هماهنگی ساده تری دارد. و کار را عادلانه تقسیم می کند و مشکل نقطه شکست(single point of failure) دارد. الگوریتم های نظیر به نظر به یک سرور مرکزی بستگی ندارند، و هر سایت داده را می گیرد و کار خودش را انجام می دهد.

مدل تکی در مقابل فرا-یادگیری

مدل تکی یک الگوریتم داده کاوی داده ها را به بخش های کوچکتر که در هر سایت توزیع شده اند تقسیم می کند. الگوریتم میتواند داده، نتایج میانی، مدل های پیشگوی، نتایج نهایی یک الگوریتم داده کاوی را انتخاب کند و منتقل کند. از طرف دیگر، فرا-یادگیری، به صورت یادگیری از دانش یادگرفته شده تعریف شده است، تکنیک دیگری است که سر و کارش با مشکل محاسبه یک مدل سراسری از پایگاه های داده ایی ذاتا توزیع شده و بزرگ است. هدف فرا-یادگیری محاسبه تعدادی مدل مستقل(classifier) با اعمال برنامه های یادگیری به طور موازی بدون انتقال یا دسترسی مستقیم به سایت است.

داده همگن در مقابل داده ناهمگن

در یک سیستم دیتابیس رابطه ای توزیع شده، اطلاعات باید در سایت های مختلف ذخیره شود. هر سایت جداول رابطه ای کامل دارد، داده همگن برای هر سایت ذخیره می شود. که هر سایت اطلاعات جداول مختلف را ذخیره می کند، در این حالت مسئله داده ناهمگن داریم. اکثر داده کاوی موجود، داده همگن را در سایت های مختلف در نظر می گیرد. در حالت داده توزیع شده ناهمگن، ما فقط داشت ناقص را درباره دیتابیس کامل مشاهده می کنیم. مدل های محلی مختلف دید محلی از کل مسئله دارند داده کاوی توزیع شده باید یک مدل سراسری از آن مدل های محلی مختلف ایجاد کند. پس، کاوش دادگان ناهمگن، چالش برانگیز است.

داده کاوی توزیع شده و داده کاوی موازی

هم داده کاوی توزیع شده و هم داده کاوی موازی فرآیند داده کاوی را سرعت می بخشدند اما از جهاتی متفاوت هستند. DDM اغلب الگوریتم های کاوش متفاوت یا یکسانی را اعمال می کنند برای داده محلی و بین واحدهای پردازش ارتباط برقرار می کند و کشف الگوی محلی را ترکیب می کند با الگوریتم های داده کاوی محلی از پایگاه های داده ای محلی به یک راه حل داشت جهانی. در این حالت، داشت کشف شده اغلب با داشت کشف شده توسعه اعمال الگوریتم های داده کاوی به کل دیتابیس فرق دارد. دقت یا کارایی DDM تقریباً پیشگویی اش سخت است چون بستگی دارد به افزار داده، زمان بندی کار، و سنتز جهانی. در مقابل با DDM، یک الگوریتم داده کاوی موازی همان داشتی را پیدا می کند که با الگوریتم ترتیبی اش یافت شده به علت اعمال الگوریتم موازی جهانی روی کل دیتابیس. دقت آن نسبت به DDM بیشتر ضمانت شده است.

الگوریتم کامینز

الگوریتم کامینز استاندارد

در این الگوریتم، داده ها را به K خوشه مجزا تقسیم کنیم. این الگوریتم به دو فاز اول برای هر خوشه یک نقطه را به عنوان نقطه ثقل خوشه(centroid) یا نقطه مرکزی خوشه بدست می آوریم و در فاز بعدی بدست می آوریم که هر نقطه از مجموعه داده به کدام مرکز خوشه(centriod) نزدیک تر است و آن نقطه را به خوشه مربوطه نسبت می دهیم. در حالت کلی برای بدست آوردن فاصله بین نقاط داده و مراکز خوشه ها، از فاصله اقلیدسی استفاده می شود. زمانی که تمام نقاط در خوشه ها قرار گرفتند مرحله اول به اتمام رسیده و خوشه بندی اولیه انجام شده است. سپس دوباره برای خوشه ها مراکز جدیدی را بدست می آوریم و فاصله هر نقطه را نسبت به این نقاط مرکزی اندازه می گیریم تا خوشه ها به روز شوند و این کار تا زمانی ادامه پیدا می کند که خوشه ها همگرا شوند. [2]

Algorithm 1: The k-means clustering algorithm

Input:

$D = \{d_1, d_2, \dots, d_n\}$ //set of n data items. k // Number of desired clusters

Output:

A set of k clusters.

Steps:

1. Arbitrarily choose k data-items from D as initial centroids;
2. Repeat

 Assign each item d_i to the cluster which has the closest centroid;

 Calculate new mean for each cluster;

 Until convergence criteria is met.

اما اشکال عمده این الگوریتم این است که با توجه به مقدار مراکز اولیه، خوشه های متفاوتی تولید می شود و در نتیجه کیفیت خوشه های نهایی شدیداً به انتخاب مراکز اولیه خوشه ها مابتنی است. این الگوریتم احاظ محسنه را تعداد نقاشه، تعداد خمشه ها و تعداد تکرارها نیاز به زمان دارد. در قسمت بعد الگوریتم اصلاح شده کامینز را بررسی می



الگوریتم کامینز بهبود یافته

همانطور که بیان شد، الگوریتم خوش بندی کامینز استاندارد از لحاظ محاسباتی سنتی و کیفیت نتایج خوش های آن، شدیداً به انتخاب مراکز اولیه خوش ها وابسته است. یه همین دلیل محققان سعی کرده اند تا با ارائه روش هایی این نقایص را برطرف کنند و الگوریتم کامینز را بهبود ببخشند. در روش کامینز بهبود یافته، هر دو فاز الگوریتم کامینز برای بهبود دقت و کارایی آن اصلاح شده است. در مرحله اول، مراکز اولیه خوش ها برای تولید خوش ها با دقت بالاتر از یک روش سیستماتیک به جای انتخاب تصادفی استفاده میکنند. مرحله دوم با تشکیل خوش های اولیه بر اساس فاصله نسبی هر نقطه (داده) از مراکز اولیه خوش ها شروع می شود. این خوش ها متعاقباً بوسیله یک روش اکتشافی تنظیم می شوند، بنابراین کارایی بهبود می یابد.^{[3],[4]} هر دو مرحله الگوریتم، در دو مرحله 2 و 3 توضیح داده شده که مطابق زیر است :

Algorithm 2: Finding the initial centroids

Input:

$D = \{d_1, d_2, \dots, d_n\}$ // set of n data items k // Number of desired clusters

Output:

A set of k initial centroids .

Steps:

1. Set $m = 1$;
2. Compute the distance between each data point and all other data-points in the set D ;
3. Find the closest pair of data points from the set D and form a data-point set Am ($1 \leq m \leq k$) which contains these two data-points, Delete these two data points from the set D ;
4. Find the data point in D that is closest to the datapoint set Am , Add it to Am and delete it from D ;
5. Repeat step 4 until the number of data points in Am reaches $0.75*(n/k)$;
6. If $m < k$, then $m = m+1$, find another pair of datapoints from D between which the distance is the shortest, form another data-point set Am and delete them from D . Go to step 4;
7. For each data-point set Am ($1 \leq m \leq k$) find the arithmetic mean of the vectors of datapoints in Am , these means will be the initial centroids.

Algorithm 3: Assigning data-points to clusters

Input:

$D = \{d_1, d_2, \dots, d_n\}$ // set of n data-points.

$C = \{c_1, c_2, \dots, c_k\}$ // set of k centroids

Output:

A set of k clusters

Steps:

1. Compute the distance of each data-point di ($1 \leq i \leq n$) to all the centroids cj ($1 \leq j \leq k$) as $d(di, cj)$;
 2. For each data-point di , find the closest centroid cj and assign di to cluster j .
 3. Set $ClusterId[i] = j$; // j : Id of the closest cluster
 4. Set $Nearest_Dist[i] = d(di, cj)$;
 5. For each cluster j ($1 \leq j \leq k$), recalculate the centroids;
 6. Repeat
 7. For each data-point di ,
 - 7.1 Compute its distance from the centroid of the present nearest cluster;
 - 7.2 If this distance is less than or equal to the present nearest distance, the data-point stays in the cluster;
 - Else
 - 7.2.1 For every centroid cj ($1 \leq j \leq k$)

Compute the distance $d(di, cj)$;
 - Endfor;
 - 7.2.2 Assign the data-point di to the cluster with the nearest centroid cj
 - 7.2.3 Set $ClusterId[i] = j$;
 - 7.2.4 Set $Nearest_Dist[i] = d(di, cj)$;
 - Endfor;
 8. For each cluster j ($1 \leq j \leq k$), recalculate the centroids;
- Until the convergence criteria is met.

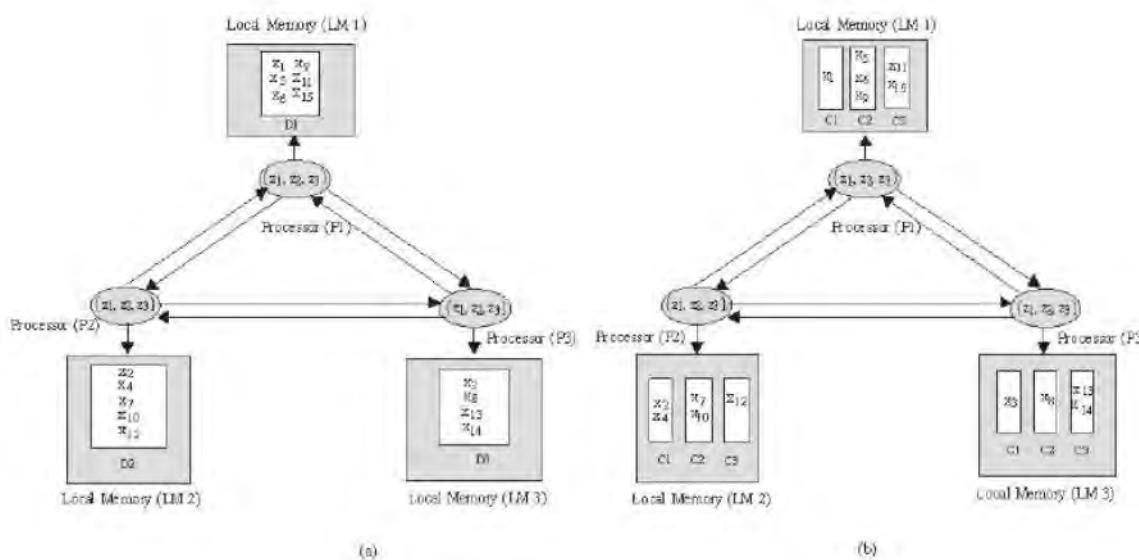


الگوریتم های کامپیوتری توزیع شده

همانطور اشاره شد، با توجه به توزیع شدگی داده ها در مجموعه سایت های مختلف و عدم مرکز کردن آنها در یک مجموعه، نیاز به الگوریتم هایی می باشد که در محیط های توزیع شده به خوبی عمل کنند. الگوریتم کامپیوتری کامپیوز کی از الگوریتم هایی است که می تواند در این موارد نتایج مورد قبولی را ارائه دهد، این الگوریتم به صورت استاندارد قادر خواهد بود در یک مجموعه مرکز کار کند، لکن با کمی تغییر در ساختار این الگوریتم و ترکیب آن با دیگر الگوریتم ها می توان آن را به صورت توزیع شده استفاده کرد. در زیر به ذکر نمونه هایی از این الگوریتم که در محیط های توزیع شده کار می کنند، می پردازم.

الگوریتم کامپیوز توزیع شده در محیط های موازی

اگر به الگوریتم کامپیوز استاندارد نگاه کنیم میبینیم که توزیع شدگی به طور ذاتی درون آن نهفته است. فرایند با یک مجموعه داده آغاز می شود و سپس مراکز اولیه خوش ها انتخاب می گردند. برای اجرای الگوریتم کامپیوز استاندارد به یک سیستم پردازنده واحد احتیاج ندارد. برای هر داده در هر تکرار اعضیت آن در خوش ها محاسبه می شود. در این وضعیت، پردازنده همه ویژگی های ساختاری الگوریتم کامپیوز را در حافظه محلی نگهداری می کند و تکرار گام های الگوریتم کامپیوز نقاطه مرکزی نهایی Z را محاسبه می کند. در محیط های توزیع شده ما از پردازنده های شبکه ای برای الگوریتم کامپیوز استفاده میکنیم. فرض میکنیم که مجموعه داده ها در پردازنده های شبکه توزیع شده اند. این پردازنده ها، فرایند را در یک روش تعاملی اجرا می کنند. شکل (a) و (b) توزیع داده را در یک شما توزیع شده در قبل و بعد از اجرای الگوریتم نوزیع شده نشان می دهد.^[5]



شکل (1) شما از اجرای الگوریتم k-means در محیط های توزیع شده (a) قبل از خوش بندی (b) بعد از خوش بندی

فرض میکنیم که داده های اولیه توزیع شده، مطلقاً دلخواه و مستقل هستند. هر پردازنده P_i به طور تصادفی، یک مجموعه برداری از مراکز $Z_i = \{Z_k \text{ for } k=1 \text{ to } K\}$ را مقدار دهنده می کند. این مراکز، مراکز قدیمی نام گذاری می شوند $Z_k^{(old)}$. بعد از این مقدار دهنده مراکز خوش های محلی از هر پردازنده i به دیگر پردازنده ها پخش می شود. در داخل حلقه تکرار اصلی الگوریتم کامپیوز توزیع شده، داده ها در هر پردازنده برای محاسبه فاصله از مراکز خوش های مربوطه به صورت موادی استفاده می شود. بعد از توزیع محلی، بردار مراکز جدید $Z_k^{(new)}$ با استفاده از داده ها محلی طبقه بندی شده و مراکز قدیمی، محاسبه می شود. این کار برای همه K ها انجام میگیرد. هم چنین دقت داشته باشید که از مراکز قدیمی خوش های نیز برای محاسبه مراکز جدید استفاده می شود بخاطر اینکه از ایجاد خوش های خالی در هر پردازنده جلوگیری شود.

گام محاسبه مرکز، یکی از مهمترین بخش های این الگوریتم است. از نظر ریاضیاتی، تفاوت بین محاسبه مراکز در الگوریتم کامپیوز معمولی و کامپیوز توزیع شده به صورت زیر است :

محاسبه مرکز در الگوریتم کامپیوز استاندارد :

$$Z_k^{(new)} = \frac{1}{n_k} \{ \sum_{x_j \in c_j} (X_j) \} \quad (1)$$

محاسبه مرکز در الگوریتم کامپیوز توزیعی :

$$Z_k^{(new)} = \frac{1}{n_k+1} \{ \sum_{x_j \in c_j} (X_j) + Z_k^{(old)} \} \quad (2)$$

بردار مراکز جدید در بین همه پردازنده ها، توزیع شده است. برای هر پردازنده i ، بردار مرکز جدید بصورت میانگین نقاط و مراکز قدیمی برای همه مقادیر $k=1$ تا K بدست می آید. می کنند.

همایش ملی الکترونیکی دستاوردهای نوین در علوم مهندسی و پایه

National e-Conference on Advances in Basic Sciences and Engineering

WWW.AEBSCONF.IR



فرایند بالا ادامه می‌باید تا زمانی که دیگر بردار مرکز، تغییری نکند و ثابت بماند. استراتژی محاسبه مراکز جدید باعث می‌شود این الگوریتم از دیگر الگوریتم‌های کامپیوتر متغیر باشد.

غیر اینصورت همان الگوریتم کامپیوتر اصلی می‌باشد.

این الگوریتم بصورت زیر است:

Algorithm 4: The Distributed k -means Algorithm

Input: a set S of d-dimensional data and an integer K.

Output :K clusters

Begin

for each processor P_i do in parallel initialize a set of center vectors z_k for $k=1$ to K ,
(these centers are called old centers $\{Z_k^{(old)}\}$)

Repeat

for each processor P_i do in parallel begin distribute local data in P_i into k classes according to
minimum distance from z_k for $k=1$ to K compute new center vectors $\{Z_k^{(new)}\}$ for $k=1, \dots, K$
considering old centers $\{Z_k^{(old)}\}$ as data items

end

for each processor P_i broadcast $\{Z_k^{(new)}\}$, for $k=1, \dots, K$, to all other processors

for each processor P_i do in parallel begin

for $k=1, \dots, K$ do

compute average of $\{Z_k^{(new)}\}$ from self and those received from other processors and replace $\{Z_k^{(old)}\}$ with this average

end

until center vectors are stable

end

الگوریتم خوشبندی نرم توزیع شده گروه محور (DSCA)

الگوریتم کامپیوتری خوشبندی توزیع شده (DKM) گروه محور، یکی دیگر از الگوریتم‌های خوشبندی توزیع شده است. این الگوریتم به منظور خوشبندی کردن داده‌های محلی استفاده می‌شود. نقاط مرکزی خوشبندی محلی، در یک مکان مرکزی ترکیب می‌شوند و با استفاده از الگوریتم کامپیوتری خوشبندی می‌شوند که نقاط مرکزی سراسری را تولید می‌کنند. از این نقاط سراسری برای به روز کردن خوشبندی محلی استفاده می‌شود. به منظور تولید خوشبندی نرم سراسری، الگوریتم کامپیوتری خوشبندی اصلاح شده است. [6]

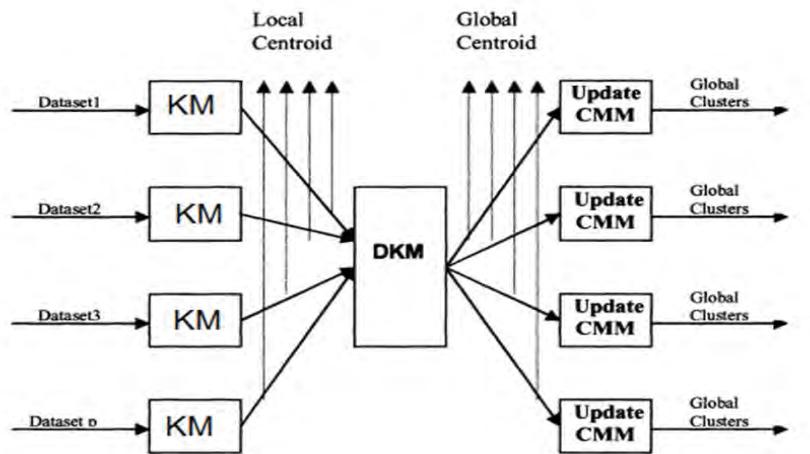
خوشبندی نرم و خوشبندی سخت

به طور طبیعی خوشبندی نرم یا سخت باشند. در خوشبندی های معمول، اشیاء مشابه در یک خوشبندی جمع آوری می‌شوند در حالی که اشیاء با ویژگی‌های مختلف در خوشبندی دیگر هستند. به این خوشبندی که نقطه مشترک ندارند خوشبندی های سخت می‌گویند. اما در خوشبندی نرم اشیاء می‌توانند در بیش از یک خوشبندی عضو باشند. خوشبندی نرم ممکن است بصورت فاری باشد و یا مزه‌های ناهماور داشته باشد.

روش معمول کامپیوتری، که در اکثر برنامه‌های داده کاوی پیدا می‌شود، همیشه منجر به خوشبندی های سخت می‌شود. اما در عمل، ویژگی‌هایی که یک شی ممکن است در دسته‌های مختلف به نمایش در بیاید. در چنین مواردی شی باید به چندین خوشبندی متعلق باشدو در نتیجه لزوماً مزه‌های خوشبندی دارند. الگوریتم های خوشبندی فازی، مانند الگوریتم C-Means این احتمال را ایجاد می‌کنند که یک شی به چندین خوشبندی متعلق داشته باشد. [7]

تشریح الگوریتم

الگوریتم خوشبندی نرم توزیع شده گروه محور (DSCA) در شکل 2 نمایش داده شده است. در ابتدا مجموعه داده‌های محلی، با استفاده از الگوریتم k-means ماتریسی از مقادیر عضویت خوشبندی (CMM) و نقاط مرکزی محلی بدست آید. در این مرحله، فقط نقاط مرکزی محلی به یک مکان مرکزی منتقل می‌شوند و در آنجا با هم ترکیب شده و به اصطلاح CentroidDataset را نامیده می‌شوند. در این زمان، الگوریتم کامپیوتری (DKM) به کار می‌رود تا نقاط مرکزی محلی (local centroid) را که مشابه هستند گروه بندی کرده و سپس میانگین آنها، نتایج مورد نیاز برای فرایند خوشبندی توزیعی نقاط سراسری (global Centroid) است. در گام آخر نیز، ماتریس CMM که در اولین مرحله تولید شد را توسط این نقاط سراسری (Global Centroid) به روز می‌کنیم تا در نهایت خوشبندی نهایی بدست آید. بدون در نظر گرفتن تعداد خوشبندی های تولید شده در مجموعه داده های محلی، تعداد K گروه سراسری بدست می‌آید.



شکل 2 چارجوب برای خوشبندی نرم توزیع شده

الگوریتم k-means توزیع شده نرمال نرمال سازی

پیش پردازش اغلب قبل از استفاده از هر الگوریتم داده کاوی برای بهبود کارایی نتیجه نیاز است. نرمال سازی داده یکی از روش‌های پیش پردازشی در داده کاوی است، که صفت درون یک محدوده خاص قرار می‌گیرد مثل $-1.0 \leq A \leq 1.0$ یا $0.0 \leq A \leq 1.0$. نرمال سازی قبل از خوشبندی برای معیار فاصله استفاده می‌شود، مثل فاصله اقلیدسی، که به اختلاف در اندازه یا مقیاس صفات حساس است. در اپلیکیشن‌های واقعی، به دلیل تفاوت در محدوده مقدار صفات، یک صفت باید به دیگری غلبه کند. نرمال سازی مانع از زیادی وزن دادن به صفات با محدوده بزرگ مثل حقوق نسبت به صفات با محدوده کوچکتر می‌شود مثل سن، هدف، مساوی سازی اندازه یا بزرگی و تنوع این صفات است.

روش‌های مختلفی برای نرمال سازی داده وجود دارد، نرمال سازی Z-score، Min-max و مقیاس ده دهی (decimal scaling). نرمال سازی Min-max یک تغییر خطی روی داده اصلی ایجاد می‌کند. فرض کنید A مقادیر کمینه و بیشینه برای صفت A هستند. نرمال سازی Min-max یک مقدار $\bar{A} = \frac{A - min_A}{max_A - min_A}$ را در محدوده $(0, 1)$ اینگونه نگاشت می‌کند:

$$v' = \frac{v - min_a}{(max_a - min_a)} \quad (3)$$

در نرمال سازی Z-score، مقادیر برای یک صفت A بر اساس میانگین و انحراف معیار A نرمال می‌شود. مقدار \bar{A} به v نرمال می‌شود اینگونه:

$$v' = \frac{v - \bar{A}}{\sigma_A} \quad (4)$$

که \bar{A} و σ_A میانگین و انحراف معیار صفت A هستند. این روش نرمال سازی مفید است وقتی که مقدار واقعی کمینه و بیشینه صفت A مجهول هستند، یا وقتی که پرتهایی هستند که بر نرمال سازی \bar{A} غلبه می‌کند. نرمال سازی مقیاس ده دهی مقادیر صفت A را منتقل می‌کند. تعداد نقاط دهی منتقل شده به مقدار محض بیشینه A بستگی دارد. مقدار \bar{A} به v نرمال می‌شود اینگونه:

$$v' = \frac{v}{10^j} \quad (5)$$

که j کوچکترین عدد صحیحی است که $1 < |v'| < Max(|v'|)$.

الگوریتم توزیع شده نرمال کامینز (NDKM) ابتدا در سایت‌های محلی با الگوریتم k-means خوشبندی را انجام می‌دهد، سپس تمام مقادیر مرکزی سراسری، مجدداً از k-means بدست می‌آید. NDKM سخه توسعه یافته DKM هست، که برای استاندارد سازی اشیا داده ای به یک محدوده انجام می‌شود. روش گام به گام الگوریتم پیشنهادی در زیر آمده است. ابتدا مقادیر بیشینه و کمینه بردارهای ویژگی از تمام دیتابسته‌های محلی استخراج می‌شود و به یک مکان مرکزی ارسال می‌شود، که مقادیر بیشینه و کمینه مشخص شده‌اند. این دو مقدار برای انجام نرمال سازی سراسری با روش نرمال سازی min-max به سایت‌های محلی منتقل می‌شوند. سپس، اشیا نرمال شده برای ایجاد ماتریس مرکز و اندیس خوشبندی برای هر دیتابسته با k-means کامل شوند. تمام مرکزهای محلی با الگوریتم k-means، برای گروه کردن مرکزهای مشابه و کسب مرکزهای سراسری دسته بندی و ادغام می‌شوند. مرکزهای سراسری به سایت محلی منتقل می‌شود، که فاصله اقلیدسی هر شی از مجموعه سراسری مرکزها محاسبه می‌شوند و به مرکز نزدیک‌ترین خوشبندی منتب می‌شود. الگوریتم زیر با تغییر دو گام اول، می‌تواند برای انواع دیگر روش نرمال سازی پیاده سازی شود.

Algorithm 5: NDKM

Input: Homogeneous p datasets, each with d dimensions

Output: Global partitions of p datasets



Procedure

- Step 1: Find Maximum And Minimum Values of each feature from each local dataset and Transmit them into Central Place
- Step 2: Compute Maximum And Minimum Values at Central Place
- Step 3: Normalize real Scalar Values of Local datasets with Global Maximum And Minimum Values
- Step 4: Cluster each Local dataset By K-Means Algorithm and obtain Centroid Matrix Along with cluster index for each dataset
- Step 5: Merge cluster centroids of local datasets into a single dataset named as centroids Dataset at Central Place
- Step 6: Cluster centroids dataset using K-Means Algorithm to obtain Global Centroids
- Step 7: Update local Cluster indices by assigning each object to nearest cluster centroid, After Computing Euclidean distance between the object and global centroid

نتیجه گیری

با توجه به رشد داده های توزیع شده و عدم جایه جایی این داده ها به یک مجموعه مرکزی ، داده کاوی توزیع شده می تواند بسیار مفید باشد. در این میان الگوریتم های متفاوتی برای این زمینه وجود دارد. در این مقاله روش های داده کاوی توزیع شده با الگوریتم کامپینز را بررسی کردیم. برخی از این الگوریتم ها به صورت موادی اجرا می گردند می توانند سرعت داده کاوی را به طور قابل ملاحظه ای افزایش دهند. هم چنین بنابر اهمیت امنیت داده در محیط های توزیع شده، الگوریتم های گروه محور می توانند بدون نیاز به جایه جا کردن اطلاعات در این محیط ها، داده کاوی را انجام دهند.

منابع

- [1] Rekha Sunny T, Survey on Distributed Data Mining in P2P Networks
- [2]K. A. Abdul Nazeer, Improving the Accuracy and Efficiency of the k-means Clustering Algorithm, 2009 Vol I
- [3] Fahim A.M, Salem A. M, Torkey A and Ramadan M. A, "An Efficient enhanced k-means clustering algorithm," Journal of Zhejiang University, 10(7):1626–1633, 2006
- [4] Yuan F, Meng Z. H, Zhang H. X and Dong C. R, "A New Algorithm to Get the Initial Centroids," Proc. of the 3rd International Conference on Machine Learning and Cybernetics, pages 26–29, August 2004.
- [5] Malay K. Pakhira, Clustering Large Databases in DistributedEnvironment,2009
- [6] N. KarthikeyaniVisalakshi, ENSEMBLE BASED DISTRIBUTED SOFTCLUSTERING,IEEE,2008
- [7] Zhou A, Cao F, Van Y, Sha C, He X, "Distributed Data StreamClustering: A Fast EM-based Approach", ICDE 2007, IEEE 23rdInternational Conference on Data Engineering, pp. 736-745 ,April 2007
- [8]SouptikDatta, Chris R. Approximate Distributed K-Means Clustering over a Peer-to-Peer Network,2009
- [9] M.C. Naldi, Evolutionary k-means for distributed datasets,30-42,2014
- [10] JiGenlin, Ling Xiaohan, Ensemble learning based distributed clustering, Lecture Notes In Computer Science, Springer-Verlag, Vol. 4819, pp. 312-321, November 2007.
- [11] Amir Ben-Dor, Ron Shamir and Zohar Yakini, Clustering Gene Expression Patterns, Journal of Computational Biology, 6(3/4): 281-297, 1999
- [12] Margaret H. Dunham, Data Mining- Introductory and Advanced Concepts, Pearson Education, 2006.