

مقایسه ی عملکرد پرس و جو های KNN بر روی پایگاه داده های مکانی

مهدی سامانی^۱، مهدی فولادگر^۲ و مژگان چاکری^۳
دانشگاه صنعتی اصفهان، m.samani@ec.iut.ac.ir
دانشگاه صنعتی اصفهان، m.fooladgar@ec.iut.ac.ir
دانشگاه صنعتی اصفهان، m.chakeriyazdi@ec.iut.ac.ir

چکیده – امروزه پایگاه داده های مکانی به منظور انجام پرس و جو های مکانی نظیر یافتن فاصله، یافتن K نزدیک ترین همسایه (KNN) و یافتن محیط استفاده می شوند. پرس و جو KNN یک پرس و جو کاربردی است که زیربنای بسیاری از نرم افزار های مبتنی بر داده های مکانی را تشکیل می دهد. این پرس و جو ها زمانی داده ها در حجم بالا مورد استفاده قرار می گیرد بسیار زمانبر خواهد شد. در این مقاله هدف مقایسه ای میان عملکرد پایگاه داده های مکانی $GeoSpark$ ($Non\ Relational, In-memory$)، $MemSQL$ ($Relational, In-memory$) و $PostGIS$ ($Relational, Disk\ database$) برای انجام پرس و جو KNN می باشد. همچنین به معرفی پایگاه داده های مختلف و بررسی ساختار ایندکس ها مکانی مانند $R-Tree$ برای تسریع در سرعت اینگونه پرس و جو ها در پایگاه داده های مختلف انجام شده است.

کلید واژه – KNN , $GeoSpark$, $MemSQL$, $PostGIS$, $R-Tree$

و تجزیه و تحلیل داده های مکانی در مقیاس بزرگ را در سراسر ماشین ها به طور موثر انجام می دهد.

۱- مقدمه

$AsterixDB$ پایگاه داده ای است که توسط شرکت $Apache$ توسعه داده شده و پشتیبانی می شود. این پایگاه داده به کاربران اجازه می دهد از بسترهای فایل سیستم های توزیع شده^۱ استفاده کرده و پایگاه داده ای نیم-ساختاریافته^۲ را با زبانی مشابه با SQL مدیریت کرده و به اجرای جست و جوی خود پردازند. این پایگاه داده به صورت پیش فرض مفهوم نقطه و ایندکس های مکانی را پشتیبانی می کند.

$Neo4j$ یک کتابخانه و پلتفرم مبتنی بر جاوا برای پیاده سازی و مدل سازی گراف هاست. $Neo4JSpatial$ یک افزونه

امروزه حجم داده های مکانی به سرعت رو به افزایش است و در نتیجه آن، پرس و جو های مکانی نیز در حال توسعه و گسترش می باشند. پرس و جو یافتن k نزدیک ترین همسایه (KNN) پرس و جویی است که شاکله ی اصلی بسیاری از برنامه های کاربردی و پرس و جو های مختلف، بر مبنای آن قرار گرفته است. به عنوان نمونه می توان برنامه هایی را در نظر گرفت که نیاز به یافتن k نزدیک ترین خودرو در اطراف یک نقطه ی مشخص دارد. در این صورت زمانی که با حجم انبوهی از داده سروکار داریم علاوه بر یافتن صحیح نقاط مورد نظر زمان یافتن این نقاط نیز بسیار اهمیت خواهد داشت.

امروزه پایگاه داده های مکانی مختلفی برای انجام چنین پرس و جوهایی به کمک کاربران آمده اند و با ارائه ی الگوریتم های متنوع و استفاده از ساختارهای ایندکس گوناگون تلاش به سرعت بخشیدن به این پرس و جوها کرده اند. از جمله این پایگاه داده ها می توان به $GeoSpark$ ، $AsterixDB$ ، $Neo4j$ ، $MemSQL$ و $PostGIS$ اشاره کرد.

$GeoSpark$ یک سیستم محاسباتی خوشه ای برای پردازش داده های مکانی در مقیاس بزرگ است. $GeoSpark$ یک توسعه ای از $Apache Spark$ است که با بهره گیری از $Spatial Resilient Distributed Datasets$ ($SRDDs$) بارگذاری، پردازش

^۱ Distributed

^۲ Semi-Structured

برای Neo4j است که توابع مورد نیاز برای پرس و جوهای مکانی را فراهم آورده است.

MemSQL یک پایگاه داده مقیم در حافظه^۲ رابطه‌ای است این پایگاه داده می‌تواند داده‌ها را به صورت سطری و یا ستونی ذخیره نماید. این پایگاه داده با پیاده‌سازی توابع مکانی به بستر مناسبی برای اجرای پرس و جوهای مکانی تبدیل شده است.

GeoMesa یک پایگاه داده ی توزیع شده ی مکانی-زمانی متن باز است که بر روی سیستم های ذخیره سازی ابری مانند HBase, Accumulo, Cassandra کار می کند. در واقع GeoMesa با استفاده از پردازش موازی و استراتژی ایندکس گذاری همانند عملکرد PostGIS بر روی PostgreSQL، پرس و جو های مکانی-زمانی را بر روی این سیستم های توزیع شده فراهم می سازد. (متاسفانه مشکلاتی که در اجرای این پایگاه داده بر روی این سیستم های نام برده شده به وجود آمد باعث شد که آن را از آزمایشات و مقایسه ها کنار گذاریم. این مشکلات بیشتر در زمینه ی راه اندازی و اجرای پرس و جو های مکانی بر روی این بستر همچنین عدم وجود داکيومنت مناسب می باشد.)

PostGIS نرم افزاری منبع باز است که به منظور بهره گیری از داده های مکانی به پایگاه داده رابطه ای PostgreSQL اضافه می شود. در واقع postgis با پشتیبانی از داده های جغرافیایی، این امکان را برای کاربر فراهم می کند تا پرس و جوهای مکانی را در SQL اجرا کند .

در این مقاله پس از معرفی ابزار ها و مدل های لازم برای بررسی قدرت پایگاه داده های مختلف در اجرای پرس و جوهای مکانی، به معرفی و بررسی الگوریتم های KNN و Similarity Path پرداخته می شود و در نهایت در پایگاه داده های مختلف با توجه به خصوصیات آن ها و بر اساس پارامتر های متفاوت مورد مقایسه و ارزیابی قرار می گیرند.

^۲ In Memory

۲- معرفی ابزار و مدل ها

پس از بررسی های اولیه پایگاه داده های متنوعی نظیر Spark, PostgreSQL, AsterixDB, Neo4JSpatial و ... به عنوان کاندیدهای اولیه برای اجرای پرس و جوهای KNN مطرح شدند. که هر یک از این پایگاه های داده ویژگی های منحصر به فرد برای نوع خاصی از داده ها و پرس و جوها را تامین می کنند. در نهایت پس از بررسی از میان این پایگاه داده ها سه پایگاه داده به عنوان کاندیدا مقایسه انتخاب شدند.

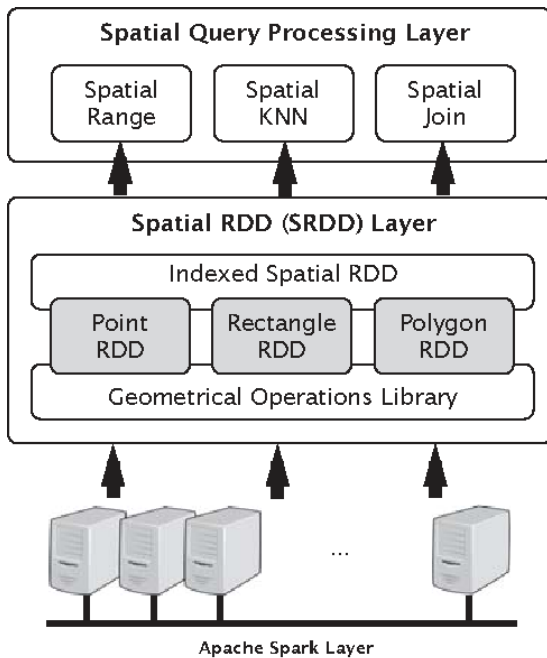
در این قسمت ما ابتدا به معرفی دیتاستی انتخابی خود برای انجام مقایسات و سپس به بررسی هر یک از این پایگاه های داده و مراحل پیاده سازی و اجرای پرس و جوهای KNN در آن ها خواهیم پرداخت.

۲-۱- Spark

Apache Spark یک سیستم محاسباتی خوشه ای در حافظه است. اسپارک یک سطح انتزاع از داده را فراهم می کند که resilient distributed datasets نام دارد و داده ها بر روی تمامی ماشین ها توزیع شده اند. RDD های موجود در حافظه باعث شده که اسپارک نسبت به رقیب های خود در امر پردازش های توزیع شده مانند MapReduce بهتر عمل کند [1]. متاسفانه spark خود به تنهایی داده های مکانی را پشتیبانی نمی کند و برای کار با داده های مکانی نیاز به یک افزونه مانند GeoSpark است.

۲-۲- GeoSpark

GeoSpark یک فریم ورک برای پردازش داده های مکانی با حجم بسیار زیاد است که داده ها را در حافظه ی اصلی ذخیره می کند و این باعث تسریع در امر محاسبات و پردازش ها خواهد شد. GeoSpark علاوه بر پشتیبانی از داده های مکانی از ایندکس و سایر عملیات بر روی داده های مکانی مانند KNN و... پشتیبانی می کند. اسپارک با پشتیبانی از ساختار Spatial Resilient Distributed Dataset (SRDD) باعث می شود ذخیره و پردازش داده های مکانی بر روی بستر اسپارک به راحتی انجام شود. همچنین GeoSpark از ساختار های R-tree و Quad-tree برای ایندکس گذاری بر روی پارتیشن های داده استفاده می کند که باعث سرعت بخشیدن به پردازش های مکانی خواهد شد. در شکل ۱ لایه های تشکیل دهنده ی GeoSpark نشان داده شده است. لایه ی اول Apache Spark را نشان می دهد که در قسمت Spark توضیح داده شد و عملیات ذخیره و بازیابی داده



شکل ۱ لایه های تشکیل دهنده ی GeoSpark

ها در این لایه انجام می شود. لایه ی دوم (SRDD) در واقع همان توسعه ی اسپارک برای ذخیره سازی داده های مکانی می باشد. PointRDD، RectangleRDD و PolygonRDD از جمله ساختارهای قابل تعریف در این فریم ورک می باشد. در نهایت در لایه ی آخر که Spatial Query Processing نام دارد پرس وجوهای مکانی مانند range، join و KNN را در بر میگیرد. برای پردازش پرس و جو KNN مکانی، GeoSpark با استفاده از الگوریتم top-k مبتنی بر heap که شامل دو مرحله عمل میکند: انتخاب^۴ و ادغام^۵. برای محاسبه نزدیک ترین اجسام اطراف نقطه P در مرحله انتخاب، برای هر پارتیشن SRDD، GeoSpark فاصله ی بین هر شی با نقطه داده شده را محاسبه می کند، سپس با افزودن یا حذف عناصر بر اساس فاصله، یک هیپ^۶ محلی به وجود می آورد. در این هیپ k نزدیک ترین اجسام به نقطه ی P وجود خواهند داشت [1].

۳-۲- PostGIS

همانگونه که قبلا ذکر شد، postgis، پرس و جو های مکانی روی پایگاه داده های رابطه ای را ممکن می سازد بدین منظور کافی است ابتدا در پایگاه داده PostgreSQL، PostGIS extension را نصب نمود و جداول، شماها، توابع، ایندکس ها و trigger های مکانی را ساخت و پرس و جو های مکانی مورد نظر را اعمال نمود. پس از ساخت جداول و توابع و ... می توان از فایل هایی با فرمت Text، Binary و CSV داده ها را با تعیین پارامتر های مختلف چون Delimiter، OID و ... به جداول وارد نمود. برای کار با داده های مکانی عموما نیاز است که در این داده های وارد شده تغییراتی اعمال شود تا برای اعمال پرس و

^۴ selection

^۵ merge

^۶ heap

جو های مکانی آماده سازی شوند. برای مثال در فایل وارد شده، مکان در هر رکورد به صورت نقطه ذخیره نشده است بلکه طول و عرض جغرافیایی آن به صورت جداگانه قرار داده شده است. حال برای آماده سازی این نوع فایل برای اعمال پرس و جو های مکانی کافیست با استفاده از توابع موجود در PostGIS دو attribute طول و عرض جغرافیایی را به یک attribute مکانی^۷ نقطه جغرافیایی تبدیل کرد و بدین ترتیب فایل داده ها را برای اعمال پرس و جو های مکانی آماده می شود. به طور کلی، در PostGIS می توان attribute های مکانی مختلف از جمله Point، Polygon، MultiLineString و ... را ساخت و استفاده نمود.

برای درک بیشتر از مفاهیم PostGIS، ساختار آن به طور خلاصه در ادامه بیان می شود. همانطور که گفته شد، PostGIS در واقع پلاگینی است بر روی پایگاه داده PostgreSQL که کتابخانه های مختلفی دارد که عبارتست از :

- GEOS: شامل توابع پیش فرض، نوع داده ها و عملگرهای مکانی چون Union، Intersects، Point، Distance و....
- Proj4: شامل فیلتر های استاندارد UNIX که برای تبدیل مختصات جغرافیایی به مختصات دکارتی مورد استفاده قرار می گیرد.
- LibXML2: کتابخانه ای که می تواند برنامه به زبان های مختلف C++، python و ... را به یکدیگر مرتبط نماید.

در این گزارش، از کتابخانه GEOS این پلاگین استفاده شده است و از توابع و عملگر های پیش فرض آن استفاده شده است. در ادامه این بخش، به معرفی چند تابع و نوع داده و عملگر پیش فرضی که در این گزارش از آن استفاده شده است پرداخته می شود.

ST_MakePoint: تابعی که با گرفتن دو مقدار اعشاری، نقطه ای به مختصات آن می سازد. اما در PostGIS باید به هر نقطه ای که ساخته می شود یک شناسه خاص یا SRID نسبت داد که این کار با تابع ST_SetSRID و با در نظر گرفتن یک ID مثلا ۴۳۲۶ برای نقطه خاصی ممکن می شود. اما برای تبدیل یک نقطه که به صورت متنی از کاربر گرفته می شود، شناسه SRID با استفاده از تابع ST_GeomFromText ممکن است.

ST_Distance: با گرفتن دو مقدار Geometry، فاصله آن دو را بر می گرداند. که این فاصله بر حسب متر می باشد.

ST_Makeline: این تابع با گرفتن چندین نقطه یک خط که می تواند باز یا بسته باشد (مثل چندضلعی) را بسازد. این تابع برای ساخت مسیر در این گزارش استفاده شده است.

ST_HausdorffDistance: برای برگرداندن فاصله Hausdorff دو داده geometry استفاده می شود که برای اندازه گیری شباهت یا عدم شباهت دو داده مقدراری اعشاری را بر می گرداند.

توابع پیش فرض زیادی در این پلاگین پیاده سازی شده است اما در اینجا چند تا از توابع معروف و استفاده شده در این گزارش بیان شد. [3]

۴-۲- AsterixDB

همانگونه که قبلا ذکر شد، AsterixDB، به ما اجازه می دهد تا پرس و جو های مکانی را در محیط های توزیع شده و بر روی داده های نیم-ساختاریافته اجرا کرده و نتیجه را ببینیم. علاوه بر آن CloudBerry نیز با برقراری ارتباط با این پایگاه داده به ما در نمایش^۸ در آوردن داده ها کمک خواهد کرد. این پایگاه داده برای اجرای پرس و جو های پیشرفته ای نظیر KNN راه حل آماده ای ندارد لذا یکی از ساده ترین راه حل ها تعریف تابع فاصله و اجرای پرس و جوی Top-K بر روی آن هاست.

^۸ Visualize

^۷ Geography type

۵-۲- MemSQL

MemSQL یک پایگاه داده مقیم در حافظه است. این پایگاه داده با قابلیت ذخیره سازی داده ها به دو شکل سطری و ستونی و پیاده سازی توابع مکانی به گزینه خوبی برای اجرا پرس و جویهای مکانی تبدیل شده است. از طرف دیگر در اختیار داشتن محیط مناسبی برای مانیتور کردن منابع مصرفی این پایگاه داده را به ابزاری منحصر به فرد تبدیل می کند.



شکل ۳ محیط اجرای MemSQL

اما متاسفانه روش های اندیس گذاری در این پایگاه داده تنها محدود به B-Tree و Hash هستند که این موضوع کار را برای اندیس گذاری داده های مکانی سخت می کند. توابع مکانی پیاده سازی شده در این پایگاه داده عبارتند از:

- GEOGRAPHY_AREA
- GEOGRAPHY_POINT
- GEOGRAPHY_DISTANCE
- و ...

برای ورود اطلاعات دیتاست به پایگاه داده نیازمند ابزاری برای خواندن فایل های CSV و ورود آن ها به پایگاه داده ایم. خوشبختانه MemSQL به طور پیش فرض این ابزار را فراهم کرده است. پس ابتدا باید جدول مورد نیاز برای داده ها را بسازیم برای این کار پرس و جوی زیر کار است:

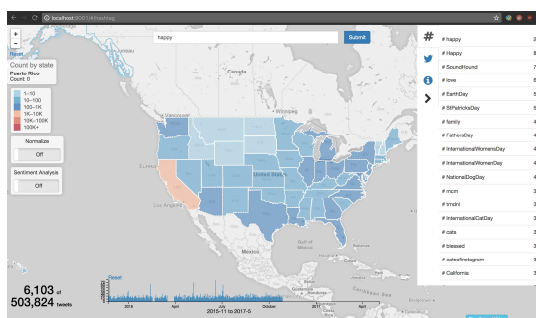
```
CREATE TABLE IF NOT EXISTS ADB (
  id BIGINT AUTO_INCREMENT,
  tid BIGINT,
  time CHAR(15),
  lat FLOAT,
  lon FLOAT,
  key(id)
);
```

شکل ۴ ساخت جدول

بدیهی است که این راهکار برای این حجم از داده راهکاری معقول و کارا نیست.

در این پایگاه داده می توان با یکی از دو زبان SQL++ و یا AQL پرس و جویهای مورد نیاز را نوشت. هر دو زبان توسط سازندگان AsterixDB گسترش یافته اند، لذا مستندات این زبان ها تنها در سایت رسمی AsterixDB موجود است و همچنین هنوز به بلوغ خود نرسیده است. به طور مثال در این زبان ابزار لازم برای تجزیه^۹ زمان و نقطه از فایل های CSV هنوز پیاده سازی نشده است.

یکی از نقاط مثبت این پایگاه داده وجود پلتفرم های نمایش داده ها با نام CloudBerry است. این پلتفرم محیطی ساده و شکیل برای نمایش داده ها در اختیار کاربران قرار می دهد. از طرف دیگر وجود ابزارهایی برای ورود استریم داده های نظیر توئیت های توئیتر نیز از دیگر نقاط مثبت این پایگاه داده است.



شکل ۲ نمونه ای از محیط اجرای AsterixDB

در مورد ایندکس های پیاده سازی شده در این پایگاه داده می توان روش اندیس گذاری R-Tree را نام برد.

^۹ Parse

باشند. به همین دلیل در پایگاه داده ها، ایندکس های مختلفی پیاده سازی شده است که یکی از این ایندکس ها، R-Tree می-باشد. برای پیاده سازی این نوع ایندکس در بعضی پایگاه داده ها مثل PostGIS نیاز است که ابتدا ایندکسی به نام GiST ساخته شود و سپس ایندکس R-Tree پیاده شود. در ادامه ساختار این نوع دو ایندکس مورد بررسی قرار گرفته می‌شود.

۳-۱- R-Tree Index

این نوع ایندکس برای نخستین بار توسط گوتمن^{۱۰} و برکلی^{۱۱} برای تمام داده های مکانی و چند بعدی در سال ۱۹۸۴ ارائه شد. ساختار سلسله مراتبی و متعادل این نوع ایندکس شبیه ساختار ایندکس B-Tree می‌باشد اما برای ذخیره داده، روشی متفاوت را پیش می‌گیرد. ایده‌ی اصلی این ساختار، دسته بندی اشیا و قرار دادن اشیا نزدیک بهم در یک ناحیه مستطیلی می-باشد. برای ساخت این درخت ابتدا اشیایی که فاصله شان از یکدیگر از یک آستانه مشخص کمتر باشد، در کوچکترین مستطیل ممکن قرار می‌گیرند و تمام اشیا موجود در آن مستطیل، در سطوح بعدی درخت با همان ناحیه در نظر گرفته شده، شناخته می‌شوند سپس مرز های ساخته شده نزدیک به هم در یک مستطیل بزرگتر قرار داده می‌شوند و به همین ترتیب این کار ادامه داده می‌شود تا یک مرز کلی ایجاد شود. در شکل ۷ روند ساخت این درخت قابل مشاهده است.

این ساختار ایندکس در هر نوع فضای مکانی قابل پیاده سازی می‌باشد. این بدین معنی است که اگر فضای مکانی دو بعدی باشد و داده ها را بتوان در دایره هایی با شعاع های مختلف در نظر گرفت، ساخت مستطیل های مرزی همچنان ممکن است و کوچکترین مستطیل های مرزی را می‌توان به اندازه قطر دایره داده ها، در نظر گرفت. نکته لازم به ذکر این است که ابعاد مرز-

پس از آن با اجرای دستوری داده‌ها را داخل پایگاه‌داده وارد کرده و با ویرایش، آن‌ها را تبدیل به داده‌های از نوع مکانی می‌کنیم:

```
LOAD DATA INFILE '<path-to-data>'
INTO TABLE ADB
COLUMNS TERMINATED BY ',';

ALTER TABLE ADB ADD COLUMN loc GEOGRAPHYPOINT;

UPDATE ADB SET loc=GEOGRAPHY_POINT(lat, lon);
```

شکل ۵ وارد کردن داده در MemSQL

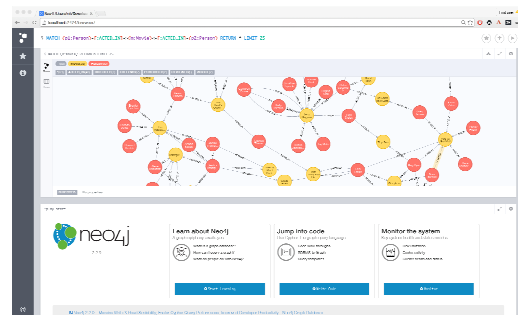
و پس از اجرای این پرس‌جو داده‌ها آماده پردازش هستند.

۶-۲- Neo4j

یکی دیگر از ابزارهای قابل استفاده برای تحلیل و اجرای پرس‌وجو بر روی داده‌های مکانی کتابخانه مبتنی بر جاوا Neo4j است. این کتابخانه به مدل‌سازی مسائل و نگاشت آن‌ها به گراف کمک می‌کند.

از طرف دیگر وجود ابزار Neo4jSpatial با پیاده‌سازی توابع مکانی و روش اندیس‌گذاری RTree برای داده‌های مکانی باعث شده تا بتوانیم از این کتابخانه نیز برای پرس‌وجوهای مکانی استفاده نماییم.

Neo4j با فراهم آوردن محیط گرافیکی مناسب و توضیحات و مستندات مناسب به نظر محیطی کارا برای دانشمندان علوم کامپیوتر و داده است اما از نظر یک متخصص پایگاه‌داده انعطاف لازم را ممکن است نداشته باشد.



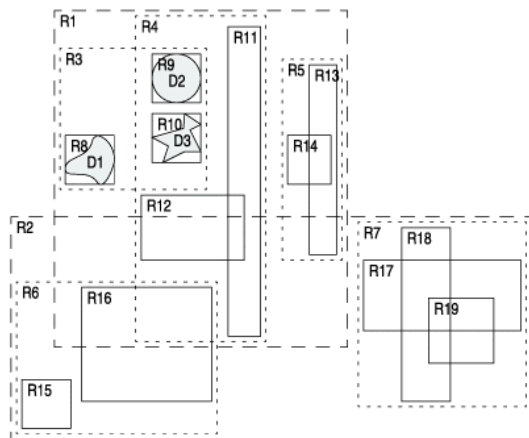
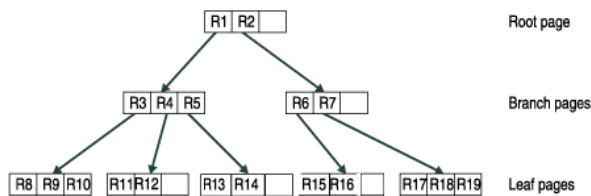
شکل ۶ نمونه‌ای از محیط Neo4j

۳- مدل‌ها و ایندکس‌ها

با توجه به گسترش داده‌ها و منابع اطلاعاتی، دیتاست‌های مورد استفاده نیز ممکن است بسیار بزرگ باشند و گاهی انجام پرس و جوها و بررسی مستقیم تمام رکورد ها نیز زمان بر

^{۱۰} Antonin Guttman

^{۱۱} UC Berkley



شکل ۷ ساختار R-Tree برای ایندکس مکانی

های تشکیل داده شده، می‌تواند متغیری به جز مکان، مانند زمان نیز باشد.

هرچند در ساختار تئوری R-Tree Index، در گره‌های برگ، داده‌ها باید ذخیره شوند اما در نسخه جدید پایگاه داده‌ها به صورت پیش فرض اینگونه نیست و کپی از داده‌ها در برگ‌نگه‌داری نمی‌شود بلکه محدوده آن در برگ ذخیره می‌شود که در این صورت به این نوع ایندکس، bounding-box-only R-tree گفته می‌شود که مزیت‌های این نوع ایندکس عبارتست از: استفاده بهینه از حافظه و زمان، کاهش فضای مورد نیاز برای به روز رسانی درخت و افزایش کارایی.

همانطور که ذکر شد، ساختار این نوع ایندکس، یک درخت متعادل است که این بدین معنی است که از ریشه تا برگ، هر مسیری طی شود، از تعداد گره‌های یکسانی عبور می‌شود. یا به عبارتی دیگر، تمام برگ‌ها در یک سطح قرار دارند.

ساختار R-Tree index، که در این گزارش مورد استفاده قرار گرفته است در گره‌های برگ علاوه بر کپی از کلید یا داده، اشاره‌گری به گره‌های قبلی نیز ذخیره می‌شود. همچنین در گره‌های میانی و ریشه، مرز شامل مرزهای کوچکتر و اشاره‌گر به گره‌های سطوح پایین‌تر ذخیره می‌شوند. بدین صورت پیمایش بالا به پایین درخت و پیمایش پایین به بالای درخت فراهم می‌شود. در این نوع ایندکس، با توجه به ضریب انشعاب^{۱۲} حداکثر و حداقل تعداد گره‌های ذخیره شده در هر سطح مشخص می‌شود. [1] [2]

۲-۳- GiST Index

برای پیاده‌سازی R-Tree در postgis لازم است که ابتدا GiST بر روی پایگاه داده ساخته شود و سپس روی آن، R-Tree را می‌توان در نظر گرفت. GiST، اشیاء را براساس اینکه موقعیت مکانیشان نسبت به ناحیه مورد نظر، چگونه است تقسیم بندی می‌کند و بدین ترتیب اشیایی که هیچ تداخلی با ناحیه ندارند در یک دسته، اشیایی که با ناحیه هم پوشانی دارند در دسته دیگر و اشیایی که درون ناحیه هستند نیز در دسته سوم قرار می‌گیرند.

نکته مثبت پیاده‌سازی R-Tree بر روی GiST را می‌توان

در دو مورد زیر خلاصه نمود:

۱- Null Safe: ایندکس بر روی داده‌هایی که ستون مورد

نظرشان Null است، نیز ساخته می‌شود.

۲- قابلیت Lossines: این قابلیت زمانی اهمیت پیدا می‌کند

که داده‌های GIS بزرگتر از 8K باشند. بدین ترتیب در

این مواقع، در ساخت ایندکس، قسمت‌های مهم داده

در ایندکس ذخیره می‌شود.

^{۱۲} Branching factor

۴- پرس و جوی KNN

پارامتر ورودی false در این تابع نشان دهنده ی اجرای پرس و جوی بدون استفاده از ساختار ایندکس می باشد.

در حالت بعدی همان روند با استفاده از ایندکس بر روی داده ها انجام می شود. همانطور که پیش از این نیز بیان شد در GeoSpark ایندکس ها بر روی هر RDD به صورت مستقل تعریف می شود. برای ساخت ایندکس بر روی داده ها از تابع buildIndex استفاده شده است و نوع ایندکس استفاده شده R-Tree می باشد.

نتایج حاصل از این پرس و جوی در قسمت ارزیابی و نتایج نشان داده شده است.

۴-۲ MemSQL

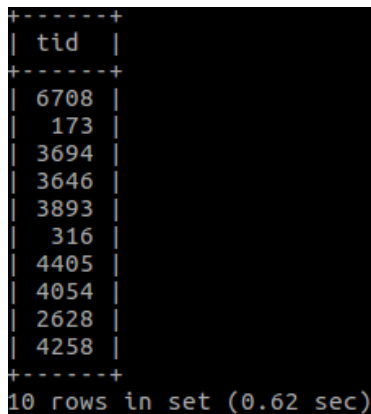
برای یافتن نزدیک ترین همسایگان به یک نقطه در MemSQL باید در وهله اول می توان از یکی از ابتدایی ترین روش ها یعنی مرتب سازی بر اساس فاصله تمامی نقاط دیتاست و یافتن K کمترین آن ها را مهم را تحقق بخشید. پرس و جوی لازم برای این کار با چند خط قابل به پیاده سازی است.

```
SELECT DISTINCT tid
FROM
  ADB
ORDER BY GEOGRAPHY_DISTANCE(Loc, GEOGRAPHY_POINT(-72.1235, 42.3521)) ASC
LIMIT 10;
```

شکل ۱۰ پرس و جوی KNN بر روی MemSQL

نمونه خروجی این پرس و جوی در شکل زیر قابل مشاهده

است:



tid
6708
173
3694
3646
3893
316
4405
4054
2628
4258

10 rows in set (0.62 sec)

شکل ۱۱ خروجی پرس و جوی KNN بر روی MemSQL

۴-۱ GeoSpark

پس از انجام مراحل نصب و راه اندازی GeoSpark با استفاده از زبان scala کوئری KNN برای دیتاست مربوطه را در مرحله ی اول بدون استفاده از ایندکس (شکل ۸) و در مرحله ی بعد با استفاده از ساختار ایندکس R-Tree (شکل ۹) اجرا می کنیم.

```
import org.datasyslab.geospark.spatialOperator.KNNQuery;
import org.datasyslab.geospark.spatialRDD.PointRDD;
import com.vividsolutions.jts.geom.GeometryFactory;
import com.vividsolutions.jts.geom.Point;
import com.vividsolutions.jts.geom.Coordinate;
import org.datasyslab.geospark.enums.FileDataSplitter;

val fact=new GeometryFactory();
val queryPoint=fact.createPoint(new Coordinate(-109.73, 35.08));
/* Range query window format: X Y */
val objectRDD = new PointRDD(sc, "/home/huser/Desktop/all.csv", 2, FileDataSplitter.CSV, false);
/*
 * 2 is the starting column of spatial data in the input file
 * FileDataSplitter.CSV means the data format is CSV. We CSV, TSV, MKT, GeoJSON and self-defined format mapper.
 * false means each spatial object doesn't need to carry the original input tuple with it.
 */
val resultSize = KNNQuery.SpatialKnnQuery(objectRDD, queryPoint, 5, false).size();
/* The number 5 means 5 nearest neighbors
 * The false means don't use spatial index.
 */
```

شکل ۸ پرس و جوی KNN بر روی بستر GeoSpark با استفاده از زبان Scala بدون استفاده از ساختار ایندکس

```
import org.datasyslab.geospark.spatialOperator.KNNQuery;
import org.datasyslab.geospark.spatialRDD.PointRDD;
import com.vividsolutions.jts.geom.GeometryFactory;
import com.vividsolutions.jts.geom.Point;
import com.vividsolutions.jts.geom.Coordinate;
import org.datasyslab.geospark.enums.FileDataSplitter;
import org.datasyslab.geospark.enums.IndexType;

val fact=new GeometryFactory();
val queryPoint=fact.createPoint(new Coordinate(-109.73, 35.08));
/* Range query window format: X Y */
val objectRDD = new PointRDD(sc, "/home/huser/Desktop/all.csv", 2, FileDataSplitter.CSV, false);
/*
 * 2 is the starting column of spatial data in the input file.
 * FileDataSplitter.CSV means the data format is CSV. We CSV, TSV, MKT, GeoJSON and self-defined format mapper.
 * false means each spatial object doesn't need to carry the original input tuple with it.
 */
objectRDD.buildIndex(IndexType.RTREE, false);
/*
 * false means just build index on original spatial RDD instead of spatial partitioned RDD.
 */
val resultSize = KNNQuery.SpatialKnnQuery(objectRDD, queryPoint, 5, true).size();
/* The number 5 means 5 nearest neighbors
 * The true means use spatial index.
 */
```

شکل ۹ پرس و جوی KNN بر روی بستر GeoSpark با استفاده از زبان Scala با استفاده از ساختار ایندکس

قبل از اجرای پرس و جوی باید فایل دیتاست را در محیط هادوپ اضافه کنیم که چون محیط هادوپ به گونه ای تنظیم شده که فایل ها را به صورت MB۱۲۸ ذخیره می کند پس در این صورت فایل دیتاست (MB۷۵۲) به صورت ۶ فایل MB ۱۲۸ تبدیل خواهد شد.

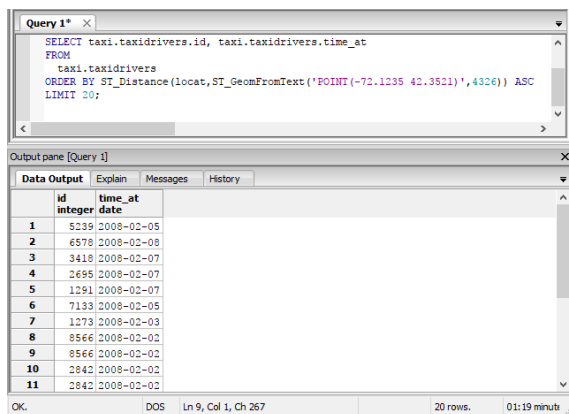
ابتدا در حالت بدون استفاده از ساختار ایندکس کتابخانه های مورد نظر را به برنامه اضافه کرده سپس نقطه ای که می خواهیم حول آن کوئری KNN را اجرا کنیم با استفاده از تابع createpoint تعریف می کنیم. در مرحله ی بعد RDD را با استفاده از فایل حاوی نقاط ایجاد می کنیم. عدد ۲ در ورودی این تابع نشان دهنده ی ستون شروع شونده ی داده های مکانی می باشد. در نهایت با استفاده از تابع SpatialKNNQuery از داده های موجود کوئری می گیریم. عدد ۵ در ورودی این تابع نشان دهنده ی تعداد K در این پرس و جوی می باشد. همچنین

۳-۴ - PostGIS

ساخته شده در ایندکس R-Tree در postgis می‌باشد. طبق توضیح ساختار R-Tree، در این نوع ایندکس از مرزها برای تعیین محدوده استفاده می‌کند که به همین دلیل نزدیکترین خودروها بدون تعیین هیچ فاصله دقیقی، به دست می‌آید. برای پیاده سازی KNN با ایندکس R-Tree ابتدا لازم است با دو عملگر فاصله ای که در مرتب سازی استفاده می‌شوند آشنا شد.

<>: عملگری است که فاصله دو بعدی دو geometry را بر می‌گرداند. این عملگر همراه با Order By می‌آید و کمترین/بیشترین فواصل را بر می‌گرداند. در واقع این عملگر از محاسبه تک تک نقاط موجود با نقطه داده شده و سپس مرتب کردن فواصل و برگرداندن بیشترین/کمترین فواصل جلوگیری می‌نماید. این فاصله همان فاصله مراکز دو داده geometry می‌باشد.

<#>: فاصله بین لبه های جعبه های در نظر گرفته شده در ساختار ایندکس را بر می‌گرداند. با استفاده از این دو عملگر می‌توان ایندکس مکانی مناسب را پیاده کرد و پرس و جوی مورد نظر را بر اساس آن پیاده سازی کرد. بدین ترتیب، زمان اجرای پرس و جوی از مرتبه $O(N)$ به مرتبه $O(\log^2 b)$ کاهش پیدا می‌کند. (b تعداد مرزهای در نظر گرفته شده می‌باشد).



The screenshot shows a PostgreSQL query window with the following SQL query:

```
SELECT taxi.taxidriversons.id, taxi.taxidriversons.time_at
FROM
  taxi.taxidriversons
ORDER BY ST_Distance(Locat, ST_GeomFromText('POINT(-72.1235 42.3521)', 4326)) ASC
LIMIT 20;
```

The output pane shows the following results:

id	time_at
integer	date
1	5239 2008-02-05
2	6578 2008-02-08
3	3418 2008-02-07
4	2695 2008-02-07
5	1291 2008-02-07
6	7133 2008-02-05
7	1273 2008-02-03
8	8566 2008-02-02
9	8566 2008-02-02
10	2842 2008-02-02
11	2842 2008-02-02

شکل ۱۲ Simple KNN

برای پیاده سازی KNN با استفاده از ایندکس در postgis، ایندکس هایی بجز R-Tree مانند GiST، B-Tree و BRIN و ... وجود دارد. که R-Tree و GiST در بخش های قبلی مورد بررسی قرار گرفته شد.

اگر پیدا کردن نزدیکترین خودروهای موجود به جایگاه خاصی مورد نظر باشد، کافی است پرس و جویی مبتنی بر الگوریتم KNN را روی دیتابیس مورد نظر اعمال کرد که بدین منظور راهکار های مختلفی وجود دارد، که در ادامه به بررسی آن ها در PostGIS پرداخته می‌شود.

ساده ترین راهکار این است که جدول یا جداول ترکیبی^{۱۳} بر اساس فاصله تعریف شده در پرس و جوی مرتب شود و K تا اولین رکورد ها را به عنوان نتیجه پرس و جوی برگردانده شود. این الگوریتم در شکل ۱۲ قابل مشاهده است. اما باید توجه کرد در این روش، اگر تعداد رکورد های جدول یا جداول های مورد نظر، زیاد باشند یا اینکه K را تعداد بالایی در نظر گرفت، همین پرس و جوی ساده ممکن است مدت زمان زیادی طول بکشد.

راهکار بعدی که برای پیاده سازی KNN می‌توان در نظر گرفت روش محدودیت ایندکس^{۱۴} می‌باشد که در آن یک عدد شاخص یا به اصطلاح عدد جادویی^{۱۵} نیاز است که با استفاده از آن، جعبه ای در محدوده ای که در پرس و جوی بیان شده است، ساخته می‌شود و خودروهای قرار گرفته در آن جعبه به عنوان نتایج پرس و جوی نشان داده می‌شود. این روش نیز در شکل ۱۳ آورده شده است. اما این روش نیز بدون نقص نیست و اگر رکوردی وجود نداشته باشد که در محدوده جعبه ساخته شده با عدد جادویی مورد نظر قرار گرفته باشد، هیچ نتیجه ای نشان داده نمی‌شود در حالیکه ما قصد داریم که نزدیکترین خودروها را پیدا کنیم هر چند نزدیکترین خودرو در فاصله بسیار زیادی از آن باشد.

اما روش اصلی که برای پیاده سازی الگوریتم KNN در postgis استفاده می‌شود، ارزیابی فاصله در محدوده مرز های

^{۱۳} Joined Tables

^{۱۴} Index constraint

^{۱۵} Magic number

شکل ۱۵ بررسی شباهت مسیر دو ماشین با شناسه ۱۰۰ و ۵۰۰

۶- ارزیابی و نتایج

۶-۱- دیتاست

یکی از مهم‌ترین و تاثیرگذارترین پارامترهای مقایسه عملکرد پایگاه‌داده‌ها برای پرس‌وجوهای مختلف انتخاب پایگاه‌داده‌های مناسب است.

پس از بررسی دیتاست‌های مختلف در نهایت یک دیتاست با حجم حدوداً ۸۰۰ مگابایت از مسیریابی عبور ۱۰۳۵۷ تاکسی در بازه ۲ الی ۸ فوریه سال ۲۰۰۸ میلادی که توسط شرکت مایکروسافت جمع‌آوری و منتشر شده است، به عنوان دیتاست تست انتخاب شد.

در این دیتاست حدوداً ۱۵ میلیون رکورد، در مجموع ۹ میلیون کیلومتر مسیر پیموده شده، متوسط زمان نمونه‌برداری ۱۷۷ ثانیه و متوسط فاصله بین هر دو نمونه‌برداری ۶۲۳ متر بوده است.

این دیتاست در هر رکورد شامل شناسه تاکسی، زمان نمونه‌برداری و محل نمونه‌برداری بوده است. با خصوصیات ذکرشده برای این دیتاست پرس‌وجوهای نظیر KNN و شباهت‌سنجی مسیرها برای این دیتاست پرس‌وجوهای معقولی به نظر می‌آید.

۶-۲- مشخصات کلاستر

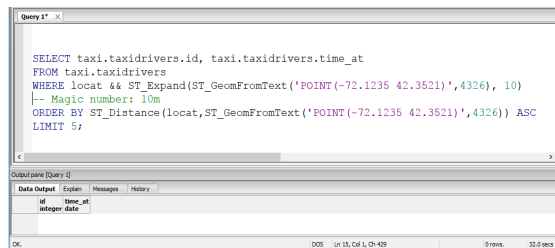
مشخصات کامپیوترهای استفاده شده در این آزمایش به شرح زیر می‌باشد:

- CPU: تک هسته ای GH۲,۴
- RAM: ۲ گیگابایت
- HDD: ۲۰ GB SSD

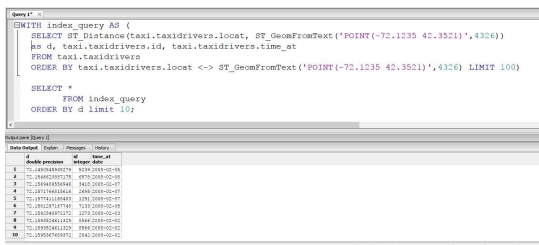
برای اجرای پرس و جو های PostGIS و MemSQL از یک کامپیوتر با مشخصات فوق استفاده شده است و برای اجرای پرس و جو ها بر روی بستر GeoSpark از یک کلاستر با دو عدد WorkerNode با مشخصات فوق استفاده شده است.

۶-۳- نمودار ها و مقایسه ها

پس از اجرای پرس و جو ها بر روی بستر GeoSpark و محاسبه ی زمان اجرای آنها نتایج زیر در به دست آمد. در این جدول زمان اجرای هر RDD، حجم آن و محل اجرای آن



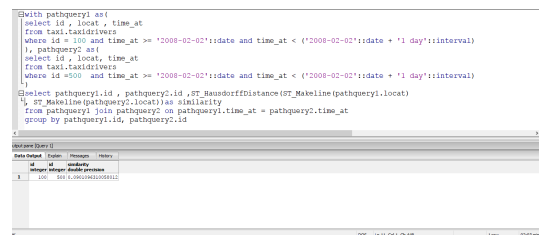
شکل ۱۳ Index Constraint KNN



شکل ۱۴ Index Based KNN

۵- پرس و جوی شباهت سنجی مسیر ها

در پایگاه داده های مکانی، یکی از مواردی که می‌تواند مورد پرس و جو قرار گیرد، مسیرهای مختلف و میزان شباهت مسیرها می‌باشد. راهکاری که در پایگاه داده های مختلف از جمله PostGIS به کار برده می‌شود، این است که ابتدا مسیرهای مورد نظر با استفاده از داده های geometry چون خط، چند ضلعی و ... ساخته شود و سپس به نحوی فاصله نقطه به نقطه این دو مسیر محاسبه شود و در نهایت فاصله کلی این دو مسیر و در واقع میزان اختلاف دو مسیر برگردانده شود. برای این منظور در PostGIS از تابع ST_HausdorffDistance که در بخش قبل، توضیح داده شد استفاده می‌شود. برای مثال در یک پایگاه داده حاوی اطلاعات مکانی ماشین‌های مختلف در زمان‌های مختلف، ابتدا مسیرهای دو ماشین را در یک روز به دست آورده و سپس با استفاده از تابع گفته شده، اختلاف دو مسیر بر حسب عدد اعشاری بر می‌گرداند. نمونه این پرس و جو و نتایج حاصل از آن را می‌توان در شکل ۱۵ مشاهده نمود.



شکل ۱۶ مقایسه ی زمان اجرای پرس و جوی KNN بدون استفاده از ایندکس

مشخص شده است. می توان مشاهده کرد که بیشترین زمان اجرا برای یک RDD در حدود ۱۲ ثانیه بوده است در نتیجه با افزایش حجم داده ها و اضافه شدن Worker node ها به راحتی می توان scalable بودن این روش را مشاهده کرد.

۷- نتیجه گیری

در این گزارش ابتدا به معرفی پایگاه داده ها و مدل های ایندکس مورد استفاده پرداخته شد. سپس الگوریتم KNN با مثالی از پایگاه داده تاکسی ها نشان داده شد و نکاتی در حین پیاده سازی به تفکیک پایگاه داده های مختلف بیان شد. پس از آن الگوریتم های شباهت سنجی مسیر ها و نکات پیاده سازی آن در پایگاه داده های مختلف، بیان شد و در آخر هر کدام از الگوریتم های KNN و شباهت سنجی مسیریها با استفاده از ایندکس و بدون ایندکس در سه نوع پایگاه داده رابطه ای (PostgreSQL)، رابطه ای و درون حافظه ای (MemSQL) و غیررابطه ای و درون حافظه ای (GeoSpark) مورد ارزیابی و مقایسه قرار گرفت که MemSQL به عنوان بهترین پایگاه داده در دیتاست مورد استفاده، شناخته شد. البته ذکر این نکته حائز اهمیت است که در صورت نیاز به فضایی برای پردازش داده های بزرگ که بر روی یک کامپیوتر قابل ذخیره سازی نیست بهترین روش استفاده از الگوریتم های مقیاس پذیر^{۱۶} است که بر روی بستر spark پیاده سازی شده اند.

مراجع

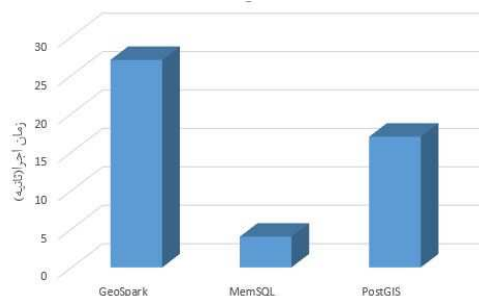
- 1] I. Team, "IBM Informix R-Tree Index User's Guide," IBM, 2008.
- 2] Antonin Guttman, Michael Stonebraker, "R-TREES: A DYNAMIC INDEX STRUCTURE FOR SPATIAL SEARCHING," College of Engineering University of California, Berkeley, 1983.

شماره	محل اجرا	زمان اجرا(S)	حجم(MB)
۱	Worker 2	12	128.1
۲	Worker 1	۱۲	128.1
۳	Worker 2	۷	128.1
۴	Worker 1	۷	128.1
۵	Worker 2	۶	128.1
۶	Worker 1	۶	112.2

در مجموع با استفاده از جدول زیر می توان نتیجه گرفت که هر کدام از worker node ها زمانی برابر با ۲۷ ثانیه کار کرده و حجم داده ی پردازش شده توسط هر کدام در این جدول نمایش داده شده است.

شماره	Node	زمان اجرا(S)	حجم(MB)
۱	Worker 1	27	384.2
۲	Worker 2	27	368.3

زمان اجرای پرس و جو ها به صورت کلی در نمودار شکل ۱۶ قابل مشاهده است. این نمودار زمان اجرای پرس و جوی KNN بدون استفاده از ایندکس را بر روی بستر های مختلف با یکدیگر مقایسه می کند. همانطور که مشاهده می شود برای دیتاست استفاده شده بهترین نتیجه با استفاده از پایگاه داده ی MemSQL به دست آمده است که زمانی برابر با ۴ ثانیه داشته است.



^{۱۶} scalable

- P. team, "PostGIS 2.1.10dev Manual," Postgis, 3] 2015.
- "MemSQL," [Online]. Available: 4] <https://docs.memsql.com>.
- M. Zaharia, M. Chowdhury, T. Das, A. Dave, 5] J. Ma and M. Justin, "Resilient distributed datasets: A fault-tolerant abstraction for in-memory cluster computing," in *Proceedings of the 9th USENIX conference on Networked Systems Design and Implementation*, 2012.
- N. Roussopoulos, S. Kelley and F. Vincent, 6] "Nearest neighbor queries," in *ACM sigmod record*, 1995.
- "GeoSpark," [Online]. Available: 7] <http://geospark.datasyslab.org/>.