

IBM SPSS Analytic Server
Version 2.1.0.1

Overview

IBM

Note

Before using this information and the product it supports, read the information in "Notices" on page 5.

Product Information

This edition applies to version 2.1.0.1, release 1, modification 0 of IBM SPSS Analytic Server and to all subsequent releases and modifications until otherwise indicated in new editions.

Contents

Overview	1	Notices	5
Architecture	2	Trademarks.	7
Spark and Analytic Server	2		
What is new in version 2.1.	3		

Overview

IBM® SPSS® Analytic Server is a solution for big data analytics that combines IBM SPSS technology with big data systems and allows you to work with familiar IBM SPSS user interfaces to solve problems on a previously unattainable scale.

Why big data analytics matters

Data volumes collected by organizations are growing exponentially; for example, financial and retail businesses have all customer transactions for a year (or two years, or ten years), telco providers have call data records (CDR) and device sensor readings, and internet companies have the results of web crawls.

Big data analytics is needed where there exists:

- A large volume of data (terabytes, petabytes, exabytes), especially when it is a mixture of structured & unstructured data
- Rapidly changing/accumulating data

Big data analytics also assists when:

- A large number (thousands) of models are being built
- Models are frequently built/refreshed

Challenges

The same organizations that collect large volumes of data often have difficulty actually making use of it, for a variety of reasons:

- The architecture of traditional analytic products are not suited to distributed computation, and
- Existing statistical algorithms are not designed to work with big data (these algorithms expect the data to come to them, but big data is too costly to move), thus
- Performing state of the art analytics on big data requires new skills and intimate knowledge of big data systems. Very few analysts have these skills.
- In-memory solutions work for medium-size problems, but do not scale well to truly big data.

Solution

Analytic Server provides:

- A data-centric architecture that leverages big data systems, such as Hadoop Map/Reduce with data in HDFS.
- A defined interface to incorporate new statistical algorithms designed to go to the data.
- Familiar IBM SPSS user interfaces that hide the details of big data environments so that analysts can focus on analyzing the data.
- A solution that is scalable to any size problem.

Architecture

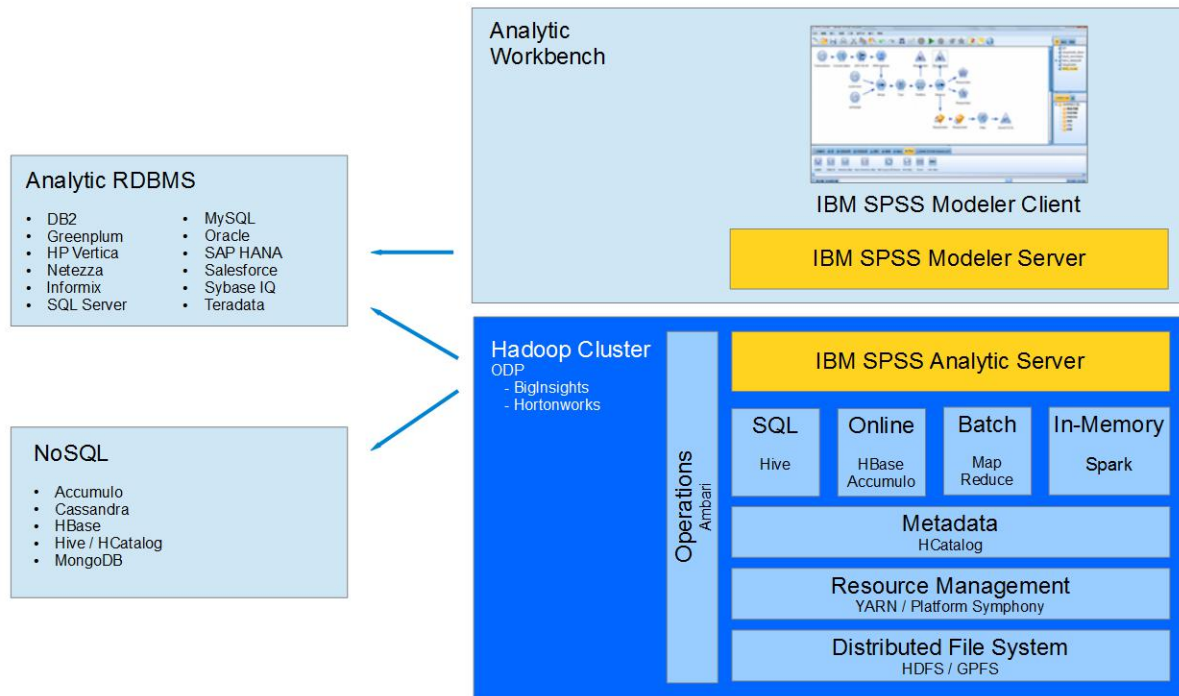


Figure 1. Architecture

Analytic Server sits between a client application and Hadoop cloud. Assuming that the data resides in the cloud, the general outline for working with Analytic Server is to:

1. Define Analytic Server data sources over the data in the cloud.
2. Define the analysis you want to perform in the client application. For the current release, the client application is IBM SPSS Modeler.
3. When you run the analysis, the client application submits an Analytic Server execution request.
4. Analytic Server orchestrates the job to run in the Hadoop cloud and reports the results to the client application.
5. You can use the results to define further analyses, and the cycle repeats.

Spark and Analytic Server

Analytic Server integrates with Apache Spark to increase performance.

When Spark is and is not used

If Spark is installed as an Ambari service on the Hadoop cluster, then Analytic Server uses it to process big data jobs. The following guidelines apply to determine when Spark is not used.

1. If the data set is smaller than 128MB, then Analytic Server uses the embedded MapReduce function in the Analytic Server JVM and does not utilize Spark or the Hadoop cluster.
2. If Spark is not installed on the cluster, then Analytic Server uses MapReduce v2.
3. Analytic Server uses MapReduce v2 to build PSM models. When a job ends with a PSM model build, Analytic Server uses Spark to process the job through all steps leading to the model build, then write

to disk, and then use MapReduce to build the PSM model. For example, if a job includes a join followed by a PSM model build, the join runs in Spark and the PSM runs on the joined data in MapReduce.

How Spark is used

When the Analytic Server service is started and discovers that Spark is available, it initializes a "Spark Hadoop job" that allows communication between distributed tasks across the cluster. This job runs for as long as the Analytic Server service runs, and is used for all Analytic Server executions. This approach improves performance relative to orchestrating multiple MapReduce Hadoop jobs, because it eliminates the overhead of reloading all Analytic Server components for each Hadoop Job.

Spark is capable of running MapReduce jobs. This allows Analytic Server to use "native" Spark algorithms such as join, sort, and union where available. At the same time, Analytic Server can run existing SPSS Map and Reduce algorithms in Spark, and without directly using the Hadoop API.

What is new in version 2.1

Analytics

Spark integration

When Analytic Server is installed as an Ambari service in a Hadoop cluster where Spark is installed, all Analytic Server jobs run in Spark, rather than MapReduce.

Additionally, you can include PySpark scripts in Modeler streams through the Custom Dialog Builder. See the Modeler Extensions for details on how to create custom Modeler nodes that call PySpark scripts.

Support for new SPSS Modeler functionality

Generalized linear models

Added support for distributed building and scoring of generalized linear models. See the GLE node in SPSS Modeler.

Support vector machines (SVM)

Added support for distributed building and scoring of linear support vector machines. See the LSVM node in SPSS Modeler.

Random trees

Added support for distributed building and scoring of random trees. See the Random Trees node in SPSS Modeler.

Improved support for existing SPSS Modeler functionality

Restructure

When large numbers of new fields are created, the Restructure node performs faster than before.

Distinct

The Distinct node performs faster than before.

Predictor importance

Models that compute predictor importance as part of the output perform faster than before.

2.1.0.1 only. Merge

Merging datasets in 2.1.0.1 performs faster than in previous versions, with greater increases in performance as the size of the data increases.

2.1.0.1 only. Random trees

Both building and scoring random trees performs faster than in previous versions. Scoring performance is comparable to native Spark scoring.

Analytic Server console

Data sources

Relational databases

You can define data sources for Amazon Redshift, if Analytic Server has been configured to be able to use this data source.

2.1.0.1 only. Analytic Server now uses the Redshift driver instead of Postgres.

NoSQL databases

The names of the storage handlers have changed since version 2.0. For resources on setting up external Hive tables, see Using HCatalog data sources

Reader role

Within a tenant, you can assign users and groups to a Reader role that cannot log in to the Analytic Server, but can read Analytic Server data sources through the Analytic Server Source node in Modeler.

Installation and configuration

Analytic Server is now installed and runs as an Apache Ambari service. This speeds and simplifies the installation and management of Analytic Server relative to earlier versions.

Platform

Support for operating systems and Hadoop distributions is changed from version 2.

Operating systems

Analytic Server now runs on Red Hat Enterprise Linux (Power LE) in addition to existing operating system support.

Hadoop distributions

Analytic Server runs with Big Insights and Hortonworks, the distributions supported by Ambari. On these distributions, Analytic Server is configured as an Ambari service.

2.1.0.1 only. Analytic Server runs with Cloudera and MapR. On Cloudera, Analytic Server is configured using Cloudera Manager. On MapR, Analytic Server installation and configuration is a manual process.

Metadata repository

Analytic Server no longer supports Derby as the default metadata repository, and instead uses MySQL. DB2 is still supported as an alternative repository.

Database data sources

Data sources can be defined for Amazon Redshift, in addition to existing database support.

For the most up-to-date system requirements information, use the Detailed system requirements reports at the IBM Technical Support site: <http://publib.boulder.ibm.com/infocenter/prodguid/v1r0/clarity/softwareReqsForProduct.html>. On this page:

1. Type SPSS Analytic Server as the product name and click **Search**.
2. Select the wanted version and scope of report, then click **Submit**.

Notices

This information was developed for products and services offered in the U.S.A.

IBM may not offer the products, services, or features discussed in this document in other countries. Consult your local IBM representative for information on the products and services currently available in your area. Any reference to an IBM product, program, or service is not intended to state or imply that only that IBM product, program, or service may be used. Any functionally equivalent product, program, or service that does not infringe any IBM intellectual property right may be used instead. However, it is the user's responsibility to evaluate and verify the operation of any non-IBM product, program, or service.

IBM may have patents or pending patent applications covering subject matter described in this document. The furnishing of this document does not grant you any license to these patents. You can send license inquiries, in writing, to:

IBM Director of Licensing
IBM Corporation
North Castle Drive
Armonk, NY 10504-1785
U.S.A.

For license inquiries regarding double-byte (DBCS) information, contact the IBM Intellectual Property Department in your country or send inquiries, in writing, to:

Intellectual Property Licensing
Legal and Intellectual Property Law
IBM Japan Ltd.
1623-14, Shimotsuruma, Yamato-shi
Kanagawa 242-8502 Japan

The following paragraph does not apply to the United Kingdom or any other country where such provisions are inconsistent with local law: INTERNATIONAL BUSINESS MACHINES CORPORATION PROVIDES THIS PUBLICATION "AS IS" WITHOUT WARRANTY OF ANY KIND, EITHER EXPRESS OR IMPLIED, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF NON-INFRINGEMENT, MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE. Some states do not allow disclaimer of express or implied warranties in certain transactions, therefore, this statement may not apply to you.

This information could include technical inaccuracies or typographical errors. Changes are periodically made to the information herein; these changes will be incorporated in new editions of the publication. IBM may make improvements and/or changes in the product(s) and/or the program(s) described in this publication at any time without notice.

Any references in this information to non-IBM Web sites are provided for convenience only and do not in any manner serve as an endorsement of those Web sites. The materials at those Web sites are not part of the materials for this IBM product and use of those Web sites is at your own risk.

IBM may use or distribute any of the information you supply in any way it believes appropriate without incurring any obligation to you.

Licensees of this program who want to have information about it for the purpose of enabling: (i) the exchange of information between independently created programs and other programs (including this one) and (ii) the mutual use of the information which has been exchanged, should contact:

IBM Software Group
ATTN: Licensing
200 W. Madison St.
Chicago, IL; 60606
U.S.A.

Such information may be available, subject to appropriate terms and conditions, including in some cases, payment of a fee.

The licensed program described in this document and all licensed material available for it are provided by IBM under terms of the IBM Customer Agreement, IBM International Program License Agreement or any equivalent agreement between us.

Any performance data contained herein was determined in a controlled environment. Therefore, the results obtained in other operating environments may vary significantly. Some measurements may have been made on development-level systems and there is no guarantee that these measurements will be the same on generally available systems. Furthermore, some measurements may have been estimated through extrapolation. Actual results may vary. Users of this document should verify the applicable data for their specific environment.

Information concerning non-IBM products was obtained from the suppliers of those products, their published announcements or other publicly available sources. IBM has not tested those products and cannot confirm the accuracy of performance, compatibility or any other claims related to non-IBM products. Questions on the capabilities of non-IBM products should be addressed to the suppliers of those products.

All statements regarding IBM's future direction or intent are subject to change or withdrawal without notice, and represent goals and objectives only.

All IBM prices shown are IBM's suggested retail prices, are current and are subject to change without notice. Dealer prices may vary.

This information is for planning purposes only. The information herein is subject to change before the products described become available.

This information contains examples of data and reports used in daily business operations. To illustrate them as completely as possible, the examples include the names of individuals, companies, brands, and products. All of these names are fictitious and any similarity to the names and addresses used by an actual business enterprise is entirely coincidental.

Each copy or any portion of these sample programs or any derivative work, must include a copyright notice as follows:

This information contains examples of data and reports used in daily business operations. To illustrate them as completely as possible, the examples include the names of individuals, companies, brands, and products. All of these names are fictitious and any similarity to the names and addresses used by an actual business enterprise is entirely coincidental.

Each copy or any portion of these sample programs or any derivative work, must include a copyright notice as follows:

© your company name) (year). Portions of this code are derived from IBM Corp. Sample Programs.

© Copyright IBM Corp. _enter the year or years_. All rights reserved.

If you are viewing this information softcopy, the photographs and color illustrations may not appear.

Trademarks

IBM, the IBM logo, and ibm.com are trademarks or registered trademarks of International Business Machines Corp., registered in many jurisdictions worldwide. Other product and service names might be trademarks of IBM or other companies. A current list of IBM trademarks is available on the Web at "Copyright and trademark information" at www.ibm.com/legal/copytrade.shtml.

Adobe, the Adobe logo, PostScript, and the PostScript logo are either registered trademarks or trademarks of Adobe Systems Incorporated in the United States, and/or other countries.

IT Infrastructure Library is a registered trademark of the Central Computer and Telecommunications Agency which is now part of the Office of Government Commerce.

Intel, Intel logo, Intel Inside, Intel Inside logo, Intel Centrino, Intel Centrino logo, Celeron, Intel Xeon, Intel SpeedStep, Itanium, and Pentium are trademarks or registered trademarks of Intel Corporation or its subsidiaries in the United States and other countries.

Linux is a registered trademark of Linus Torvalds in the United States, other countries, or both.

Microsoft, Windows, Windows NT, and the Windows logo are trademarks of Microsoft Corporation in the United States, other countries, or both.

ITIL is a registered trademark, and a registered community trademark of The Minister for the Cabinet Office, and is registered in the U.S. Patent and Trademark Office.

UNIX is a registered trademark of The Open Group in the United States and other countries.

Cell Broadband Engine is a trademark of Sony Computer Entertainment, Inc. in the United States, other countries, or both and is used under license therefrom.

Linear Tape-Open, LTO, the LTO Logo, Ultrium, and the Ultrium logo are trademarks of HP, IBM Corp. and Quantum in the U.S. and other countries.



Printed in USA