



ارائه شده توسط:

سایت ترجمه فا

مرجع جدیدترین مقالات ترجمه شده

از نشریات معتبر

طراحی پایگاه داده توزیع شده: یک مطالعه موردی

چکیده

تخصیص داده از مسائل مهم در طراحی پایگاه داده توزیع شده است. به طور کلی، الگوریتم های تکاملی برای تعیین تکالیف قطعات برای سایت استفاده می شوند. الگوریتم های تخصیص داده باید تکرار، تعداد پرس و جو، کیفیت سرویس (QoS)، ظرفیت های استناد، هزینه های به روز رسانی جدول، هزینه های انتخاب و طرح ریزی را هدایت نمایند. بسیاری از الگوریتم ها در نوشته ها به یک یا چند جزء از مسئله روی آورده اند. در این مقاله، ما یک مطالعه موردی را با در نظر گرفتن تمام این ویژگی ها ارائه می دهیم. مدل ارائه شده از برنامه نویسی انتگرال خطی برای فرمولاسیون این مسئله استفاده می کند.

واژه های کلیدی: پایگاه داده توزیع شده، تکرار، تخصیص داده ها

1. مقدمه

دفعات پاسخ پرس و جو، کیفیت سرویس (QoS)، ثبات و تمامیت داده ها، در کاربردهای سیستم مدیریت پایگاه داده توزیع شده (DDBMS) بسیار مهم است. در یک DDBMS، جداول و قطعات بر روی سایت های مختلف توزیع می شوند. هر پرس و جو از یک سایت اجرا می شود. کل هزینه شامل اجرای برنامه پرس و جو و هزینه دسترسی های جدول / قطعه از طریق شبکه می شود. مسئله تخصیص داده، NP-کامل است. بنابراین، الگوریتم های تکاملی به طور کلی برای پیدا کردن یک راه حل با هزینه حداقل برای مسئله استفاده می شوند. الگوریتم های تخصیص داده برای به حداقل رساندن هزینه دسترسی جدول / قطعه از پرس و جوها تلاش می نمایند. آنها یک تخصیص بهینه از جداول / قطعات را برای سایت ها پیدا می کنند. آنها همچنین پارامترهایی مانند داده های افزونه، هزینه های به روز رسانی جدول، و ظرفیت های سایت را در نظر می گیرند. عوامل متعددی در هنگام طراحی یک DDBMS باید در نظر گرفته شوند. پرس و جوهای به کارگرفته شده می توانند دارای وظایف مشترک باشند و پرس و جوهای مشابه

ممکن است از سایت های مختلف سرچشمه گیرند. ظرفیت های سایت، عناصر پردازش، ذخیره سازی و دفعات پاسخ پرس و جو در یک زمان به کار گرفته می شوند. بنابراین، این مسئله، ماهیت یک مسئله بهینه سازی چند هدف را نشان می دهد. ما یک مدل با برنامه نویسی انتگرال خطی (ILP) را طراحی نمودیم که دارای توانایی صدور هر یک از این عوامل به عنوان محدودیت ها است. توپولوژی شبکه، تکرار، هزینه های جدول / بروز رسانی قطعه، سایت های منشا، ظرفیت سایت، فرکانس پرس و جو، همه و همه می توانند به عنوان محدودیت در این فرمول تعریف شوند.

2. کار مرتبط

الگوریتم های ژنتیکی،¹،²،³ بازپخت شبیه سازی شده و بازپخت میدانی متوسط³ و ابتکارات کلونی مورچه ها⁴ برخی از روش ها حل مسئله تخصیص داده ها در متون مختلف هستند. همه این روش ها، یک یا چند ویژگی را از مسئله حذف می کنند.¹ Corcoran و² Frieder ظرفیت های سایت و داده های افزونه را در نظر نمی گیرند. الگوریتم ژنتیک، بازپخت شبیه سازی شده و راه حل های بازپخت میدانی متوسط پیشنهاد شده توسط³ Ahmad، داده های غیرافزونه را در نظر می گیرد. روش کلونی مورچه های پیشنهاد شده توسط⁴ Adl، این مسئله را به عنوان یک مسئله انتساب درجه دوم مدلسازی می کند. با این حال، هزینه های به روز رسانی و هزینه های تکرار در این کار به کار گرفته می شوند. چند الگوریتم با برنامه نویسی انتگرال خطی⁵،⁶،⁷ مطرح شدند. این الگوریتم ها به طور کلی ساده هستند و یک طرح پرس و جوی واقع بینانه و یا توپولوژی شبکه را پوشش نمی دهند. این فرمول ها تنها بخش کوچکی از مسئله را مد نظر قرار می دهند. آنها بخش خاصی از مسئله مانند تخصیص قطعات به سایت های افقی / عمودی یا تخصیص داده های غیرافزونه را در نظر می گیرند. Cornell و⁵ Yu، یک روش برای تعیین روابط و پیوستن عملیات ها به سایت ها را ارائه نمود. الگوریتم آنها برای به حداقل رساندن هزینه های ارتباطی تلاش می کند و هدف آن، استفاده از منابع و در عین حال اختصاص دادن قطعات به سایت ها و اجرای وظایف پیوسته است. این الگوریتم فاقد تجسم مسئله به عنوان ترکیبی از بهینه سازی پرس و جو، استفاده از شبکه و تخصیص داده هاست. روش پیشنهادی برای حل این مسائل به طور جداگانه تلاش می کند. علاوه بر این، فرمولاسیون برنامه نویسی خطی عدد صحیح پیشنهادی پیچیده است. Ailamaki و⁸ Papadomanolakis نیز از ILP برای نشان دادن مرز

کارآمد و واقع بینانه برای انتخاب شاخص استفاده نموده اند. آنها ادعا می کنند که رویکرد پیشنهادی، راه حل های محدود را می یابد.

3. طراحی پایگاه داده توزیع شده و برنامه نویسی خطی عددی

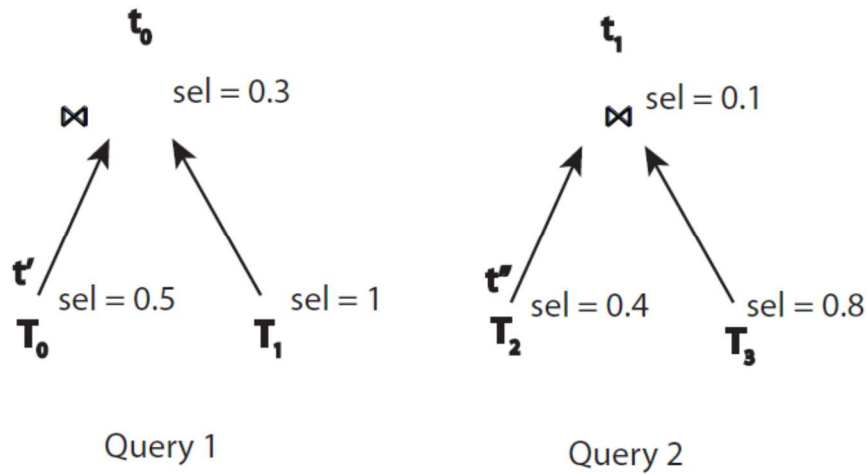
الگوریتم اول، تمام فواصل میان سایت ها را توسط الگوریتم کوتاه ترین مسیر⁹ Dijkstra محاسبه می نماید. فرض می شود که همه پرس و جوهای ورودی عمیق هستند. جداول پایه به عنوان برگ های درختان پرس و جو در نظر گرفته می شوند. درختان پرس و جو نمونه که در مورد مطالعه ما استفاده می شوند در شکل 1 نشان داده شده است. جداول پایه با حرف بزرگ T نشان داده شده اند و جداول به صورت T_0 تا T_n نامیده شده اند که در آن n تعداد جداول است. به طور مشابه، پیوندها به صورت t_0 تا t_n نامیده می شوند. عامل انتخاب، نسبت داده هایی است که باید بعد از عمل الحاق منتقل شود. همچنین جداول پایه می توانند توسط یک عامل انتخاب کوتاه شوند.

مدل شبکه ما شامل سرعت های ارتباط لینک مختلف می شوند، همانطور که در شکل 2 نشان داده شده است. سه سایت S_0 ، S_1 و S_2 وجود دارند. سایت ها دارای ظرفیت $C_0 = 18MB$ ، $C_1 = 15MB$ and $C_2 = 10MB$ هستند. لینک ها دارای سرعت های ارتباط از 100 کیلوبیت بر ثانیه، 200 کیلوبیت بر ثانیه و 500 کیلوبیت بر ثانیه هستند. پرس و جوهای انجام شده در شکل 1 نشان داده شده است. سه سایت و چهار جدول پایه وجود دارند. به منظور نشان دادن تکالیف جداول پایه برای سایت ها، ما از رسمی سازی در جدول 1 استفاده نمودیم. تعداد کل متغیرها، 12 برای تکالیف سایت-جدول داده شده برای مثال مسئله مثال است. 30 محدودیت وجود دارند و 8 تا از آنها برای محدودیت هایی هستند که بیان می کنند آیا المثنی ها برای هر رابطه با دادن تعداد المثنی به عنوان محدودیت مجاز هستند (به عنوان مثال 1 به معنی عدم تکرار برای جدول مربوطه). در مرحله بعد، 4 محدودیت داده می شود تا اطمینان حاصل شود که الزامات ذخیره سازی کلی برای جداول منسوب به سایت های خاص از ظرفیت های ذخیره سازی هر سایت تجاوز نمی کند.

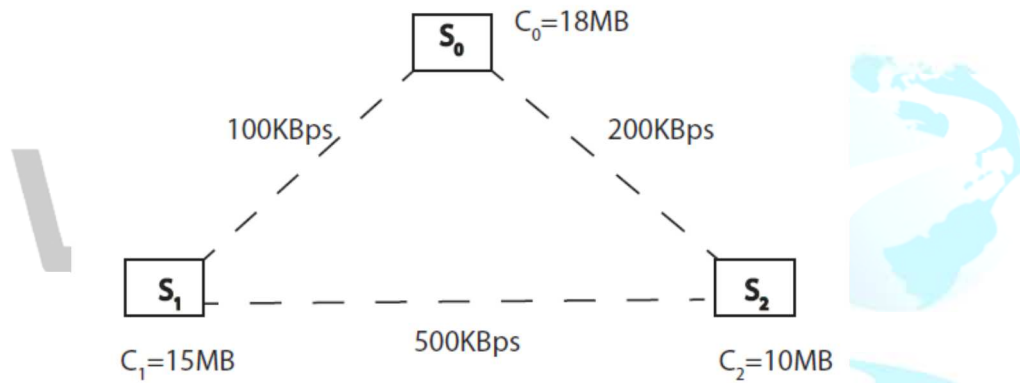
29 معادلات مورد استفاده برای مشخص کردن گره هایی که هر عمل را از پرس و جوهای داده شده انجام می دهند وجود دارند و تابع هدف شامل هزینه ارتباطات برای هر مسیر ممکن انتخاب شده وجود دارد. 2 پرس و جو مورد

استفاده در نمونه های ما وجود دارند. پرس و جوی 1، 100 بار از سایت منشاء S_0 اجرا می شود. پرس و جوی 2، 20 بار از سایت S_2 را اجرا می شود. در مثال های ما، حداکثر 2 کپی را برای همه جداول در نظر می گیریم. پارامترهای نسبت های به روز رسانی به صورت 0.1، 0.05، 0.5 و 0.1 برای جداول T_0-T_3 انتخاب می شوند. هزینه های به روز رسانی ها توسط ضرب نسبت در اندازه جدول محاسبه می شوند. تابع هدف مسئله بهینه سازی، به حداقل رساندن مجموع هزینه های انتقال جداول پایه و نتایج متوسط مورد استفاده توسط پرس و جویها برای سایت های S_0 ، S_1 و S_2 در حین اجرای پرس و جویهاست. اندازه جدول برای $T_0 = 10MB$ و $T_1 = 8MB$ که $10MB \times 0.5 = 5MB$ را برای T_0 و $8MB \times 1 = 8MB$ را برای T_1 ارائه می دهد که در آن 0.5 و 1، مقادیر انتخاب جدول هستند. هنگام انجام عملیات پیوند، رابطه متوسط حاصل $t\alpha$ به عنوان $5MB \times 8MB \times 0.3 = 12MB$ محاسبه می شود که در آن 0.3 انتخاب پیوند است. شبیه به پرس و جوی 1، پرس و جوی 2 دارای دو جدول $T_2 = 6MB$ و $T_3 = 5MB$ است. این به $6MB \times 0.4 = 2.4MB$ و $5MB \times 0.8 = 4MB$ منجر می شود که در آن 0.4 و 0.8 مقادیر انتخاب جدول هستند. هنگام انجام عملیات پیوستن، اندازه نتیجه t_1 رابطه متوسط به صورت $2.4MB \times 4MB \times 0.1 = 0.96MB$ محاسبه می شود که در آن 0.1 انتخاب پیوند است.

در مجموع 48 متغیر مورد استفاده برای نشان دادن مسئله بهینه سازی به عنوان یک مدل برنامه نویسی خطی وجود دارد. 12 مقدار، نشان دهنده جدول تکالیف سایت است، همانطور که در جدول 1 نشان داده شده است در حالی که 36 متغیر نشان دهنده هزینه های ارتباطی و هزینه های به روز رسانی. 8 معادلات مربوط به تکرار فرمول وجود دارند. معادله 1 تا معادله 4 نشان دهنده حداقل تعداد جداولی است که باید در سایت درج شود.



شکل 1. درختان پرس و جوی مورد استفاده در نمونه ما.



شکل 2. ظرفیت سایت: 18 مگابایت، 15 مگابایت، 10 مگابایت-ها: 100 کیلوبایت، 200 کیلوبایت، 500 کیلوبایت

TarjomeFa.Com

جدول 1. جدول برای متغیرهای انتساب سایت مورد استفاده در مدل ما.

Table \ Site	S_0	S_1	S_2
T_0	x_1	x_2	x_3
T_1	x_4	x_5	x_6
T_2	x_7	x_8	x_9
T_3	x_{10}	x_{11}	x_{12}

ما از این محدودیت ها استفاده نمودیم، زیرا برنامه نویسی خطی با هدف تعیین متغیرها صورت می گیرد. معادله 5 تا رابطه 8 نشان دهنده حداکثر تعداد جداولی است که باید در سایت درج شود. به طور کلی این سیستم سعی می کند

تا از المثنی ها استفاده نماید، زمانی که هزینه های به روز رسانی صفر هستند. هر پرس و جو برای بهره برداری از جداولی که در سایت خود منشاء استفاده می کند، تلاش می نماید.

$$\begin{array}{ll} x_1 + x_2 + x_3 \geq 1 & (1) \quad x_4 + x_5 + x_6 \geq 1 & (2) \\ x_7 + x_8 + x_9 \geq 1 & (3) \quad x_{10} + x_{11} + x_{12} \geq 1 & (4) \\ x_1 + x_2 + x_3 \leq 2 & (5) \quad x_4 + x_5 + x_6 \leq 2 & (6) \\ x_7 + x_8 + x_9 \leq 2 & (7) \quad x_{10} + x_{11} + x_{12} \leq 2 & (8) \end{array}$$

ظرفیت های سایت با معادله 9 تا معادله 11 نشان داده شده اند. هزینه های ارتباطی، مهم ترین بخش از سیستم هستند. مهم ترین مسئله با متغیرهای هزینه های ارتباطی، انتخاب متغیر برای بخش هایی از پرس و جو است که باید سازگار باشد. پرس و جوی 2 نشأت گرفته از S_2 توسط معادلات 12 تا معادله 20 بیان می شود.

$$\begin{array}{ll} 10x_1 + 8x_4 + 6x_7 + 5x_{10} \leq 18 & (9) \quad 10x_2 + 8x_5 + 6x_8 + 5x_{11} \leq 15 & (10) \\ 10x_3 + 8x_6 + 6x_9 + 5x_{12} \leq 10 & (11) \quad -x_{12} + x_{42} + x_{45} + x_{48} = 0 & (12) \\ -x_{11} + x_{41} + x_{44} + x_{47} = 0 & (13) \quad -x_{10} + x_{42} + x_{23} + x_{46} = 0 & (14) \\ -x_{33} - x_{36} - x_{39} + x_{46} + x_{47} + x_{48} = 0 & (15) \quad -x_{32} - x_{35} - x_{38} + x_{43} + x_{44} + x_{45} = 0 & (16) \\ -x_{31} - x_{34} - x_{37} + x_{40} + x_{41} + x_{42} = 0 & (17) \quad -x_9 - x_{39} - x_{38} - x_{37} = 0 & (18) \\ -x_8 - x_{36} - x_{35} - x_{34} = 0 & (19) \quad -x_7 + x_{33} + x_{32} + x_{31} = 0 & (20) \end{array}$$

تابع هدف که متشکل از هزینه های به روز رسانی و هزینه های ارتباطی است به شرح زیر برای شکل 2. محاسبه شده اند. لینک های ارتباطات دارای هزینه 10 ثانیه برای S_1-S_0 ، 5 ثانیه برای S_2-S_0 و در نهایت 2 ثانیه برای لینک S_1-S_2 می باشد. این هزینه ها، هزینه های متوسط برای انتقال 1 مگابایت اطلاعات بین دو سایت می باشند. ما می دانیم که نسبت به روز رسانی برای جداول مربوطه 0.1، 0.05، 0.5 و 0.1 هستند. در نهایت، تابع هدف برای شکل 2، معادله 21 است. پس از اجرای این طرح مثال ما، جداول T_0 و T_1 در سایت S_0 قرار دارد و جداول T_2 و T_3 در سایت S_1 قرار داده می شوند.

$$\begin{aligned} & x_1 + x_2 + x_3 + 0.4x_4 + 0.4x_5 + 0.46x_6 + 3x_7 + 3x_8 + 3x_9 + 0.5x_{10} + 0.5x_{11} + 0.5x_{12} + 0x_{13} + 5000x_{14} + \\ & 2500x_{15} + 5000x_{16} + 0x_{17} + 1000x_{18} + 2500x_{19} + 1000x_{20} + 0x_{21} + 0x_{22} + 17000x_{23} + 8500x_{24} + 5000x_{25} + \\ & 12000x_{26} + 7000x_{27} + 1000x_{28} + 13000x_{29} + 6000x_{30} + 0x_{31} + 480x_{32} + 240x_{33} + 480x_{34} + 0x_{35} + 960x_{36} + \\ & 240x_{37} + 96x_{38} + 0x_{39} + 96x_{40} + 518.4x_{41} + 240x_{42} + 576x_{43} + 96x_{44} + 96x_{45} + 336x_{46} + 192x_{47} + 0x_{48} \end{aligned} \quad (21)$$

4. نتیجه گیری

در این مقاله، یک فرمولاسیون برنامه نویسی انتگرال خطی برای حل مسئله تخصیص داده در پایگاه داده توزیع شده ارائه شده است. مدل ارائه شده دقیقاً مسائلی مانند ظرفیت سایت، فراوانی پرس و جو و هزینه های ارتباطی را هدایت

می کند. این مدل با قطعه قطعه شدن و همان پرس و جویهای نشات گرفته از سایت های مختلف سرو کار ندارد. انتخاب توپولوژی شبکه مناسب، هزینه های عملیات شبکه و دفعات پاسخ پرس و جو نیز از عوامل دیگر هستند که باید در یک طراحی واقعی به کار گرفته شوند. ما قصد داریم این الگوریتم را برای پرس و جویهای وظیفه-به اشتراک گذاشته و مدیریت قطعه در آینده گسترش دهیم. توازن بار و اجرای وظیفه همزمان، معیارهای دیگر هستند که باید در کار آینده به کار گرفته شوند.

References

1. A.L. Corcoran, and J. Hale, "A Genetic Algorithm for Fragment Allocation in a Distributed Database System," *In Proc. 1994 Symp. on Applied Computing*, pp. 247-250, 1994
2. O. Frieder, and H. T. Siegelmann, "Multiprocessor Document Allocation: A Genetic Algorithm Approach," *Transactions on Knowledge and Data Engineering*, vol.9, no.4 , 1997, pp.640642
3. I. Ahmad, K. Karlapalem, Y. Kwok, and S. So, "Evolutionary algorithms for allocating data in distributed database systems," *International Journal of Distributed and Parallel Databases*, vol. 11, no. 1, pp. 532, 2002
4. R.K. Adl, and S.M.T.R. Rankoohi, "A new ant colony optimization based algorithm for data allocation problem in distributed databases," *Knowledge and Information Systems*, vol. 20, no. 3, pp. 349-372, 2009.
5. D.W. Cornell and P.S. Yu, "Site assignment for relations and join operations in the distributed transaction processing environment," *In Proc. Fourth Int. Conf. on Data Eng.*, pp. 100-108, 1988.
6. B.Gavish and H. Pirkul, "Computer and database location in distributed computer systems," *IEEE Transactions on Computers*, vol. C-35, no. 7, pp. 583-590, 1986.
7. S. Ram and R.E. Marsten, "A model for database allocation incorporating a concurrency control mechanism," *IEEE Transactions on Knowledge and Data Engineering*, vol. 3, no. 3, pp. 389-395, 1991.
8. S. Papadomanolakis and A. Ailamaki. "An integer linear programming approach to database design," *In Proc. of the 2007 IEEE 23rd Int. Conf. on Data Eng. Workshop*, p.442-449,2007.
9. E. W. Dijkstra, "A Note on Two Problems in Connection with Graphs," *Numeriche Mathematik*, 1:269-271, 1959.

برای خرید فرمت ورد این ترجمه، بدون واتر مارک، اینجا کلیک نمایید.



این مقاله، از سری مقالات ترجمه شده رایگان سایت ترجمه فا میباشد که با فرمت PDF در اختیار شما عزیزان قرار گرفته است. در صورت تمایل میتوانید با کلیک بر روی دکمه های زیر از سایر مقالات نیز استفاده نمایید:

لیست مقالات ترجمه شده ✓

لیست مقالات ترجمه شده رایگان ✓

لیست جدیدترین مقالات انگلیسی ISI ✓

سایت ترجمه فا ؛ مرجع جدیدترین مقالات ترجمه شده از نشریات معتبر خارجی