



بررسی الگوریتم‌های موازی و توزیع شده برای استخراج مجموعه آیتم پرتکرار

سیاوش فخری^۱، حمید رستگاری^{۲*}، مهدی شریفی^۳

۱- دانشکده مهندسی کامپیوتر، واحد نجف آباد، دانشگاه آزاد اسلامی، نجف آباد، ایران.

Siavash.Fakhri@Sco.iaun.ac.ir

۲- دانشکده مهندسی کامپیوتر، واحد نجف آباد، دانشگاه آزاد اسلامی، نجف آباد، ایران.

rastegari@iaun.ac.ir

۳- دانشکده مهندسی کامپیوتر، واحد نجف آباد، دانشگاه آزاد اسلامی، نجف آباد، ایران.

m.sharifi@pco.iaun.ac.ir

چکیده

در جهان امروز، حجم زیادی از داده‌ها به طور مداوم توسط بسیاری از برنامه‌های کاربردی علمی مانند بیوانفورماتیک یا شبکه‌ها ایجاد می‌شوند. برای استخراج ارزش از این مجموعه داده‌های پیچیده، الگوریتم‌های استخراج داده‌ی مختلفی برای کشف ارتباطات پنهان و غیر بدیهی بین داده‌ها استفاده شده است. استخراج مجموعه آیتم‌های پرتکرار یکی از مهم‌ترین تکنیک‌های استخراج دانش از داده‌ها در بسیاری از برنامه‌های دنیای واقعی است. الگوریتم‌های مختلفی به طور گسترده برای استخراج مجموعه آیتم‌های پرتکرار از مجموعه‌های داده استفاده شده است. با این حال، فرآیند استخراج مجموعه آیتم‌های پرتکرار با چالش‌های فضای حافظه محدود و قدرت محاسبات روبه‌رو است. اعمال الگوریتم‌های استخراج مجموعه آیتم‌های پرتکرار سنتی برای کار بر روی داده‌های بزرگ، بسیار وقت‌گیر است. محاسبات موازی و توزیع شده استراتژی موثر و اغلب استفاده شده برای سرعت بخشیدن به الگوریتم‌ها برای مجموعه داده‌های بزرگ است. در این مقاله، به مرور کارهای انجام شده در این زمینه پرداخته شده است.

کلیدواژه‌ها: الگوریتم‌های توزیع شده، الگوریتم‌های موازی، الگوهای پرتکرار، قوانین انجمنی، اسپارک، هدوپ



1- مقدمه

کشف قوانین انجمنی، یک روش بسیار محبوب و موثر برای استخراج اطلاعات معنی دار از مجموعه داده های بزرگ می باشد. این کار تلاش می کند که ارتباطات احتمالی میان آیتم های موجود در مجموعه داده های بزرگ مبتنی بر تراکنش را پیدا کند. برای ایجاد این انجمن ها، مجموعه آیتم های پرتکرار باید تولید شوند.

بسیاری از روش ها برای استخراج مجموعه آیتم های پرتکرار در پایگاه داده ها بر روی یک ماشین ارائه شده اند. این روش ها در مجموعه داده های کوچک به خوبی کار می کنند. با این حال، مجموعه داده ها بزرگ تر و بزرگ تر شده اند و به دلیل ظرفیت محدود حافظه و توانایی محاسبه ماشین های تکی، این روش ها برای استخراج مجموعه آیتم های پرتکرار در مجموعه داده های عظیم ناکارآمد هستند. ظرفیت های حافظه، برای مدیریت کردن مجموعه ای کامل از آیتم های کاندید مورد نظر، به سرعت منفجر می شوند و هزینه های محاسباتی می تواند در یک ماشین، گران قیمت باشد.

به طور کلی، استفاده از تکنیک های داده کاوی برای مجموعه داده های عظیم نیاز به مقدار زیادی از منابع دارد. به همین دلیل ما شاهد استفاده ی فراوان از روش های موازی و توزیع شده هستیم که معمولا بر اساس چارچوب های توزیع شده مانند هادوپ¹ و اسپارک² هستیم؛ این دو بر پایه ی چارچوب برنامه نویسی نگاشت-کاهش³، که یک مدل برنامه نویسی ساده برای حل مسائل محاسباتی در مقیاس وسیع است می باشند. در واقع چارچوب برنامه نویسی نگاشت-کاهش مجموعه ای از توابع کتابخانه را در دل خود دارد که جزئیات و پیچیدگی را از دید برنامه نویس پنهان می کند.

2- چالش ها

تکنیک های داده کاوی و الگوریتم های برای پردازش مقدار زیادی از اطلاعات به منظور استخراج اطلاعات مفید و جالب در بسیاری از زمینه های مختلف محبوب شده اند. الگوریتم هایی لازم است که داده ها را به طور خودکار و در روش های کارآمد مورد استفاده قرار دهند. با این وجود، حتی اگر عملکرد سیستم های کامپیوتری پی در پی بهبود یابد، آن ها مناسب نیستند تا با افزایش تقاضا برای برنامه های داده کاوی و اندازه

¹ Hadoop² Spark³ MapReduce



کنفرانس ملی فناوری های نوین در مهندسی برق و کامپیوتر

۲۷ دی ۱۳۹۶



داده‌ها هماهنگ شوند. علاوه بر این، حافظه اصلی سیستم‌ها ممکن است به اندازه کافی برای نگهداری تمام اطلاعات مربوط به برنامه‌های فعلی کافی نباشد.

1-2- چالش‌های الگوریتم‌های موازی

الگوریتم‌های موازی با بهره‌گیری از حافظه اصلی و قدرت پردازش پردازنده و شتاب دهنده‌های موجود در رایانه‌های موازی، می‌توانند به راحتی به مسائل مربوط به زمان اجرا و حافظه مورد نیاز، رسیدگی کنند. با این حال، موازی سازی الگوریتم‌های موجود برای به دست آوردن عملکرد خوب و مقیاس پذیری با توجه به مجموعه داده‌های عظیم، بی‌اهمیت است. در حقیقت، سازماندهی داده‌ها و استراتژی تجزیه و تحلیل به منظور متعادل سازی حجم بار کار برای به حداقل رساندن وابستگی به داده‌ها اهمیت اساسی دارد. نگرانی دیگر مربوط به به حداقل رساندن هماهنگ سازی و سربار ارتباطات است. هزینه‌های I/O نیز باید به حداقل برسد. و در نهایت، و مصرف انرژی نیز باید بهینه باشد. برای ایجاد الگوریتم‌های موازی پیشرفته برای برنامه‌های کاربردی داده کاوی با کارایی بالا، نیاز به رسیدگی به چندین مشکل محاسباتی کلیدی است که می‌تواند به راه حل‌های جدید و بینش‌های جدید در برنامه‌های بین رشته‌ای منجر شود.

2-2- چالش‌های الگوریتم‌های توزیع شده

به طور فزاینده داده‌ها در مکان‌های توزیع شده جغرافیایی مختلف پخش می‌شوند. پردازش متمرکز این داده‌ها بسیار ناکارآمد و گران است. در برخی موارد ممکن است حتی غیر عملی باشد و به خطرات امنیتی منجر شود. بنابراین، در پردازش داده‌ها، به حداقل رساندن میزان داده‌هایی که رد و بدل می‌شوند و در عین حال تضمین صحت و کارایی یک چالش بسیار مهم است. استخراج داده‌های توزیع شده، تجزیه و تحلیل داده‌ها و استخراج را باید به شیوه‌ای گسترده انجام دهد و به محدودیت منابع، بخصوص محدودیت‌های پهنای باند، نگرانی‌های حفظ حریم خصوصی و قدرت محاسبات توجه کند.



3- بررسی الگوریتم های موازی

یکی از این الگوریتم‌ها، الگوریتم $PHIKS^4$ [1] مبتنی بر هدوپ است که مجموعه‌ای از بهینه سازی‌های قابل توجه را برای محاسبه‌ی انترویی مشترک از $Miki^5$ با اندازه‌های مختلف ارائه می‌دهد. استخراج کردن $Miki$ یک مشکل کلیدی در تجزیه و تحلیل داده‌ها با تاثیر بالقوه زیاد بر روی انواع مختلف وظایف، مانند یادگیری تحت نظارت، یادگیری بدون نظارت و یا بازیابی اطلاعات است. به کمک این الگوریتم به طور چشمگیری زمان اجرا، هزینه ارتباطات و مصرف انرژی را در یک پلتفرم محاسباتی کاهش می‌دهد. الگوریتم $PaMPa-HD^6$ [2]، یک الگوریتم استخراج مجموعه آیت‌م بسته‌ی پرتکرار موازی مبتنی بر هدوپ برای مجموعه داده‌هایی با ابعاد بزرگ است. استخراج مجموعه آیت‌م بسته‌ی پرتکرار یکی از پیچیده‌ترین روش‌های کاوش در داده کاوی است. اغلب برای کشف آیت‌م‌هایی که با هم اتفاق می‌افتند با توجه به یک آستانه فراوانی ارائه شده توسط کاربر، به نام حداقل پشتیبانی مورد استفاده قرار می‌گیرد. کارایی این روش در زمان اجرا، متعادل سازی بار و استحکام حافظه است. $FiDooP-DP^7$ [3] با استفاده از مدل برنامه‌نویسی هدوپ یک رویکرد پارتیشن‌بندی برای حل مشکل ارتباطات بالا و سربار استخراج ناشی از تراکنش‌های افزونه ارسال شده در میان گره‌های محاسباتی ارائه می‌دهد. این الگوریتم، از همبستگی بین تراکنش‌ها استفاده می‌کند و با تطبیق معیار شباهت و تکنیک Hash کردن حساس به مکان، مکان‌های بسیار مشابه را به یک پارتیشن داده برای بهبود محل، بدون ایجاد تعداد بیش از حد تراکنش‌های افزونه می‌دهد.

$R\text{-Apriori}^8$ [4] یک الگوریتم Apriori موازی بر اساس چارچوب اسپارک است که تولید مجموعه کاندید را برای رسیدن به سرعت بالا با حذف کلی تولید کاندید سرعت می‌بخشد. در این الگوریتم عملکرد تکرار با حذف مرحله تولید مجموعه کاندید و با استفاده از یک فیلتر bloom به جای یک درخت Hash بهبود داده می‌شود. آیت‌م‌های پرتکرار یکتا در فیلتر bloom ذخیره می‌شوند چرا که فیلترهای bloom سریع‌تر از درخت‌های Hash هستند و می‌توانند به راحتی آیت‌م‌های با طول یک را ذخیره کنند.

⁴ Parallel Highly Informative K-ItemSet

⁵ maximally informative k-itemset

⁶ Parallel MapReduce-Based Frequent Pattern Miner for High-Dimensional Data

⁷ Data partitioning in frequent itemset mining on HaDooP clusters

⁸ Redu*ced-Apriori



⁹MR-Apriori [5] برای مشکل کمبود قدرت محاسبات مبتنی بر هدوپ و بر اساس الگوریتم Apriori است. این الگوریتم کارایی خود را در مجموعه داده‌های عظیم نشان می‌دهد.

¹⁰YAFIM [6]، یک الگوریتم دیگر Apriori موازی بر اساس چارچوب اسپارک است که در مقایسه با الگوریتم Apriori که با هدوپ پیاده سازی شده است، 18 برابر سریع تر است. ¹¹S-FPG [7] یک الگوریتم FP-Growth مبتنی بر اسپارک است که می‌تواند به خوبی و به طور قابل توجهی داده‌های بزرگ را پردازش کند. کارایی و مقیاس پذیری از ویژگی‌های این الگوریتم است.

4- بررسی الگوریتم های توزیع شده

¹²DFIMA [8] یک الگوریتم استخراج مجموعه آیتم پرتکرار توزیع شده مبتنی بر اسپارک است که با استفاده از رویکرد هرس مبتنی بر ماتریس، تعداد مجموعه آیتم‌های کاندید را کاهش می‌دهد. این الگوریتم سرعت و مقیاس پذیری بالایی را ارائه می‌دهد. ¹³FDMCN [9] یک الگوریتم سریع و توزیع شده مبتنی بر CARM برای استخراج الگوهای پرتکرار در شبکه‌های شلوغ می‌باشد. بنابراین انتقال داده‌های رد و بدل شده را کاهش می‌دهد و به سرعت اجرا و مقیاس پذیری در نظر می‌گیرد. ¹⁴DSCAN-MAX [10] برای استخراج توزیع شده مجموعه آیتم‌های پرتکرار حداکثری است. این الگوریتم از ویژگی‌های Scan-Tree استفاده کامل می‌کند که این کار به طور قابل توجهی ارتباطات مورد نیاز بین گره‌ها را کاهش می‌دهد، و بدین گونه بهره‌وری الگوریتم را افزایش می‌دهد.

¹⁵HFIM [11] الگوریتمی مبتنی بر اسپارک است که از ترتیب عمودی مجموعه داده‌ها برای حل مشکل اسکن‌ها در هر تکرار استفاده و فضای جستجو را با تولید مجموعه آیتم‌های کاندید از تراکنش‌های جداگانه کاهش می‌دهد. نتایج تجربی نشان می‌دهد که HFIM از لحاظ زمان اجرا و مصرف حافظه، عملکرد قابل قبولی دارد. ¹⁶FAMR [12] یک الگوریتم Apriori استخراج الگوهای پرتکرار مبتنی بر هدوپ است که از

⁹ MapReduce Apriori

¹⁰ Yet Another Frequent Itemset Mining

¹¹ FP-Growth Algorithm under Apache Spark

¹² Distributed Frequent Itemset Mining Algorithm

¹³ fast and distributed algorithm for mining frequent patterns in congested networks

¹⁴ Distributed Global Maximal Frequent Itemsets based on Sorted SCan-Tree

¹⁵ Hybrid Frequent Item Mining

¹⁶ Frequent Patterns Mining Algorithm Based on MapReduce



تولید کاندیدهای بی شمار پرهیز می کند و باعث استفاده ی مناسب از حافظه و ذخیره ی مقدار زیادی زمان می شود.

Dist-Eclat و BigFIM دو الگوریتمی هستند که در [13] معرفی شدند. روش Dist-Eclat، محض است که فضای جستجو را به صورت یکنواخت در بین ماشین ها توزیع می کند. این الگوریتم قادر به استخراج داده های بزرگ است، اما در هنگام برخورد با مقادیر عظیم داده ها به مشکل می خورد. الگوریتم BigFIM ابتدا از Apriori برای استخراج مجموعه های پرتکرار به طول k استفاده می کند و بعد از Eclat استفاده می کند.

5- نتیجه گیری

در این مقاله، ما یک مرور کلی بر الگوریتم های موازی و توزیع شده و چالش های موجود در این الگوریتم ها انجام دادیم. با توجه به رشد عظیم داده های خام و بزرگ شدن پایگاه داده ها، تغییر رویکردها از الگوریتم های سنتی و متمرکز به الگوریتم های موازی و توزیع شده، بیشتر و بیشتر احساس می شود. استفاده از این رویکردها، علاوه بر برطرف کردن بسیاری از مشکلات و چالش های سیستم های متمرکز، خود نیز با چالش هایی روبه رو هستند که باید به آنها توجه ویژه ای شود. در صورت بی توجهی مشکلات چندین برابر می شود و این رویکردها کارایی خود را از دست می دهند.

6- مراجع

- [1] S. Saber, A. Reza, and M. Florent, "A highly scalable parallel algorithm for maximally informative k-itemset mining," *Knowledge and Information Systems*, journal article vol. 50, no. 1, pp. 1-26, January 01 2017.
- [2] A. Daniele, B. Elena, C. Tania, G. Paolo, P. Fabio, and M. Pietro, "A Parallel MapReduce Algorithm to Efficiently Support Itemset Mining on High Dimensional Data," *Big Data Research*, vol. 10, no. Supplement C, pp. 53-69, 2017/12/01/ 2017.
- [3] Y. Xun, J. Zhang, X. Qin, and X. Zhao, "FiDooop-DP: Data Partitioning in Frequent Itemset Mining on Hadoop Clusters," *IEEE Transactions on Parallel and Distributed Systems*, vol. 28, no. 1, pp. 101-114, 2017.
- [4] S. Rathee, M. Kaul, and A. Kashyap, "R-Apriori: An Efficient Apriori based Algorithm on Spark," presented at the Proceedings of the 8th Workshop on Ph.D. Workshop in Information and Knowledge Management, Melbourne, Australia, 2015.
- [5] X. Lin, "MR-Apriori: Association Rules algorithm based on MapReduce," in *2014 IEEE 5th International Conference on Software Engineering and Service Science*, 2014, pp. 141-144.



- [6] H. Qiu, R. Gu, C. Yuan, and Y. Huang, "YAFIM: A Parallel Frequent Itemset Mining Algorithm with Spark," in *2014 IEEE International Parallel & Distributed Processing Symposium Workshops*, 2014, pp. 1664-1671.
- [7] A. D. d. Gassama, F. Camara, and S. Ndiaye, "S-FPG: A parallel version of FP-Growth algorithm under Apache Spark," in *2017 IEEE 2nd International Conference on Cloud Computing and Big Data Analysis (ICCCBDA)*, 2017, pp. 98-101.
- [8] Z. Feng, L. Min, G. Feng, S. Weiming, S. Abdallah, and M. Yunlong, "A distributed frequent itemset mining algorithm using Spark for Big Data analytics," *Cluster Computing*, journal article vol. 18, no. 4, pp. 1493-1501, December 01 2015.
- [9] L. K. W., C. Sheng-Hao, and L. Chun-Cheng, "A fast and distributed algorithm for mining frequent patterns in congested networks," *Computing*, journal article vol. 98, no. 3, pp. 235-256, March 01 2016.
- [10] Y. Huang, J. Wang, Y. Li, and Q. Lin, "A mining algorithm for distributed global maximal frequent itemsets based on Sorted SCan-Tree," in *2016 IEEE Advanced Information Management, Communicates, Electronic and Automation Control Conference (IMCEC)*, 2016, pp. 1818-1824.
- [11] S. K. Kumar and R. Dharavath, "HFIM: a Spark-based hybrid frequent itemset mining algorithm for big data processing," *The Journal of Supercomputing*, journal article vol. 73, no. 8, pp. 3652-3668, August 01 2017.
- [12] R. Yu, M. Lee, Y. Huang, and S. Chen, "An efficient frequent patterns mining algorithm based on MaPreduce framework," *IET Conference Proceedings*, pp. 1-5
- [13] S. Moens, E. Aksehirli, and B. Goethals, "Frequent Itemset Mining for Big Data," in *2013 IEEE International Conference on Big Data*, 2013, pp. 111-118.

Archived at SID.ir