

Spaltenbasierte Datenbanken - Ein Konzept zur Handhabung großer Datenmengen

Olaf Herden

Fakultät Technik Studiengang Informatik
Duale Hochschule Baden-Württemberg
Florianstr. 15
72160 Horb
o.herden@hb.dhbw-stuttgart.de

Abstract: Bei Anwendungen in der Land-, Forst und Ernährungswirtschaft, z.B. Rückverfolgungsdaten von Lebensmitteln und Produkten oder bei Prozess- und Zustandsdaten von Landmaschinen, ist in den letzten Jahren ein rasanter Anstieg des Datenwachstums zu beobachten. Die Speicherung und Verarbeitung dieser Daten findet typischerweise in Data Warehouses statt. Um auch bei wachsendem Datenvolumen eine effiziente Verarbeitung gewährleisten zu können, bietet die spaltenorientierte Speicherung in Kombination mit Kompression eine interessante Option, die in diesem Beitrag skizziert wird.

1 Einleitung

In den letzten 15 bis 20 Jahren haben sich Data Warehouse Systeme als Architektur für entscheidungsunterstützende Informationssysteme etabliert [BG09]. Anfangs standen klassisches Berichtswesen und OLAP (Online Analytical Processing) als Anwendungen im Vordergrund, im Laufe der Jahre wurden die Systeme um Planungs- und Vorhersagekomponenten erweitert, die Analyseintervalle haben sich immer weiter verkürzt und das zu verarbeitende Datenvolumen ist rasant angewachsen. Diese Anwendungsszenarien werden häufig mit dem Schlagwort Big Data bezeichnet, wobei als Charakteristika die sog. V-Eigenschaften vorliegen. Neben der großen Datenmenge (Volume) zählen hierzu auch die Notwendigkeit einer schnellen Verarbeitung der Daten (Velocity) und das Vorhandensein strukturierter und unstrukturierter Daten (Variety).

Eine Vielzahl von Konzepten und Technologien versucht diesen neuen Anforderungen zu begegnen, z.B. der Einsatz von massiv parallelen Systemen, physisches DB-Design in Form von Indexstrukturen, Partitionierungen oder Materialisierungen oder die In-Memory-Technologie, die den gesamten Datenbankinhalt im Hauptspeicher verwaltet. Ein weiterer, in diesem Beitrag betrachteter Ansatz, ist die spaltenorientierte Speicherung in Kombination mit Kompression.

Der Rest des Papers ist folgendermaßen gegliedert: Zunächst wird in Abschnitt 2 die Kompression in Datenbanken beschrieben, Abschnitt 3 stellt das Konzept der spaltenori-

entierten Speicherung vor. Abschnitt 4 zeigt anhand praktischer Untersuchungen, welche Kompressionsraten zu erwarten sind und von welchen Faktoren diese abhängen. Der Beitrag endet mit einer Zusammenfassung und einem Ausblick.

2 Kompression in Datenbanken

Kompression ist in der Informatik ein schon seit Jahrzehnten behandeltes Thema. Einen guten Überblick über den State-Of-The-Art gibt z.B. [Sa06]. Im Rahmen von Datenbanken kommen zum Einen nur verlustfreie Kompressionstechniken zum Einsatz, zum Anderen dürfen Einfüge-, Änderungs- und Leseoperationen sowie andere Datenbank-spezifische Aktionen wie z.B. Logging nicht zu stark beeinträchtigt werden.

Abbildung 1 zeigt exemplarisch drei im Umfeld von Datenbanken populäre Kompressionstechniken. Bei der ganz links dargestellten Wörterbuchkompression werden statt des Volltextes als Zeichenkette kurze binäre Codes abgespeichert. In der Mitte ist die Lauflängencodierung dargestellt, die sich bei Werten mit wenigen Zeichen und langen Sequenzen des gleichen Zeichens anbietet, anstelle der gesamten Sequenz wird jeweils die aufeinander folgende Anzahl an Werten gespeichert. Ganz rechts in Abbildung 1 ist schließlich die Differenzmethode zu sehen, die bei Sequenzen steigender (oder fallender) ganzzahliger Werte Anwendung findet. Anstelle der einzelnen Werte zu speichern, wird jeweils die Differenz zum Vorgänger angegeben. Zu jedem der drei Verfahren ist jeweils der Kompressionsfaktor (KF) als Quotient aus unkomprimierter und komprimierter Speicherung angegeben. Der eingesparte Speicherplatz ergibt sich, indem der Kehrwert des KF von 1 subtrahiert wird.

	Wörterbuchkompression	Lauflängencodierung	Differenzspeicherung
Unkomprimierte Speicherung	Kleidung Kleidung Frucht Foto Foto Kleidung Frucht Tiefkühlwaren Tiefkühlwaren Frucht Frucht Frucht Kleidung Foto Tiefkühlwaren Tiefkühlwaren Tiefkühlwaren	aaaaabbbcccccaaa	121355 121356 121358 121361 121362 121364 121366
Komprimierte Speicherung	Wörterbuch: Kleidung 00 Frucht 01 Tiefkühlwaren 10 Foto 11 Daten: 00 00 01 11 11 00 01 10 10 01 01 01 00 11 10 10 10	5a3b7c3a	-- 1 2 3 1 2 2
Kompressionsfaktor	Unkomprimierte Speicherung: 17 Einträge, 139 Zeichen á 1 Byte Komprimiert: Pro Eintrag 2 Bit KF = 32,7	Unkomprimiert: 18 Zeichen Komprimiert: 8 Zeichen KF = 2,25	Unkomprimiert: 7 Integer-Werte á 4 Byte Komprimiert: 7 Byte-Werte KF = 4
Speicherplatzersparnis	97%	56%	75%

Abbildung 1: Kompressionsverfahren in Datenbanken¹

¹ Bei dieser Berechnung wird der Platzbedarf für das Wörterbuch vernachlässigt. Dies führt im Beispiel zu einer erheblichen Steigerung des KF (sonst 3,2). Bei einer entsprechend großen Anzahl an Einträgen ist diese Annahme aber durchaus realistisch.

4 Untersuchungen zu Kompressionsraten

Um Aussagen über erreichbare Kompressionsraten treffen zu können, wurden mit Infobright einige Untersuchungen durchgeführt [He11, HH11]. Dabei haben wir festgestellt, dass die Kompressionsrate wesentlich vom Datentyp, der Anzahl unterschiedlicher Werte einer Spalte und der Anzahl von NULL-Werten abhängt. Abbildung 3 zeigt exemplarisch die gemessenen Werte bei unterschiedlicher Anzahl von Werten in einer Spalte. Insgesamt kann festgehalten werden, dass Kompressionsfaktoren von 8 bis 10 durchaus realistisch sind, in Einzelfällen wurden auch deutlich höhere Raten gemessen.

Anzahl unterschiedlicher Wert pro Spalte	2	3	5	10	20	50	100	500
Volumen	14.585.963	23.781.491	31.226.782	43.644.058	56.146.307	72.675.515	87.024.080	116.046.409
Faktor	36,8	22,5	17,2	12,3	9,6	7,4	6,2	4,6

Jeweils 10.000.000 Datensätze, Rohvolumen 511MB.

Abbildung 3: Kompressionsfaktor in Abhängigkeit verschiedener Werte

5 Zusammenfassung und Ausblick

In diesem Beitrag wurden spaltenorientierte Datenbanken als ein möglicher Ansatz zum Beherrschen großer Datenmengen in der Datenanalyse vorgestellt. Nach einer Einführung und Motivation wurden im Datenbankumfeld relevante Kompressionsverfahren vorgestellt sowie zeilen- und spaltenorientierte Speicherung beschrieben und gegeneinander abgegrenzt. Die Resultate erzielbarer Kompressionsraten einiger Untersuchungen mit dem Open Source System Infobright wurden vorgestellt. Diese Untersuchungen sollen in Zukunft vertieft und in einem Benchmark zusammengefasst werden.

Literaturverzeichnis

- [ABH09] Abadi, D.J., Boncz, P.A., Harizopoulos, S.: Column oriented Database Systems. PVLDB 2(2), 2009; S. 1664-1665.
- [BG09] Bauer, A.; Günzel, H.: Data-Warehouse-Systeme: Architektur, Entwicklung, Anwendung. dpunkt-Verlag, Heidelberg, 2009.
- [Ex12] <http://www.exasol.com/>, Letzter Abruf am 6.11.2012.
- [He11] Herden, O.: MySQL Engine Infobright: Speicherplatz sparen und schnellere Anfragen. Proceedings DOAG-Konferenz, Nürnberg, 2011.
- [HH11] Herden, O., Haller, T.: Das spaltenorientierte MySQL-Plugin Infobright als Kern einer Open Source basierten Data-Warehouse-Infrastruktur. Proceedings ISOS-Workshop, GI-Jahrestagung, Berlin, 2011.
- [In12] <http://www.infobright.org/>, Letzter Abruf am 6.11.2012.
- [Mo12] <http://www.monetdb.org/>, Letzter Abruf am 6.11.2012.
- [Sa06] Salomon, D.: Data Compression: The Complete Reference. Springer, Berlin, 2006.
- [St05] Stonebraker, M. u.a.: C-Store: A Column-oriented DBMS. Proceedings of the 31st VLDB Conference, Trondheim (Norwegen), 2005; S. 553-564.
- [Sy12] <http://www.sybase.de/>, Letzter Abruf am 6.11.2012.