

Analytik mit R auf Oracle Exadata und OBIEE

Matthias Fuchs
ISE Information Systems Engineering GmbH
Nürnberg

Schlüsselworte

ISE, Oracle R Enterprise, Engineered System, Exadata, Datamining, CRAN, Statistik, SQL, Storage Server, Database Machine, Exalytics.

Einleitung

Datamining auf großen Datenmengen spielt eine immer größere Rolle. Oracle bietet mit der Statistik Sprache R eine Möglichkeit, Analytik in der Datenbank und Big Data Umfeld durchzuführen. Aufgrund von Praxisbeispiel wird gezeigt, wie Analysen innerhalb der Datenbanken (Exadata) implementiert und aus der Oracle Business Intelligence (OBIEE - Exalytics) verwendet werden können. Der Einsatz von Exadata und Exalytics bringt hier große Vorteile, die Beispiele können aber auch auf jedem nicht Engineered System mit Oracle DB und Oracle Business Intelligence (OBIEE) ausgeführt werden

„R“

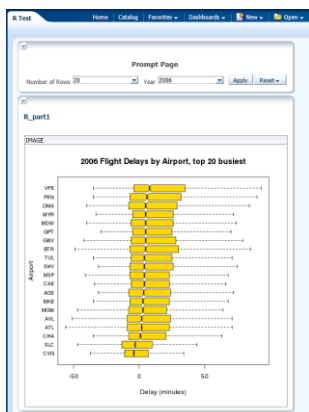
R ist eine freie Programmiersprache für statistisches Rechnen und statistische Grafiken. R gilt zunehmend als die statistische Standardsprache, sowohl im kommerziellen als auch im wissenschaftlichen Bereich. Durch den modularen Aufbau und die große Vielfalt von Erweiterungen (Paketen) bietet die Sprache viele Einsatzmöglichkeiten in der Statistik. Ob lineare oder nichtlineare Modellierung, Zeitreihenanalyse oder Clusteranalyse mit „R“ können fast alle Analysen durchgeführt werden.

„R“ und Big Data

Immer kürzere Produktlebenszyklen, der Trend zur Individualisierung sowie die fortschreitende Digitalisierung nahezu aller Geschäftsbereiche erhöhen die Menge der vorhandenen Daten und gleichzeitig die Notwendigkeit, intelligent mit dem Rohstoff Daten umzugehen. Die zu analysierenden Daten sind meist strukturiert in einer Datenbank abgelegt. Erreichen die Datenmengen mehrere Terrabyte, man kann von Big Data sprechen, kommen oft Oracle Datenbanken zum Einsatz.

Eine Kombination aus Oracle Datenbank und „R“ zur Analyse von strukturierten Daten ist daher eine Schlussfolgerung. Genau dieser Ansatz soll im Folgenden beschrieben werden.

„R“ und Oracle BI



Mit der Oracle Business Intelligence Enterprise Edition (OBIEE) ist es möglich R Berechnungen direkt auf den Analyse Daten oder über SQL in der Datenbank auszuführen. Ebenso können alle Arten von Visualisierungen aus „R“ direkt in einem Dashboard dargestellt werden. Somit werden die Möglichkeiten der Analyse deutlich erhöht.

Die Grundlage: Datamining in der Oracle Datenbank

Die Analyse von Daten umfasst mehrere Schritte. Die meiste Zeit geht vor der eigentlichen Analyse bei der Datenaufbereitung verloren. Es sind Exporte und Konvertierungen der Rohdaten durchzuführen. Diese werden dann wiederum auf weitere Systeme kopiert, um mit separaten Analysewerkzeugen direkten Zugriff zu haben. Während dieses Ablaufes geht viel Zeit verloren. Zusätzlich sind weitere Hardwareressourcen erforderlich. Diese Schritte entfallen, wenn die Daten an Ort und Stelle, in der Datenbank, verarbeitet werden. Zusätzlich greifen vorhandene Security und Compliance Richtlinien in der Datenbank und müssen nicht auf anderen Systemen repliziert werden.

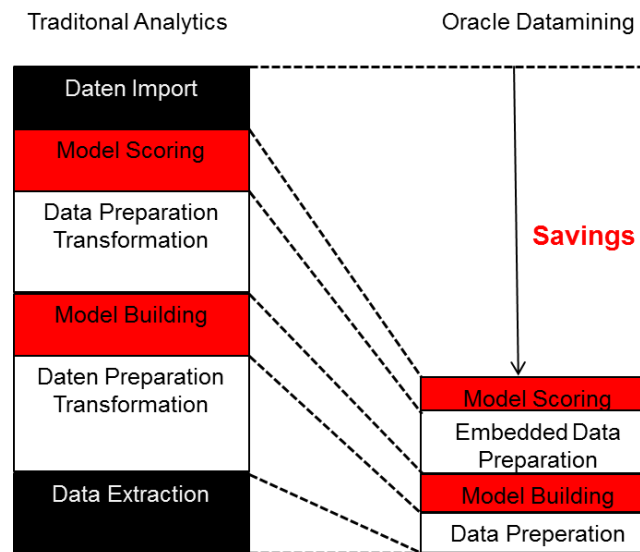


Abb. 1: Optimierungen beim in Oracle „R“ Database Datamining

Die erzielten Ergebnisse liegen ebenfalls wieder in der Datenbank und müssen nicht aufwendig importiert werden.

Zusätzlich zum einfacheren Datenhandling kommen Performancesteigerungen beim Analysieren der Daten. Je nach verwendetem Algorithmus ist mit einer deutlicher Beschleunigung beim Model Scoring oder Model Building zu rechnen. Die Verwendung von etablierten Standards für die Rechteverwaltung und Zugriffssteuerung, brauchen ebenfalls nicht verändert werden bzw. sind bereits vorhanden.

Die Erweiterung: Datamining mit Oracle „R“ Enterprise

Oracle hat die Verwendung von „R“ innerhalb der Datenbank transparent implementiert. Dies ist sowohl bei einer Installation auf Standard Hardware, als auch bei sogenannten Engineered Systems, wie die Oracle Exadata Database Maschine, möglich.

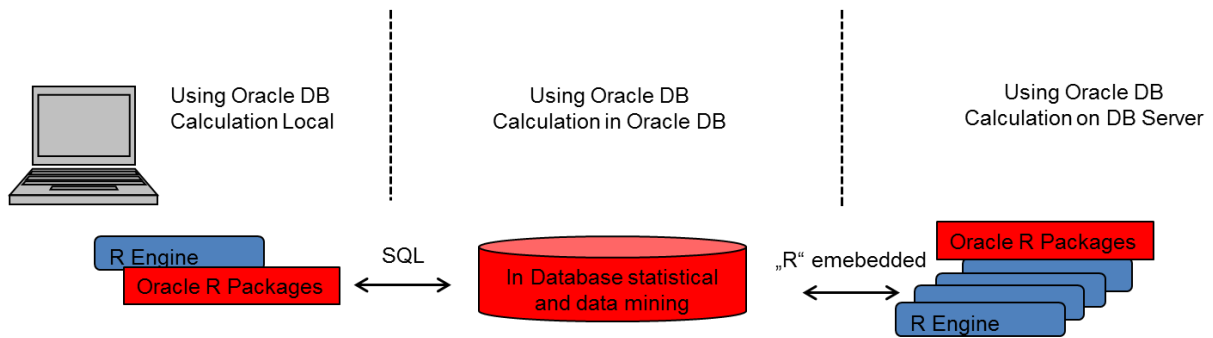


Abb. 1: Oracle „R“ Szenarios

„R“ Prozesse können auf drei Arten mit einer Oracle DB verarbeitet werden:

Die Berechnungen laufen auf einen unabhängigen Server bzw. Client und nur die Daten werden direkt aus der Oracle Datenbank geladen. Eine aufwendige Konvertierung der Daten in z.B. XML Datenstrukturen oder CSV Files entfällt. Es können alle R CRAN Pakete verwendet werden.

Alternativ kann man die Berechnungen auch direkt auf dem Datenbankserver starten. Dies erfolgt z.B. aus PL/SQL Prozeduren heraus. Der Vorteil besteht darin das keinerlei Netzwerkverkehr entsteht. Nur die Ergebnisse werden zum Client übertragen. Es können auch hier alle R Pakete verwendet werden.

Als letzte Möglichkeit Analysen mit R durchzuführen gibt es von Oracle speziell angepasste R Prozeduren. Dabei wurden die Pakete für die Ausführung in der Datenbank optimiert. Dadurch ergeben sich deutliche Performancesteigerungen gegenüber der Verwendung der „normalen“ R Pakete.

„R“ Oracle Optimierungen mit Oracle DB

Die aktuelle Version (1.3) der Oracle R Distribution enthält spezielle Pakete für Modell Erstellung und Scoring. Die Pakete setzen auf bekannten Oracle Datamining Paketen auf. Mit OREpredict ist ein hoch performantes Scoring, optimiert für Exadata, möglich. Mit OREdm sind Modellberechnungen mit Entscheidungsbäumen oder Linearer Regression möglich.

Kontaktadressen:

Matthias Fuchs
ISE Information Systems Engineering GmbH
Gewerbepark Hüll 4
D-91322 Gräfenberg

Telefon: +49 (0) 172-8288751
Fax: +49 (0) 9192-9929-22
E-Mail: matthias.fuchs@ise-informatik.de
Internet: www.ise-informatik.de

Oliver Bracht
eoda Heiko Miertzsch & Oliver Bracht GbR
Ludwig-Erhard-Straße 10
D-34131 Kassel

Telefon: +49 (0) 561/202724-40
Fax: +49 (0) 561/202724-30
E-Mail: info@eoda.de
Internet: <http://www.eoda.de>