



جامعة القاهرة
كلية الحاسبات والمعلومات
قسم علوم الحاسب



الفهرسة لتحسين تحليل البيانات الضخمة

إعداد
حسين شحاتة عبد العزيز

رسالة مقدمة إلى كلية الحاسبات والمعلومات

جامعة القاهرة

كجزء من متطلبات الحصول على درجة الماجستير

في علوم الحاسب

تحت إشراف

د/ محمد حلمي خفاجي

أ.د/ فاطمة عبدالستار حسن عمارة

قسم علوم الحاسب

قسم علوم الحاسب

كلية الحاسبات والمعلومات

كلية الحاسبات والمعلومات

جامعة الفيوم

جامعة القاهرة

كلية الحاسبات والمعلومات

جامعة القاهرة

جمهورية مصر العربية

يناير ٢٠١٦

الملخص العربى

تعد معالجة البيانات الضخمة واحدة من أعقد العمليات الحسابية حاليا نظرا للنمو الضخم فى حجم البيانات وابعادها وكذلك انماط البيانات الغير مهيكلة. الدراسات الحديثة تتوقع ازدياد مضطرد فى حجم البيانات يصل الى zettabytes. هذه البيانات يتم انتاجها عن طريق اجهزة الاستشعار , شبكات التواصل الاجتماعى وكذلك سجلات البيانات الناتجة عن مكينات المصانع.

تعد قواعد البيانات العلائقية غير مناسبة لتحليل تلك البيانات الضخمة وكذلك تطبيق الفهرسة او الارتباطات العلائقية عليها نظرا لانها تؤثر بالسلب على الوقت اللازم لتحميل ومعالجة تلك البيانات وكذلك استرجاعها.

اصبحت MapReduce من اهم الانظمة البرمجية الفعالة فى معالجة وتحليل البيانات الضخمة وعلى الجانب الاخر Hadoop يعد من اهم البرمجيات المبنية على طريقة Map/Reduce لتحليل ومعالجة البيانات الضخمة. وتعد HIVE من البرمجيات التى تعمل كمحرك قاعدة بيانات والتى تعمل بدون روابط علائقية او فهرسة للبيانات والتى بنيت كمترجم يحول جملة SQL الى مهمة MapReduce يمكنها ان تنفذ بواسطة Hadoop. لكن مراحل تنفيذ استعلامات الربط (الشبكية) من اعقد المراحل التى يقوم بها HIVE حيث انه يستهلك الوقت وكذلك مساحات تخزين مؤقتة وكذلك ذاكرة عشوائية كبيرة خاصة ب JAVA Heap والتى قد تؤدى الى حدوث تدفق زائد فى الذاكرة المحجوزة اثناء التشغيل مما يؤدى الى اعادة تشغيل المهمة مرة اخرى.

تعد قواعد البيانات النجمية من اكثر انماط قواعد البيانات التى تحتاج الى استعلامات ربط من اجل جمع المعلومات واستخراج التقارير لصانعى القرار مما يجعلها بطيئة نسبيا فى حالة استخدام HIVE كمحرك قاعدة بيانات.

يوجد العديد من الطرق التى حاولت اضافة الفهرسة الى نظام الملفات HDFS عن طريق اضافة الفهرس بطريقة ثابتة او متغيرة او اثناء عملية التشغيل او قبلها حتى تتمكن من تحسين الوقت اللازم لجمع المعلومات المطلوبة من استعلامات الربط ولكنها لا تزال تعاني من البطئ خاصة فى

حالة استعلامات الربط لأنها لا تزال تعتمد على مراحل تنفيذ HIVE لاستعلامات الربط والتي مازالت معقدة.

في هذه الرسالة سنقوم بعرض طريقتين مختلفتين كتعديل وتحسين لمراحل تنفيذ استعلامات الربط في HIVE من أجل تحسين كفاءة HIVE وكذلك تقليل الوقت المطلوب من أجل تنفيذ هذه الاستعلامات. وقد بنيت هاتين الطريقتين بأسلوب الفهرسة الثابت والمعد سلفا والذي نقوم ببناءه في مرحلة تحويل البيانات من نمط قاعدة البيانات النجمية الى نمط الجدول الكبير الذى يحوى كل بيانات قاعدة البيانات فى جدول واحد لتصبح فى حالة ربط مبدئى وكأننا قمنا بتنفيذ استعلام ربط لجمع كل البيانات فى جدول واحد كبير واضافة حقل جديد لها الجدول يمكننا من ربط الصفوف ببعضها ويمثل هذا الحقل (المؤشر او الفهرس) الذى يربط البيانات بعضها ببعض.

واحدة (Join Once) ويمكن استخدام هذه البيانات المجمعه اكثر من مرة بعد ذلك (Use Many) والتي قمنا بأختصارها الى (JOUM: Join Once Use Many).

يوفر دمج البيانات فى جدول واحد الوقت اللازم لعمليات الاستعلام وكذلك الذاكرة العشوائية المطلوبه والمساحة التخزينية المؤقتة التى يحتاجها HIVE اثنا تنفيذ استعلامات الربط ، ومن اهم مميزات الطريقتين التى نقوم بعرضهم انهما لا يقوم بتعديل فى الكود المصدري ل HIVE ولكنهما تضيف طبقات من الاكواد بعضها يسبق مراحل HIVE والاخرى تليها مما يجعلها ان تعمل مع الانظمة القديمة التى بنيت باستخدام HIVE دون اى تعديل مطلوب لهذه الانظمة مما يجعل هاتين الطريقتين تمثل اسلوب الفائز-الفائز لكل الاطراف.

تسمى الطريقة الاولى JOUM والتي تعتمد على جدول كبير واحد لكل البيانات والثانية Keys/Facts والتي تعتمد على تحويل الجدول الواحد الى جدولين احدهم يحمل المفاتيح الخاصة بالبيانات والاخرى يحمل الحقائق الخاصة بالبيانات المكررة على المفاتيح التى اضيفت من قبل فى الجدول الاول وتساعد عملية التقسيم هذه فى تحسين الوقت وتقليل المساحة التخزينية الثابتة للبيانات.

قمنا بتقييم الطريقتين السابقتين بإستخدام معيار البيانات TPC-H والتي بنيت بنمط قواعد البيانات النجمية والتي تعد من أكثر معايير البيانات المستخدمه ; وتظهر نتائج التقييم ان الطريقتين تعديا اداء

HIVE بنسبة (49.8%,53.1%) من نسبة الوقت اللازم ل HIVE وكذلك (27%,28%) من نسبة المساحة التخزينية المؤقتة اللازمة ل HIVE وكذلك (19%,26%) من نسبة المساحة التخزينية الثابتة للبيانات كحمل زائد وكذلك عدد اقل من عمليات مشاكل الذاكرة.

فى المجلد تعد الطريقتين مناسبتين لتحليل البيانات الضخمة على الرغم من وجود زيادة فى حجم البيانات عند تحويلها من نمط قاعدة البيانات النجمية الى الجدول الواحد الكبير وذلك بسبب حقل الفهرسة الذى تم اضافته الى البيانات وتعد هذه الزيادة مقبولة جدا اذا ما قورنت بنسبة التوفير فى المساحة المؤقتة والتي نحتاجها اثنا كل عملية التشغيل لاستعلام من استعلامات الربط.