



INFORMATION ASSET, LLC

Big Data Governance Tools:

*Insights from Information Asset's Enterprise
Data Management Lab*

by Sunil Soares
June 6, 2014





Introduction

Big Data Governance is part of a broader information governance program that formulates policy relating to Big Data. Big Data is often referred to in the context of the “three Vs:”

- **Volume (data at rest)**—Big Data is generally large.
- **Velocity (data in motion)**—Often time-sensitive, streaming data must be analyzed with millisecond response times to bolster real-time decisions.
- **Variety (data in many formats)**—Big Data includes structured, semi-structured, and unstructured data such as email, audio, video, clickstreams, log files, and biometrics.

Organizations need to establish policies to ensure the maximum success of their Big Data programs. In this research report, we will examine critical activities to kick-start a Big Data Governance program with a specific focus on software tools for Hadoop.

1. Define the Business Problem

Leading organizations turn to technologies such as Hadoop, NoSQL, and stream computing to solve business problems that cannot be addressed with traditional tools. The first step is to define a high-level business problem that can be addressed by Big Data. We use a three-dimensional framework for Big Data Governance as shown in Figure 1:

- *Big Data types*—Big Data governance needs a heightened focus on the data itself. We have classified Big Data into five distinct types: web and social media, machine-to-machine, big transaction data, biometrics, and human generated.
- *Information governance disciplines*—The traditional disciplines of information governance also apply to Big Data. These disciplines are organization, metadata, privacy, data quality, business process integration, master data integration, and information lifecycle management. In this research report, we examine how tools can support the execution of these disciplines for Big Data.
- *Industries and functions*—Big Data analytics are driven by uses cases that are specific to a given industry or function.

Here are a couple of examples that leverage this framework:

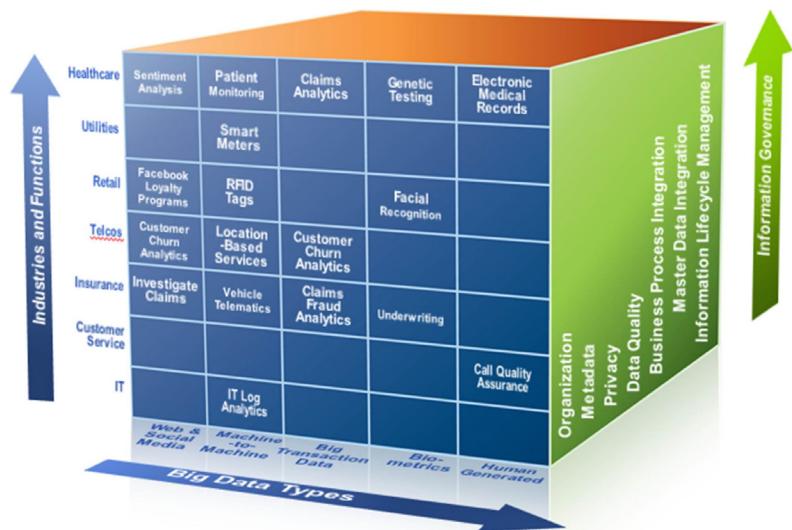


Figure 1: Big Data Governance framework.



INFORMATION ASSET, LLC

Example 1

Solution: Smart meters in utilities

Big Data type: Machine-to-machine

Disciplines: Privacy, information lifecycle management

Several utilities are rolling out smart meters to measure the consumption of water, gas, and electricity at regular intervals of an hour or less. These smart meters generate copious amounts of “interval” data that needs to be governed appropriately. Utilities need to safeguard the privacy of this interval data because it can potentially point to a subscriber’s household activities, as well as the comings and goings from his or her home. In addition, utilities need to establish policies for the archival and deletion of interval data to reduce storage costs.

Example 2

Solution: Sentiment analysis

Big Data type: Web and social media

Disciplines: Master data integration, data quality, privacy

The marketing department might want to use Twitter feeds for “sentiment analysis,” to determine what users are saying about the company and its products or services. The analytics team needs to determine if, for example, references to “@Acme” and “Acme” both refer to “Acme Corporation.” Integrating sentiment analysis with a customer’s profile can also be challenging. In addition to privacy issues, the Twitter handle might reveal the user name in only 50 to 60 percent of cases. Finally, marketing might need to answer the following question: “Do we really believe that Twitter sentiment analysis is representative if users are younger and more affluent than our typical customers?”

2. Conduct an Inventory of Data in Hadoop

Many Big Data teams start by creating a “data lake” in Hadoop to act as a landing zone for data across the organization and for Big Data analytics. Left unchecked, these data lakes can quickly become unmanageable for data scientists who might have to deal with multiple copies of the same data set. Because you cannot govern what you do not know exists, the first step to governing Big Data is to build an inventory.

In Figure 2, we show the Waterline Data Inventory for Hadoop by Waterline Data Science. The data scientist views the search results for Customer. In the left panel, the data scientist can also view the faceted search with results classified by categories. For example, based on the lineage discovered by Waterline, there are 12 files originating from SAP and three files that contain Midwestern states. On the right, the user can see profiling results and the data distribution for the fields in the selected file. In the example, the file is in CSV format without a header row, and Waterline Data Inventory used a combination of user tagging and automated data discovery to identify the meanings of the fields and to assign appropriate tags.

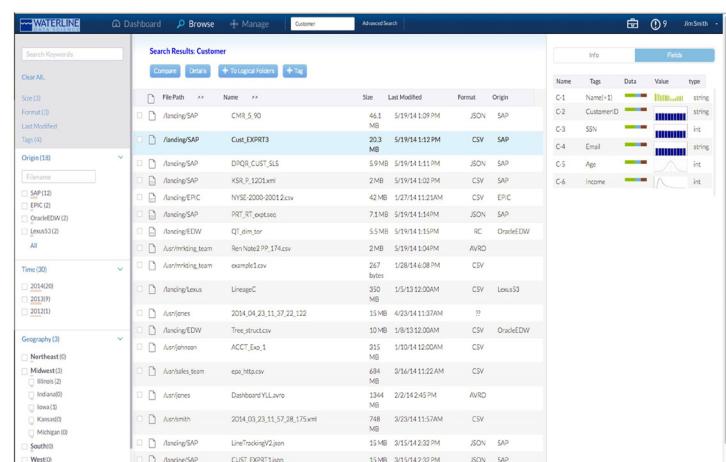


Figure 2: Hadoop Data Inventory file browser in Waterline Data Science.



3. Assign Ownership for Data in Hadoop

Data governance leads need to assign the ownership of data whether or not it resides in Hadoop.

As shown in Figure 3, a retailer uses Data Advantage Group MetaCenter to manage data ownership for all types of data. For example, the Finance Reporting application and Data Mart 1 are owned by Finance and Merchandising, respectively. On the other hand, Big Data in the form of Twitter and Facebook are primarily managed in Hadoop and are under the ownership of Marketing. Finally, Supply Chain owns radio-frequency identification (RFID) data. MetaCenter also contains the names of the Data Executive, Managing Data Steward, and Data Steward for each application or data type. The Marketing department is accountable to establish and enforce policies relating to the acceptable use of Twitter and Facebook data.

The screenshot shows the Data Advantage Group MetaCenter interface. The top navigation bar includes 'Welcome: Administrator', 'Environment: SQL SERVER', 'Date: Wed Apr 9 2014', 'View', 'Tools', 'Help', 'Logout', and 'My Tasks (0)'. A search bar at the top right says 'Data Repository Dashb...' and 'Enter search phrase here.' with an 'Exact Match' checkbox and a 'Manage Templates' link. The main area has a title 'Results by Category' and a tree view on the left with categories like Application Type, Name, Organization, Data Executive, and Data Steward. The 'Name' category is expanded, showing items like 'Data Mart1 (1)', 'Facebook (1)', 'Finance Reporting (1)', and 'Rfid (1)'. The 'Organization' category is also expanded, showing 'Marketing (2)', 'Finance (1)', 'Merchandising (1)', and 'Supply Chain (1)'. The 'Data Executive' category is expanded, showing 'Mary Jane (2)', 'Jack Murphy (1)', 'John Smith (1)', and 'Liz Shi (1)'. To the right, a table titled 'Search Results - 5 results (0.028 seconds)' lists the following data:

	Application Type	Name	Type	Organization	Data Executive	Managing Data Steward	Data Steward
1	MetaCenter	Finance Reporting	Data Repository	Finance	John Smith	Jane Smith	Mary Jane
2	MetaCenter	Data Mart1	Data Repository	Merchandising	Jack Murphy	Jill Smith	Tom Smith
3	MetaCenter	Twitter	Data Repository	Marketing	Mary Jane	Nancy Smith	Jean Hill
4	MetaCenter	Facebook	Data Repository	Marketing	Mary Jane	Nancy Smith	Jean Hill
5	MetaCenter	RFID	Data Repository	Supply Chain	Liz Shi	Helena O'Toole	Maya Danielle

Figure 3: Data ownership in Data Advantage Group MetaCenter.

Executive, Managing Data Steward, and Data Steward for each application or data type. The Marketing department is accountable to establish and enforce policies relating to the acceptable use of Twitter and Facebook data.

4. Provision a Semantic Layer for Analytics in Hadoop

Although Hadoop is a powerful analytics platform, it is meaningless without the appropriate semantic layer to provide meaning to key business terms. In the next few screenshots, we explore a simple MapReduce word count program that executes within Cloudera. The business definitions for the key terms are governed in Collibra Data Governance Center.

The screenshot shows a spreadsheet with columns 'A' and 'B'. The data consists of numbered definitions corresponding to specific sections of the Dodd-Frank Act:

1	Name definition
2	affiliate Section 2 of the Bank Holding Company ActFor purposes of this Act the term "affiliate" means any company that controls, is controlled by, or is under common control with another company (see definition)
3	agency Section 1 of the International Banking Act"Agency" means any office or any place of business of a foreign bank located in any State of the United States at which credit balances are maintained incidental to
4	bank Section 3 of the Federal Deposit Insurance Act.The term "bank"--(A) means any national bank and State bank, and any Federal branch and insured branch; and (B) includes any former savings association.
5	bank hold Section 2 of the Bank Holding Company Act(a)(1) Except as provided in paragraph (5) of this subsection, "bank holding company" means any company which has control over any bank or over any company
6	branch Section 1 of the International Banking Act "branch" means any office or any place of business of a foreign bank located in any State of the United States at which deposits are received.
7	commercial Section 602 of Dodd FrankA company is a "commercial firm" if the annual gross revenues derived by the company and all of its affiliates from activities that are financial in nature (as defined in section 4(k)
8	exposure Section 23 of the Federal Reserve ActAn insured depository institution's "exposure" to another depository institution means--A.all extensions of credit to the other depository institution, regardless of nat

Figure 4: Banking-related business terms relating to the Dodd-Frank legislation in the United States.

Figure 4 shows a list of banking terms related to the Dodd-Frank legislation in the United States.



INFORMATION ASSET, LLC

In Figure 5, we execute the MapReduce code on the .csv file that contains the Dodd-Frank terms.

In Figure 6, we view the results of the MapReduce word count program in Cloudera. The words affiliate and bank appear once each.

Finally, we can view the business definitions of these terms in Collibra Data Governance Center as shown in Figure 7. In this way, Collibra provides the semantic layer for Hadoop analytics in Cloudera.

```
[cloudera@localhost ~]$ hadoop jar wordcount.jar org.myorg.WordCount /user/cloudera/wordcount/input /user/cloudera/wordcount/output
14/04/07 06:10:27 INFO mapred.JobClient: Use SimpleOptionsParser for parsing the arguments. Applications should implement Tool for the same.
14/04/07 06:10:27 INFO mapred.JobClient: Total input paths to process : 3
14/04/07 06:10:25 INFO mapred.JobClient: Running job: job_201404070409.0005
14/04/07 06:10:26 INFO mapred.JobClient: map 0% reduce 0%
14/04/07 06:11:42 INFO mapred.JobClient: map 50% reduce 0%
14/04/07 06:12:08 INFO mapred.JobClient: map 100% reduce 0%
14/04/07 06:12:22 INFO mapred.JobClient: map 100% reduce 100%
14/04/07 06:12:29 INFO mapred.JobClient: job complete: job_201404070409.0005
14/04/07 06:12:29 INFO mapred.JobClient: Counters: 33
14/04/07 06:12:29 INFO mapred.JobClient: File System Counters
14/04/07 06:12:29 INFO mapred.JobClient: FILE: Number of bytes written=644
14/04/07 06:12:29 INFO mapred.JobClient: FILE: Number of read operations=0
14/04/07 06:12:29 INFO mapred.JobClient: FILE: Number of large read operations=0
14/04/07 06:12:29 INFO mapred.JobClient: FILE: Number of write operations=0
14/04/07 06:12:29 INFO mapred.JobClient: HDFS: Number of bytes read=2144
14/04/07 06:12:29 INFO mapred.JobClient: HDFS: Number of bytes written=541
14/04/07 06:12:29 INFO mapred.JobClient: HDFS: Number of read operations=9
14/04/07 06:12:29 INFO mapred.JobClient: HDFS: Number of large read operations=0
14/04/07 06:12:29 INFO mapred.JobClient: HDFS: Number of write operations=2
14/04/07 06:12:29 INFO mapred.JobClient: Job Counters
14/04/07 06:12:29 INFO mapred.JobClient: Launched map tasks=4
14/04/07 06:12:29 INFO mapred.JobClient: Launched reduce tasks=4
14/04/07 06:12:29 INFO mapred.JobClient: Data-local map tasks=4
14/04/07 06:12:29 INFO mapred.JobClient: Total time spent by all maps in occupied slots (ms)=176572
14/04/07 06:12:29 INFO mapred.JobClient: Total time spent by all reduces in occupied slots (ms)=19842
14/04/07 06:12:29 INFO mapred.JobClient: Total time spent by all maps waiting after reserving slots (ms)=0
14/04/07 06:12:29 INFO mapred.JobClient: Total time spent by all reduces waiting after reserving slots (ms)=0
14/04/07 06:12:29 INFO mapred.JobClient: Map input records=9
14/04/07 06:12:29 INFO mapred.JobClient: Map output records=10
14/04/07 06:12:29 INFO mapred.JobClient: Map output bytes=100
```

Figure 5: Executing MapReduce code on .csv file with Dodd-Frank business terms.

```
[cloudera@localhost ~]$ hadoop fs -cat /user/cloudera/wordcount/output/part-00000
**affiliate** 1
**bank** 1
**branch** 1
**exposure** 1
(5) 1
(1) 1
(9) 1
(9) 1
(see 1
(2 1
(15 1
(1956 1
(2 2
(23 1
(1 1
(4(R) 1
(692 1
(920 1
Act 3
Act**Agency** 1
Act(a)(1) 1
Act 1
Act.The 1
ActAn 1
ActFor 1
ActThe 1
Act 1
Bank 3
Banking 2
```

Figure 6: Viewing the results of the MapReduce word count program in Cloudera.

Dodd Frank Terms

Type: Glossary

Assets

+ Add

More...

Responsibilities

Snapshots

Tasks

Description

There is currently no description.

Default

Public. The default view to show all the assets in this vocabulary.

Hierarchy Export Discover

Save Filter

Name	Definition	Steward	Status	Type
affiliate	Section 2 of the Bank Holding C...	Candidate	Business Term	
agency	Section 1 of the International Ba...	Candidate	Business Term	
bank	Section 3 of the Federal Depositi...	Candidate	Business Term	
bank holding company	Section 2 of the Bank Holding C...	Candidate	Business Term	
branch	Section 1 of the International Ba...	Candidate	Business Term	
commercial firm	Section 602 of Dodd Frank co...	Candidate	Business Term	
exposure	Section 23 of the Federal Reser...	Candidate	Business Term	

Figure 7: Viewing definitions of business terms in Collibra Data Governance Center.



5. View the Lineage of Data In and Out of Hadoop

As Hadoop becomes mainstream, organizations will use the platform for mission-critical applications. This means that data governance teams will need to include Hadoop within their data lineage views.

In Figure 8, ASG-Rochade shows the detailed forward lineage of the **EMP_EXPENSE_FACT** table. The lineage report shows the **emp_expense_fact** Hive table as well as two sets of Sqoop ETL jobs.

Figure 9 shows how a user can use Data Advantage Group MetaCenter to browse the employee table in Hive.

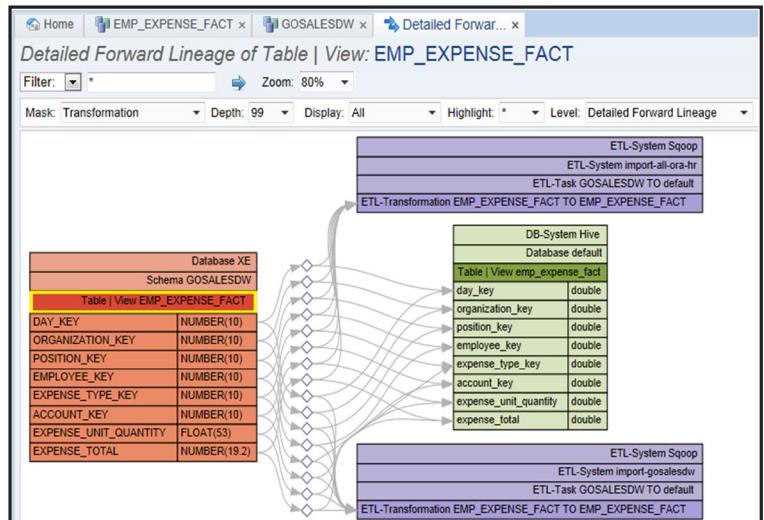


Figure 8: Detailed Forward Lineage of EMP_EXPENSE_FACT table in ASG-Rochade.

The screenshot shows the "employee" table details in the "employee" database of the "DEV_HADOOP_METASTORE" catalog. The table metadata includes:

Name	employee
Application	Hadoop
Type	Table
Location	DEV_HADOOP_METASTORE
Database Name	default
Table Name	employee
Table Type	
Owner	root
Creation Time	2013-09-25T17:36:16.000Z
Last Access Time	1970-01-01T00:00:00.000Z
Is Compressed	
SerDe Lib Name	
File System Path	hdfs://vm-hadoop-43:8020/user
Input Format	org.apache.hadoop.mapred.T
Output Format	org.apache.hadoop.hive.io.H
Partition Column Names	
Clustered By Column Names	
Sorted Column Names	

Figure 9: Browsing a Hive table in Data Advantage Group MetaCenter.



6. Manage Reference Data for Hadoop

Hadoop implementations also need high-quality reference data, as shown in the remainder of this section.

In Figure 10, the **GDP_Test.csv** file contains the Gross Domestic Product (GDP) in millions of U.S. dollars by country. This data is in unsorted format and will be used as the input file for a Pig script.

In Figure 11, we have created a Pig script to sort the countries in descending order of GDP. The script runs in Cloudera and takes the **GDP_Test.csv** file as an input file. The script places the result in a file called **output**.

Figure 12 shows the results of the Pig script in sorted format.

In Figure 13, we can see the reference data for country codes in Data Advantage Group MetaCenter. In the left panel, we can view a list of country codes. In the right panel, we can see that **FR** represents **France**.

GDP_Test.csv X	
USA,1,United States,16244699.00	
CN,2,China,8227103.00	
JPN,3,Japan,5961866.00	
IRN,21,"Iran, Islamic Rep.",552397.00	
CAN,11,Canada,1779635.00	
DEU,4,Germany,3428131.00	
GBR,5,United Kingdom,2475782.00	
ITA,9,Italy,2014679.00	
IND,10,India,1858740.00	
RUS,8,Russian Federation,2014775.00	
SWE,22,Sweden,523946.00	
ESP,13,Spain,1322955.00	
MEX,14,Mexico,1178126.00	
KOR,15,"Korea, Rep.",1129598.00	
IDN,16,Indonesia,87843.00	
COL,17,Colombia,77052.00	
NLD,18,Turkey,770555.00	
SAU,19,Saudi Arabia,711059.00	
CHE,20,Switzerland,631173.00	
BRA,7,Brazil,2252664.00	
SWE,22,Sweden,523946.00	
POL,24,Poland,489795.00	
BEL,25,Belgium,483862.00	
ARG,26,Argentina,475802.00	
AUT,27,Austria,394708.00	
PER,28,Peru,300000.00	

Figure 10: GDP_Test.csv file contains Gross Domestic Product by country.

```

data = LOAD '/user/cloudera/InputForPig/GDP_Test.csv' using PigStorage(',') AS (countrycode:chararray,ranking:int,economy_name:chararray);
describe data;
dump data;
outputdata = ORDER data BY ranking ASC;
describe outputdata;
dump outputdata;
store outputdata into '/user/cloudera/InputForPig/output/output';

```

Figure 11: Pig script to sort countries in descending order of GDP.

Actions	
View as binary	
Edit file	
Download	
View file location	
Refresh	
INFO	
Last modified	
April 14, 2014 3:38 a.m.	
User	
cloudera	
Group	
cloudera	
Size	
4 KB	
Mode	
100%44	

Warning: some binary data has been masked out with 'b6vfd'.

countrycode	economy_name	ranking
USA	United States	1.62446E7
CN	China	8227103.0
JPN	Japan	5961866.0
DEU	Germany	3428131.0
FRA	France	2014679.0
GBR	United Kingdom	2475782.0
BRA	Brazil	2252664.0
RUS	Russian Federation	2014775.0
ITA	Italy	2014078.0
IND	India	1858740.0
CAN	Canada	1779635.0
AUS	Australia	1532400.0
ESP	Spain	1322955.0
MEX	Mexico	1178126.0
KOR	"Korea"	1129598.0

Figure 12: Output of Pig script shows a sorted list of countries in descending order of GDP.

Application Search
ER
ES
ET
FI
FJ
FK
FM
FO
FR
GA
GB
GD
GE
GF
GH
GI
GL
GM
GN

FR

Name	Value
Application	Information Asset
Environment	Information Asset
Type	Code Value
Location	Information Asset / Country Codes / FR
Created By	Administrator
Created Date	2014-04-14 06:00
Last Updated By	Administrator
Last Updated Date	2014-04-21 11:06
Description	FRANCE

Figure 13: Data Advantage Group MetaCenter contains country codes relating to GDP.



INFORMATION ASSET, LLC

7. Profile Data Natively in Hadoop

Organizations also need to profile data natively in Hadoop without moving large volumes into structured data stores.

Figure 14 shows the Global IDs HDFS File Profiler. The screenshot shows a column analysis for the product data in semi-structured format in the **Inventory.txt** file. The profiler has uncovered 49 columns, including **FAMILY**, **PRODUCTGROUP**, **COMMODITYGROUPID**, **PRODUCT_TYPE**, and **SELL_SWAP_INDICATOR**.

In Figure 15, Talend Data Quality profiles data natively in Hive. By leveraging the Hadoop MapReduce distributed architecture, Talend Data Quality can cleanse, match, and de-duplicate large volumes of data without having to move it prior to processing.

Informatica also supports data profiling natively on Hadoop through Informatica Developer and the browser-based Informatica Analyst to better understand the data, identify data quality issues, collaborate on data flow specifications, and validate mapping, transformation, and rules logic.

In Figure 16, the data analyst can view unique values, null values, percent nulls, inferred data types, and other descriptive statistics to identify outliers and anomalies.

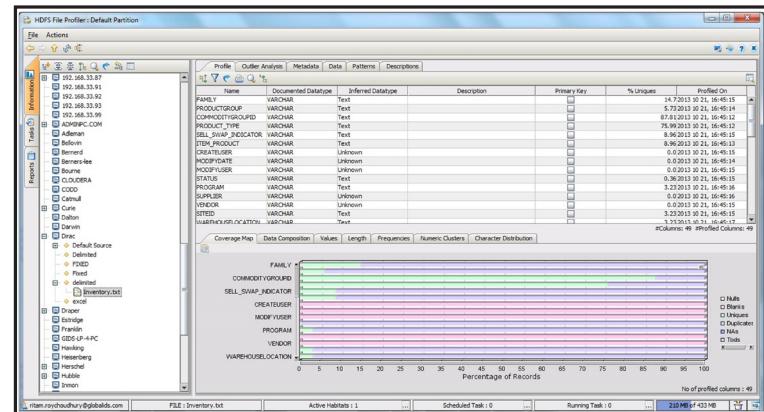


Figure 14: Global IDs HDFS File Profiler.

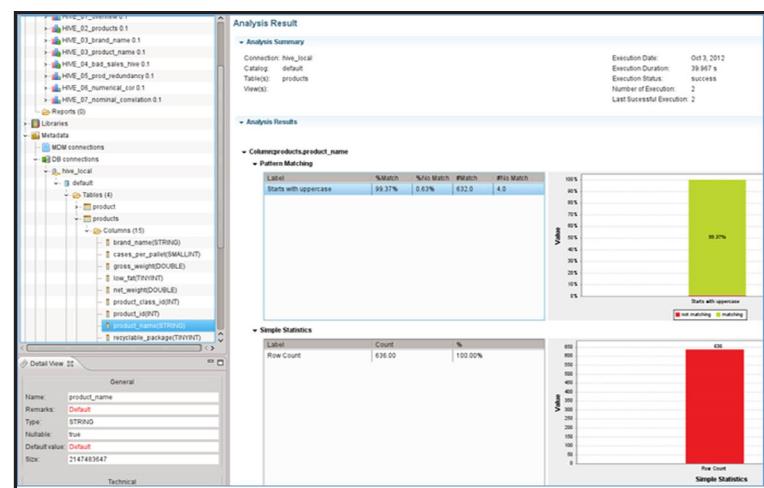


Figure 15: Talend Data Quality profiles data natively in Hive.

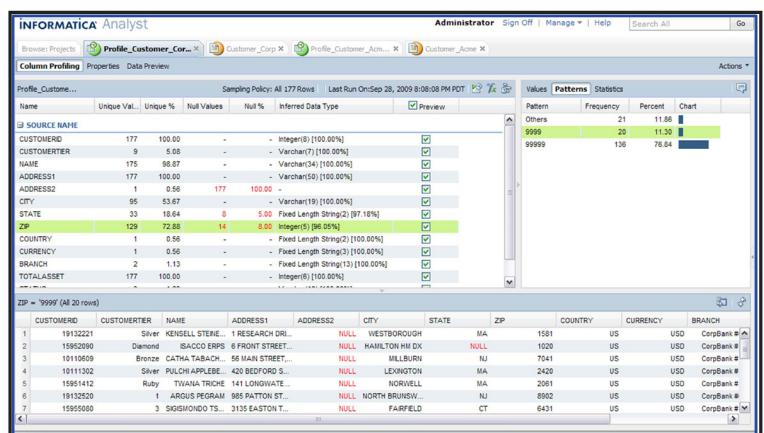


Figure 16: Informatica Data Quality Big Data Edition profiles data natively in Hadoop.



8. Discover Data Natively in Hadoop

Data profiling tools should also natively discover hidden data in Hadoop.

Figure 17 shows how Informatica Data Quality Big Data Edition can be set up to discover **IPAddress** and **StockSymbol** data within a **WebLogs** Logical Data Object in Hadoop.

As shown in Figure 18, Informatica Data Quality Big Data Edition shows the results of the data domain discovery for the **WebLogs** Logical Data Object. We can see that 100 percent of the values in the **IP_Address** column conform to the **IPAddress** pattern. The same holds true for the **Stock_Symbol** column relative to the **StockSymbol** pattern.

Name	Description	Data Domain Group
All Data Domains	Validates input for IP addr...	Undefined
<input checked="" type="checkbox"/> IPAddress	Validates input for IP addr...	Undefined
<input type="checkbox"/> Login	The login state of the user...	Undefined
<input checked="" type="checkbox"/> StockSymbol		Undefined

Figure 17: Informatica Data Quality Big Data Edition supports data domain discovery within Hadoop.

Name	% Data Conforming	% Null	Data Domain Groups	Documents
IP_Address	100.00	0.00	Undefined	string(16)
StockSymbol	100.00	0.00	Undefined	string(10)

Figure 18: Informatica Data Quality Big Data Edition shows the results of the data domain discovery exercise.



9. Execute Data Quality Rules Natively in Hadoop

Data quality tools should also execute rules natively in Hadoop.

In Figure 19, we show a sample web log file called **DailyWebLogDump.txt**.

In Figure 20, we can see the mappings in Informatica Data Quality Big Data Edition. The **DailyWebLogDump.txt** is read as a binary source file in the **Read** mapping on the left. The file is then parsed in the **POC_Session** transformation in the middle and the results passed to the Joiner mapping on the right.

Figure 21 provides further detail on how the binary file is parsed in Informatica Data Quality Big Data Edition. The parser conducts a text search for VISITOR_id. It then populates any associated text found in parentheses in the VISITOR_ID element in the VISITOR object in a predefined schema.

Figure 19: Sample web log file in semi-structured format.

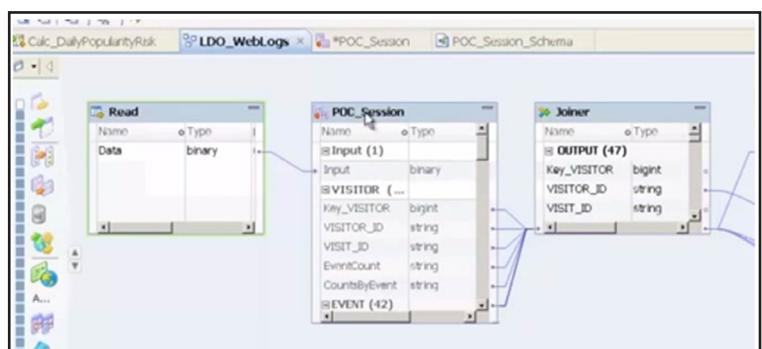


Figure 20: Mappings in Informatica Data Quality Big Data Edition.

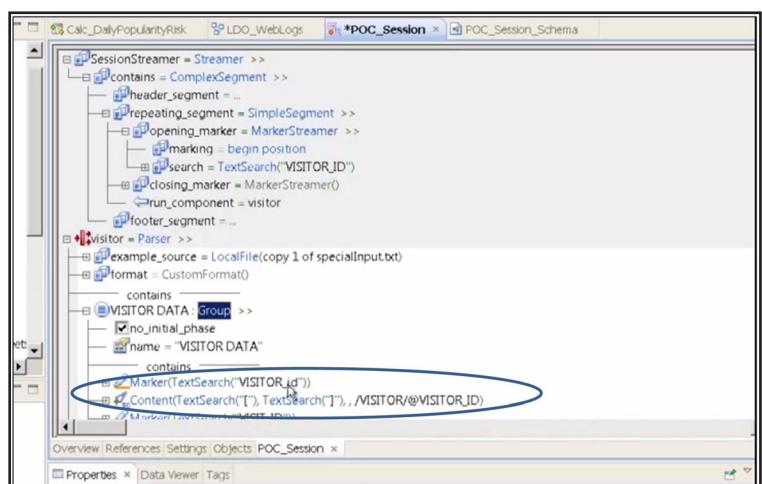


Figure 21: Complex data parsing with Informatica Data Quality Big Data Edition.



10. Integrate Master Data Management with Hadoop for Improved Performance

Organizations can also leverage the power of Hadoop to turbo charge their MDM hubs. For example, marketing organizations often need to match lists of prospects against internal records to remove any customers who have made do-not-call elections. These large data sets push the limits of existing computational resources when IT needs to match 200 million prospects against a database of 100 million customers and turn around the results to marketing in 24 hours. MDM vendors such as IBM and Informatica have integrated their MDM hubs with Hadoop.

BigInsights

Figure 22 shows Big Match from IBM. The company has embedded the IBM InfoSphere Master Data Management probabilistic matching engine within the IBM InfoSphere BigInsights Hadoop distribution.

11. Mask Sensitive Data in Hadoop

Big Data teams also need to mask sensitive data within their Hadoop data lakes. Many vendors, such as IBM and Informatica, support data masking functionality natively in Hadoop.

Figure 23 shows the batch interface in HDFS that uses IBM InfoSphere Optim Data Privacy Solution to mask sensitive data.

Big Data Matching Dashboard

Input table: person

Algorithm: mdmper

Search Type: Entity

Minimum Score: 98

Search

Entity Id	Score	First Name	Last Name	SSN	SSN Issuer	SIN	Category
387983757329571842	98	BARBARA	GONZALEZ	910-05-6610	655-42-0292	F	
387983757329571842	98	BARBARA	GONZALEZ	910-05-6610	655-42-0292	F	
387983757329571842	98	BARBARA	GONZALEZ	910-05-6610	654-52-0292	F	
387983757329571842	98	BARBARA	GONZALEZ	910-05-6610	654-52-0292	F	
387983757329571842	98	BARBARA	GONZALEZ	910-05-6610	655-42-0292	F	

Figure 22: Probabilistic matching engine within IBM InfoSphere BigInsights (“Big Match”).

Figure 23: IBM InfoSphere Optim Data Privacy Solution masks sensitive data natively in Hadoop.



12. Manage Workflows Associated with The Approval Process for Key Data Artifacts

Data governance tools should also manage the approval process to add, modify, and delete key data artifacts in Hadoop. For example, the data governance team manages a list of key business terms in Hive.

As shown in Figure 24, the team uses the Simple Approval workflow in Collibra Data Governance Center to manage the process to add a new business term. Once the term has been approved, Collibra initiates a HiveQL query to update the glossary in Hive. This process can also be used to approve other data artifacts in Hadoop, such as product lists and reference data.

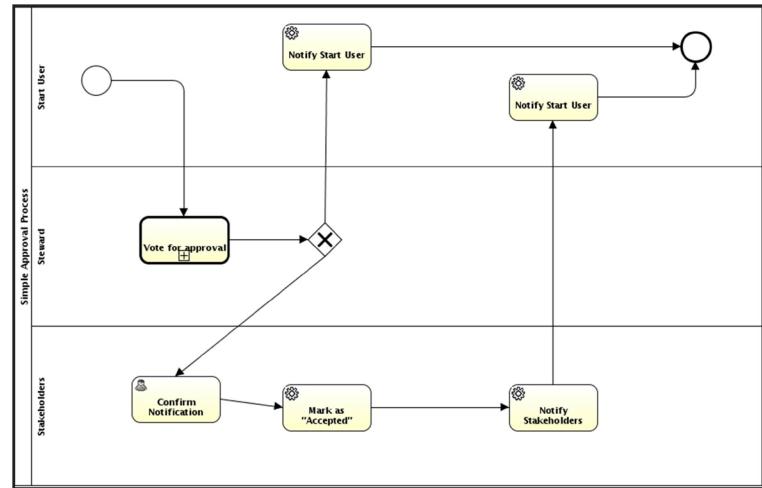


Figure 24: Simple approval workflow in Collibra Data Governance Center.

Glossary

Hadoop

Apache Hadoop is an open source software library that supports the distributed processing of large data sets across thousands of computers based on commodity hardware.

Hadoop Distributed File System (HDFS)

HDFS is a distributed file system that is designed to be highly fault tolerant and to run on low-cost hardware.

Hive

Apache Hive is a data warehousing infrastructure that sits on top of Hadoop. Hive provides a SQL-like interface called HiveQL to query large volumes of data in Hadoop. Because Hive insulates users from having to learn the intricacies of MapReduce programming in Java, it is a great transition for relational database programmers who are looking to work with Hadoop.

MapReduce

MapReduce applications can process vast amounts (multiple terabytes) of data in parallel on large clusters in a reliable, fault-tolerant manner.

NoSQL

A database management system that does not use SQL as its primary query language.



INFORMATION ASSET, LLC

Pig

Apache Pig is a platform for analyzing large semi-structured data sets in Hadoop. It uses a procedural language called Pig Latin that insulates users from learning the intricacies of MapReduce programming in Java. Hive and Pig evolved as separate Apache projects for the analysis of large data sets. Hive is better suited to users who are familiar with SQL. Pig, on the other hand, is ideal for users who are familiar with procedural programs like Microsoft Visual Basic and Python.

Sqoop

Apache Sqoop is a tool that supports the movement of massive volumes of data between Apache Hadoop and structured data stores such as relational databases. Many ETL vendors also support similar functionality.

Stream computing

A class of applications that perform high-performance, low-latency processing of large volumes of data leveraging parallel processing capabilities without landing data to disk.

About the Author

Sunil Soares is the Founder and Managing Partner of Information Asset, a consulting firm focused on Data Governance and Enterprise Data Management. He is the author of several books, including *Selling Information Governance to the Business and Big Data Governance*. His book on *Data Governance Tools* will be published in Fall 2014. For more information, please visit www.information-asset.com.

© 2014 Copyright Information Asset, LLC. All rights reserved.

THIS MATERIAL MAY NOT BE REPRODUCED, DISPLAYED, MODIFIED, OR DISTRIBUTED WITHOUT THE EXPRESS PRIOR WRITTEN PERMISSION OF INFORMATION ASSET, LLC.

Product or company names mentioned herein may be the trademarks of their respective owners.

This report is for informational purposes only and is provided “as is” with no warranties whatsoever, including any warranty of merchantability, fitness for any particular purpose, or any warranty otherwise arising out of any proposal, specification, or sample.