# Integrating Relational Data with Hadoop and Spark

**think** 2018

—

Stefan Hummel
stefan.hummel@de.ibm.com

Andreas Weininger
andreas.weininger@de.ibm.com

IBM

# Please note

IBM's statements regarding its plans, directions, and intent are subject to change or withdrawal without notice and at IBM's sole discretion.

Information regarding potential future products is intended to outline our general product direction and it should not be relied on in making a purchasing decision.

The information mentioned regarding potential future products is not a commitment, promise, or legal obligation to deliver any material, code or functionality. Information about potential future products may not be incorporated into any contract.
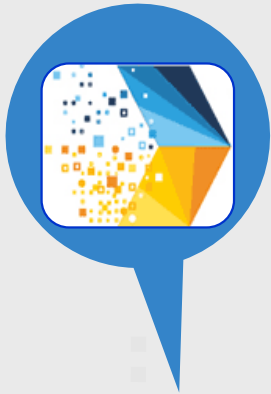
The development, release, and timing of any future features or functionality described for our products remains at our sole discretion.

Performance is based on measurements and projections using standard IBM benchmarks in a controlled environment. The actual throughput or performance that any user will experience will vary depending upon many factors, including considerations such as the amount of multiprogramming in the user's job stream, the I/O configuration, the storage configuration, and the workload processed. Therefore, no assurance can be given that an individual user will achieve results similar to those stated here.
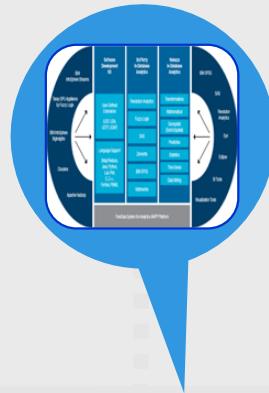
# Labs

- Lab 1: Creating Tables
- Lab 2: Loading Data
- Lab 3: Executing Queries
- Lab 4: Linear Regression in SQL
- Lab 5: Linear Regression with R
- Lab 6: Linear Regression with Python
- Lab 7: Analyzing data with R in RStudio
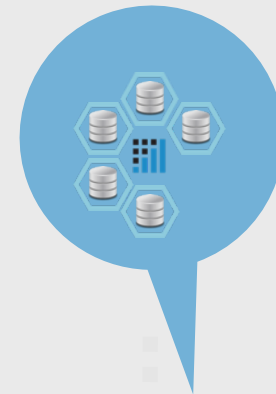
# Db2 Warehouse – Analytics Warehouse as a Service

## BLU Acceleration

- DB2 BLU columnar technology

- in-memory processing, data skipping, actionable compression, parallel vector processing, , "Load & Go" administration

## Netezza In-Database Analytics

- Netezza predictive analytic algorithms

- fully integrated RStudio & R language

## Db2 Warehouse MPP

- Massively Parallel Processing (MPP)

- Oracle compatibility

- fully-managed warehouse

# Db2 Warehouse – Key Use Cases

**Cloudant Analytics**

- Easy synchronization of JSON to structured data
- Allows analytics via standard BI tools
- In-database predictive algorithms allow greater insight for Cloudant users than ever before

**Extend or Modernize**

- Extend on-premise data warehouses to the cloud
- Flexible, cost-effective growth
- Hybrid Cloud model to support ground to cloud

**In-Database Analytics**

- Robust predictive analytic algorithms
- Integrated with R (and Python/Spark)
- Watson Analytics Ready
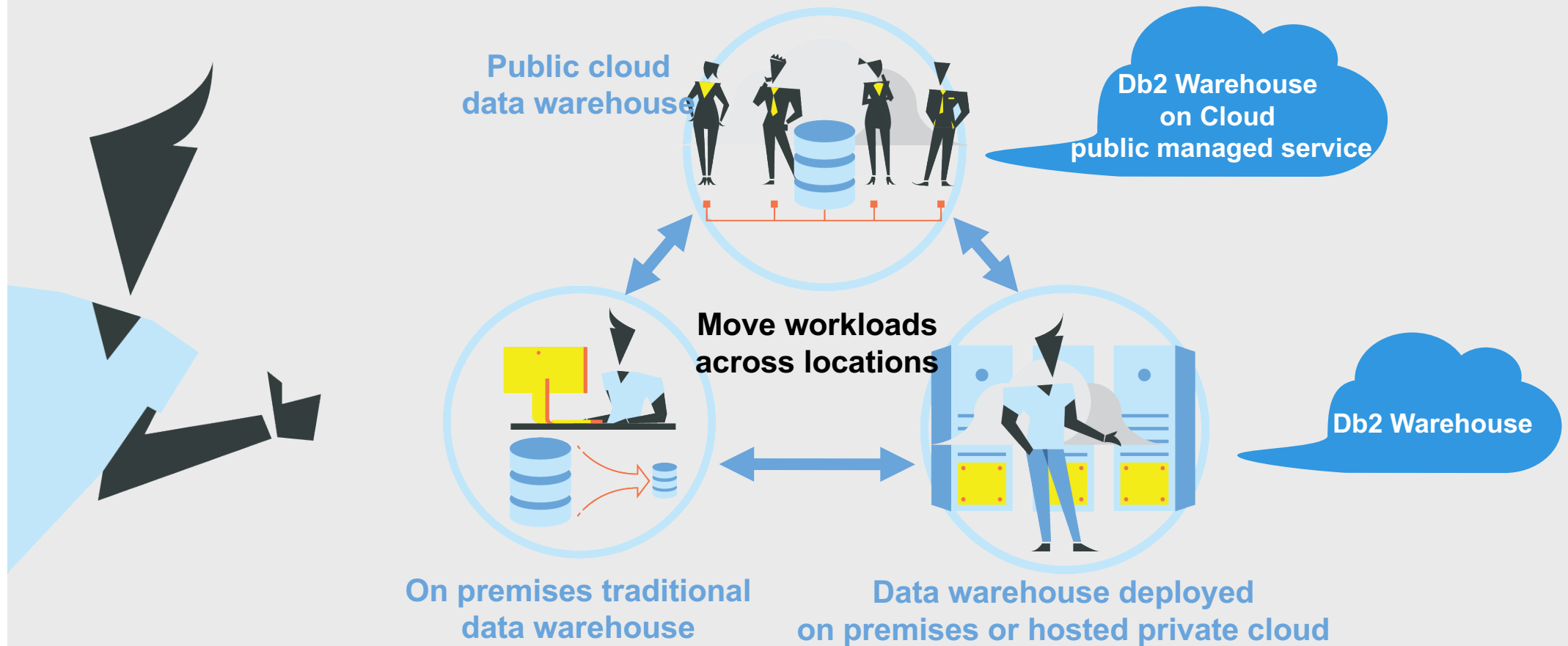- Analytics Ecosystem with Partners
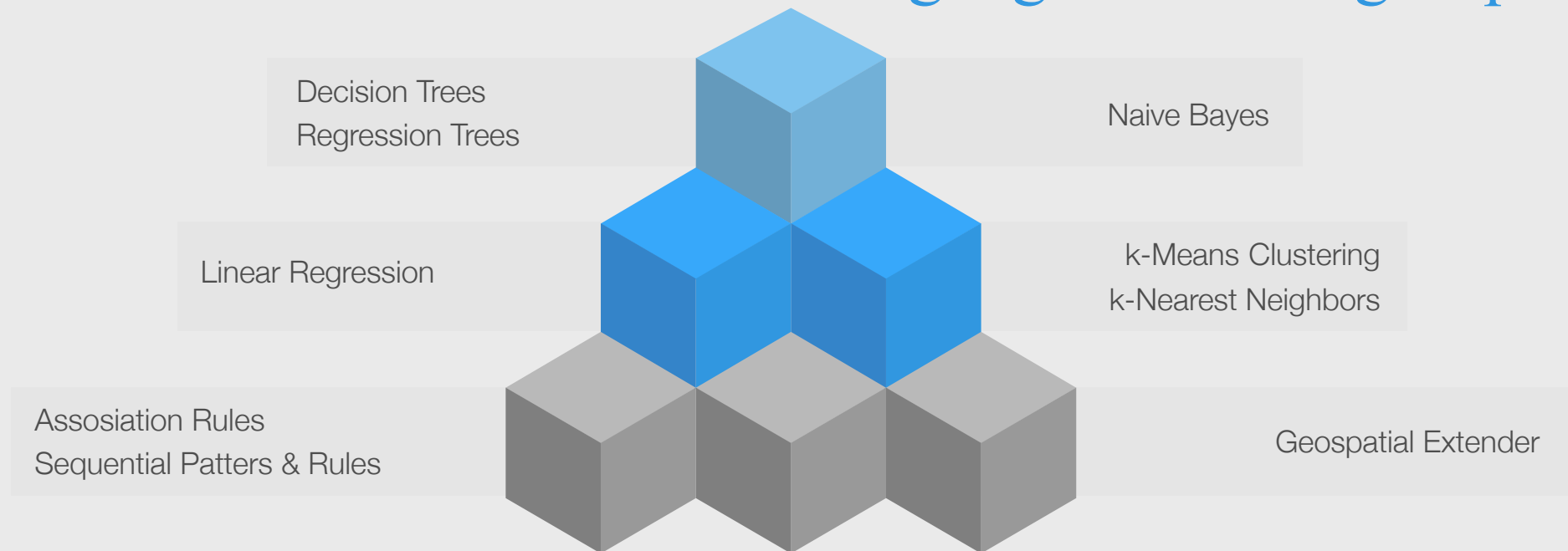
**Data Warehouse & Analytics Service**

- Data Warehousing and Analytics in the Cloud
- Cloud Agility and Flexibility
- Analytics for Cloud Data, Data Marts, and development and test environments

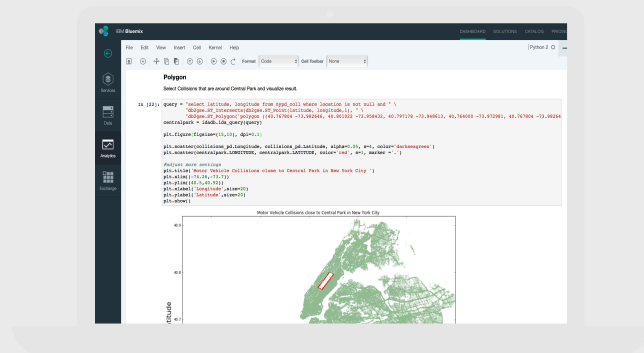# Db2 Warehouse offerings – on cloud or local

**Public cloud data warehouse**

**Db2 Warehouse on Cloud public managed service**

**Move workloads across locations**

**Db2 Warehouse**

**On premises traditional data warehouse**

**Data warehouse deployed on premises or hosted private cloud**

# Db2 Warehouse – machine learning algorithms & geospatial

Decision Trees
Regression Trees

Naive Bayes

Linear Regression

k-Means Clustering
k-Nearest Neighbors

Assosiation Rules
Sequential Patters & Rules

Geospatial Extender

# In-database analytics on Db2 Warehouse

Run machine learning algorithms and do geospatial analytics directly on the data in Db2 Warehouse leveraging the efficient parallel database engine.

# Db2 Warehouse – R and Python integration



- Integrated with Spark Service offering in Bluemix
- Access data in Db2 Warehouse from notebook interface in Spark Service or Data Science Experience GUI
- Run machine learning algorithms directly in database engine



- Access data in Db2 Warehouse from RStudio environment
- RStudio hosted on Db2 Warehouse server
- Machine learning algorithms of Db2 Warehouse available via ibmdbR library
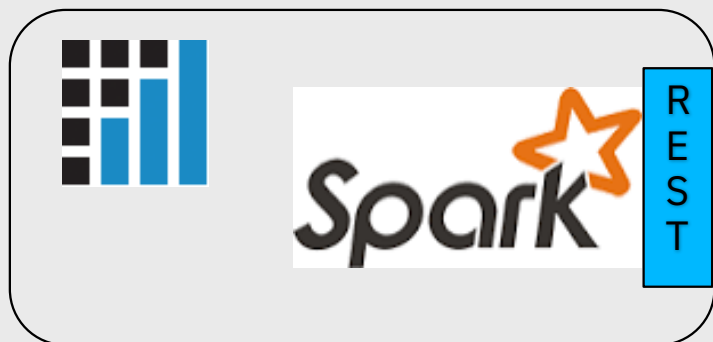- Possibility to pushdown R code to Db2 Warehouse server

# Db2 Warehouse – Spark Integration

Communication



- Spark reading and writing Db2 Warehouse tables

Administration



- Administration of Spark code via
  - REST interface
- Monitoring of Spark via
  - Db2 Warehouse Web Console or
  - REST Interface

# Db2 Warehouse Spark Integration Examples

```
CALL SPARK.SUBMIT('jarfile=myapp/app.jar
class=com.ibm.dashdb.spark.DemoJob')

CALL IDAX.GLM('model=my_model,
intable=CUST, target=CENSOR, id=ID');
```

Head/Coordinator Node

Db2 Warehouse
Relational Engine

Spark
Analytics Engine

Db2 Warehouse

**Interactive Jupyter Notebook**

jupyter

**REST API**
/dashdb-api/analytics/public/apps/submit
                              ../cancel
              ../monitoring/app_status

**spark-submit.sh**
Command Line Tool

Data Nodes

Db2 Warehouse
Relational Engine

Spark
Analytics Engine

Db2 Warehouse

Db2 Warehouse
Relational Engine

Spark
Analytics Engine

Db2 Warehouse

Db2 Warehouse
Relational Engine

Spark
Analytics Engine

Db2 Warehouse

Relational
DB Partitions

Relational
DB Partitions

Relational
DB Partitions

# Some hints for executing the labs

- Start with starting everything as described in in Appendix I
- Check with
  `docker exec Db2wh status`
  whether Db2 Warehouse is running
- Remove the directory Advanced-Analytics in the home directory of user ibmuser with
  `rm –rf Advanced-Analytics`
- Only after that execute the `git clone` command
- The Jupyter notebooks can be found in the subdirectory Advanced-Analytics/notebooks

# Notices and disclaimers

# Notices and disclaimers continued

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 about this publication 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. IBM does not warrant the quality of any third-party products, or the ability of any such third-party products to interoperate with IBM's products. **IBM expressly disclaims all warranties, expressed or implied, including but not limited to, the implied warranties of merchantability and fitness for a purpose.** The provision of the information contained herein is not intended to, and does not, grant any right or license under any IBM patents, copyrights, trademarks or other intellectual property right.

IBM, the IBM logo, ibm.com and [names of other referenced IBM products and services used in the presentation] are trademarks of International Business Machines Corporation, 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.
.

# Thank you

Stefan Hummel                    Andreas Weininger

—

stefan.hummel@de.ibm.com        andreas.weininger@de.ibm.com

+49-160-742-1795                +49-172-756-5266
ibm.com