



DATABIZ

your solution builder

Going Big Data on Apache Spark

KNIME Italy Meetup

Agenda

Introduction

Why Apache Spark?

Section 1

Gathering Requirements

Section 2

Tool Choice

Section 3

Architecture

Section 4

Devising New Nodes

Section 5

Conclusion

Let's dive

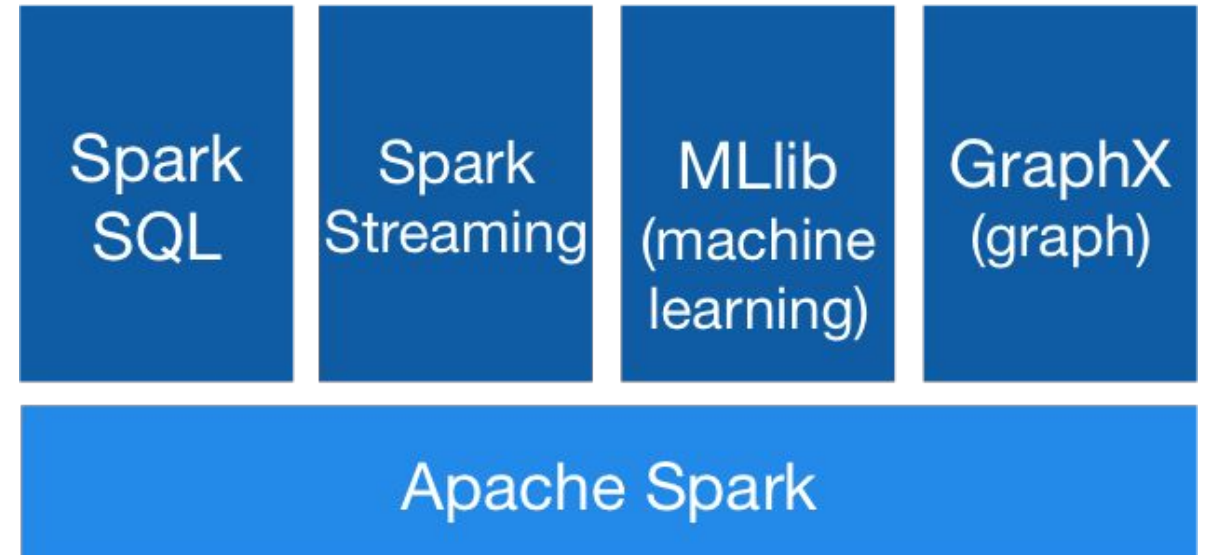


Introduction

Why Apache Spark?

Apache Spark

- Fast engine for large-scale distributed data processing
- Builds and tests predictive models in little time
- Built-in modules for:
 - SQL
 - Streaming
 - Machine Learning
 - Graph Processing



(Source: [Apache Spark](#))

Dealing with Apache Spark



Apache Spark has a steep
learning curve



Use a data scientist-friendly
interface with Spark
integration

Section 1

Gathering Requirements

Requirements



Tasks

Explore data on Hadoop Ecosystem
Leverage Spark for fast analysis
Data Preparation
Perform statistical analysis



Required Capabilities

Hadoop integration
Spark integration
Modeling
Extensibility

Section 2

Tool Choice

Figure 1. Magic Quadrant for Advanced Analytics Platforms



(Source: [Gartner](https://www.gartner.com/doc/2784441), February 2015)

Product	Suitable	Notes
Alpine	✓	
KNIME	✓	
RapidMiner	✓	
KXEN	?	Spark integration on the roadmap
SAS	?	Spark integration on the roadmap
IBM SPSS	?	Spark integration on the roadmap

What About KNIME?

- User-friendly interface
- Open Source
- Integration with Hadoop / Spark
- Clear pricing and cost effectiveness
- Rich existing feature set
- Possibility of co-development





**OK, but was it
a good choice?**

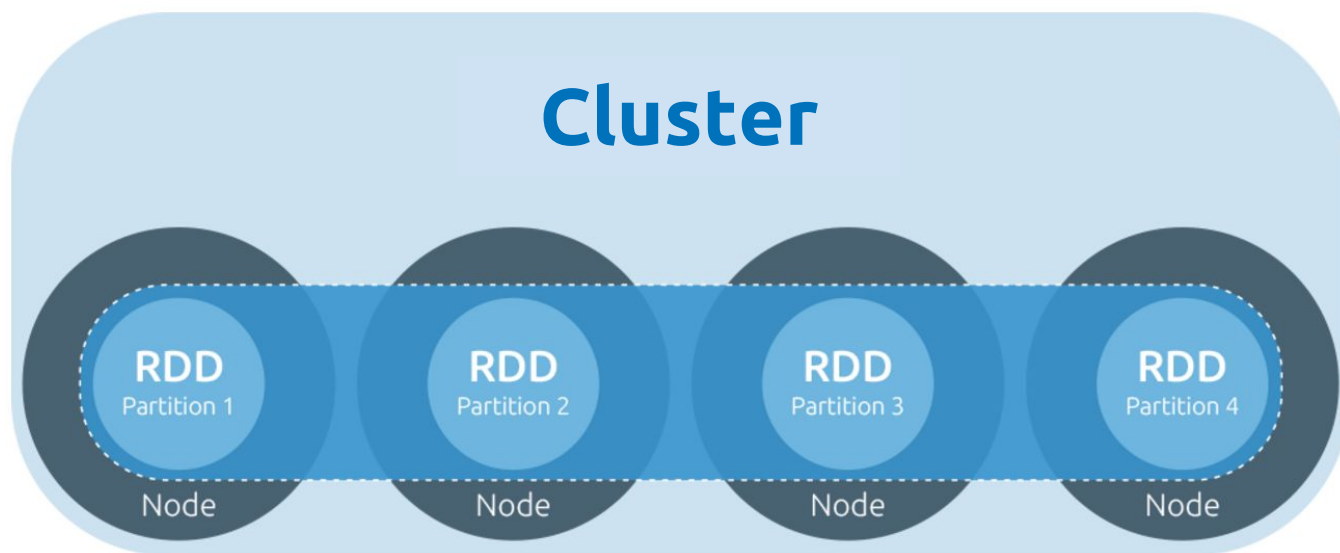


(Source: [Gartner](http://www.gartner.com), February 2015)

Section 3

Architecture

Spark's Building Block: RDD



Immutable

Each step of a dataflow will create a new RDD

Lazy

Data are processed only when results are requested

Resilient

A lost partition is reconstructed from its lineage

An underwater photograph with a blue-green tint. In the lower right, a diver's hand in a black wetsuit sleeve reaches out towards the center. The water is filled with bubbles and light rays filtering down. The text 'OK, but how can I use Spark from KNIME ?' is overlaid in white in the upper left quadrant.

OK, but how can I use
Spark from KNIME ?

From KNIME to Spark: Spark Job Server

- Simple REST interface for all aspects of Spark (job and context) management
- Support for Spark SQL, Hive and custom job contexts
- Named RDDs and DataFrames: computed RDDs and DataFrames can be cached with a given name and retrieved later on

Spark Job Creation

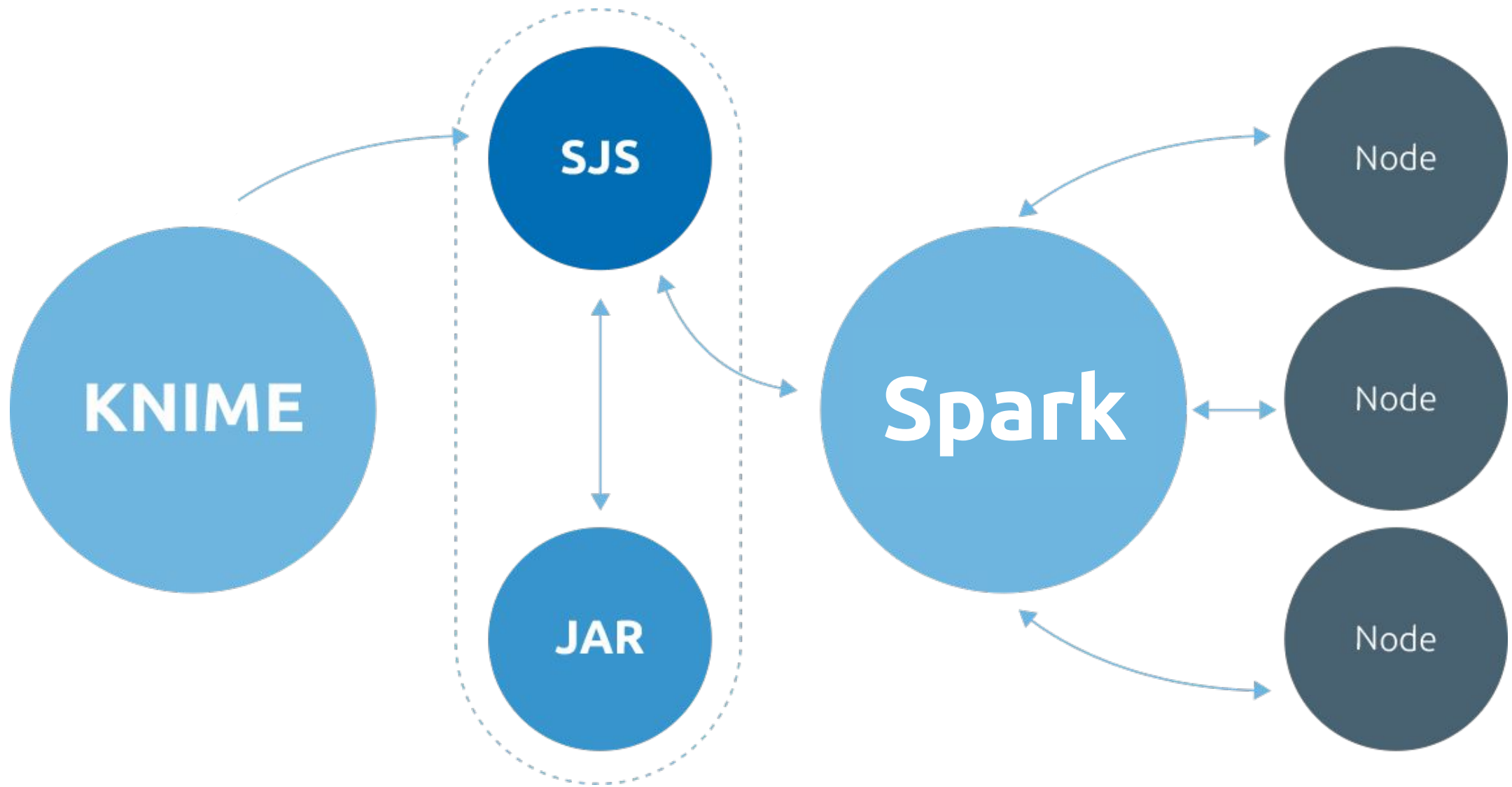
To create a job that can be submitted through the job server, the job must extend the **SparkJob** trait. Your job will look like this:

```
object SampleJob extends SparkJob {  
    override def validate(sc: SparkContext, config: Config): SparkJobValidation = ???  
    override def runJob(sc: SparkContext, jobConfig: Config): Any = ???  
}
```

Spark Job Structure

validate allows for an initial validation of the context and any provided configuration. It helps you preventing running jobs that will eventually fail due to missing or wrong configuration and save both time and resources.

runJob contains the implementation of the Job. The SparkContext is managed by the JobServer and will be provided to the job through this method.



Warning!



Spark Job Server won't send updates to KNIME regarding its internal state

This means that, if Spark Job Server is down or is restarted, you will never know it!



Due to lazy evaluation, functions are computed only when the result is actually required

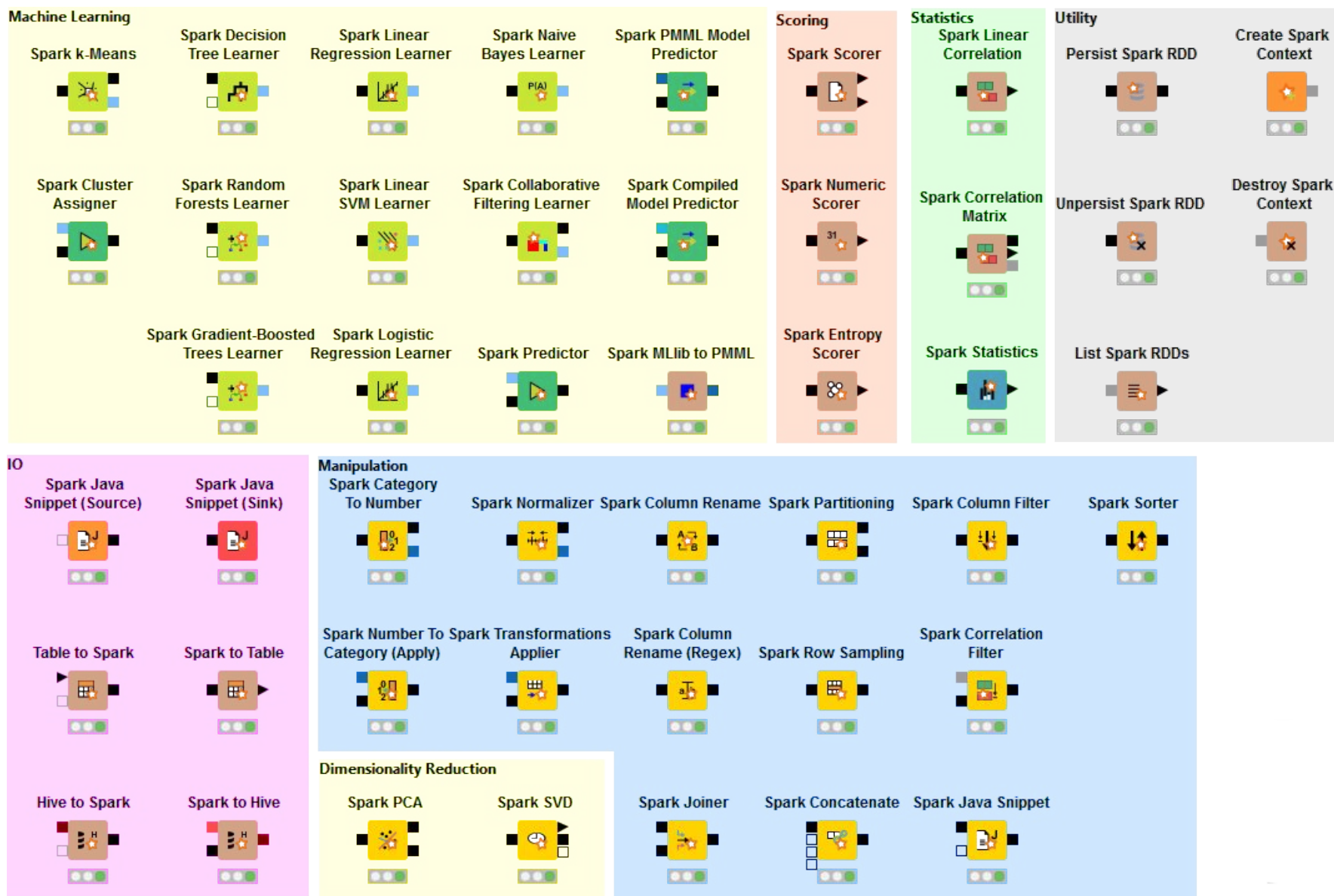
This means that you can't always trust KNIME's green light!

Section 4

Devising New Nodes

A diver in a dark wetsuit is swimming horizontally through a vast, dense school of small fish. The water is a deep blue, and the fish are silhouetted against the lighter background, creating a sense of movement and scale. The diver is positioned in the center-right of the frame, looking towards the left.

**KNIME already comes with
a lot of Spark Nodes**



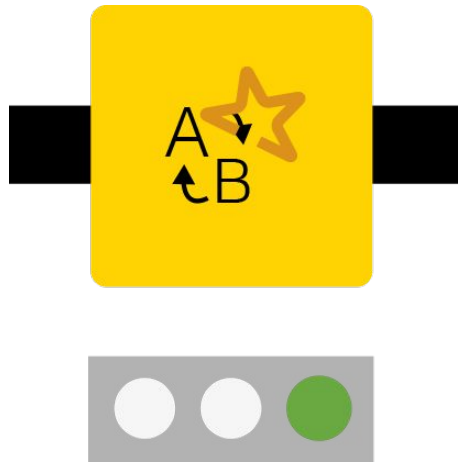
(Source: [KNIME](https://www.knime.org/))

I WANT MORE

MORE! MOAR!! MOOOOARRR!!!



Spark Default Missing Values



- Replaces all occurrences of missing values with fixed default values.
- Can be applied to:
 - Integers
 - Longs
 - Doubles
 - Strings
 - Booleans
 - Dates

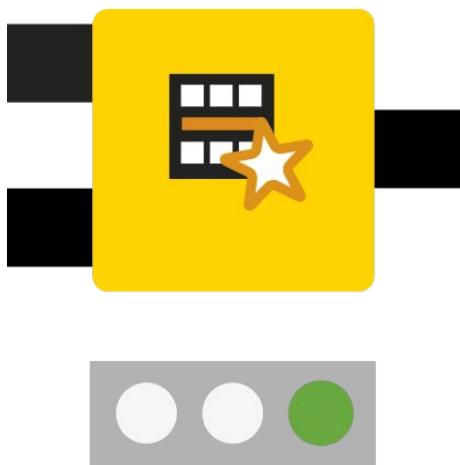
The screenshot shows a dialog box titled "Dialog - 2:1 - Spark Default Missing Values". It has a "File" menu and two tabs: "Job Manager Selection" and "Memory Policy". The "Options" tab is selected, and the "Flow Variables" section is expanded. The "Integer" section has a "Fix Value" of 0. The "Long" section has a "Fix Value" of 0. The "Double" section has a "Fix Value" of 0.0. The "String" section has a "Fix Value" field. The "Boolean" section has a "Fix Value" section with radio buttons for "True" and "False", where "False" is selected. The "Date" section has a "Fix Value" section with a date field set to "Jan 1, 1970" and checkboxes for "Hour", "Minute", and "Second", all of which are unchecked. At the bottom, there are "OK", "Apply", "Cancel", and a help button.

Spark to HBase

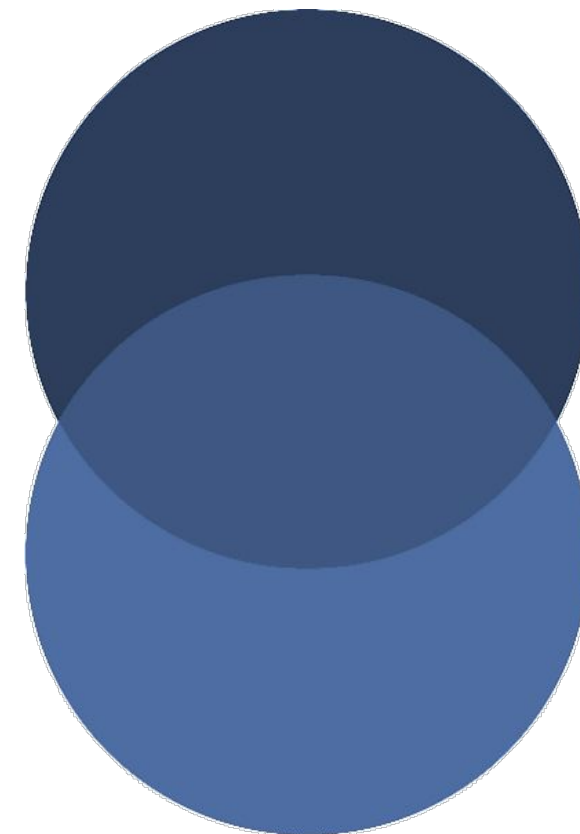


- Persists data from an incoming RDD into HBase.
- Requires:
 - the name of the table you intend to create
 - the name of the column in which the Row IDs are contained
 - the name that will be given to the Column Family

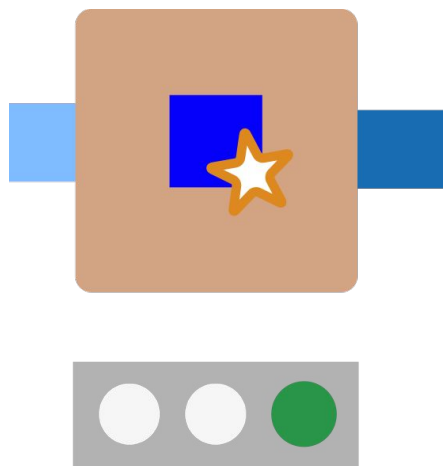
Spark Subtract



- Given two incoming RDDs, returns an element from the first input port RDD that are not contained in the second input port RDD.
- Similar to the relative complement of the set theory

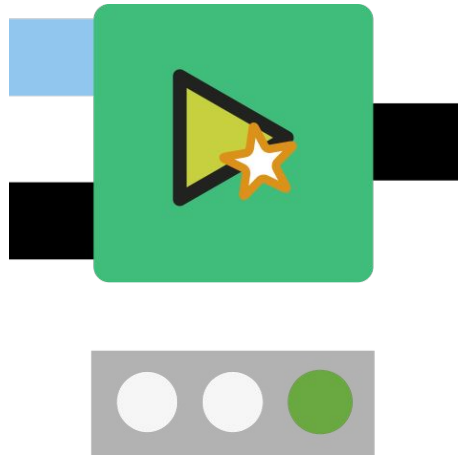


Spark MLlib to PMML

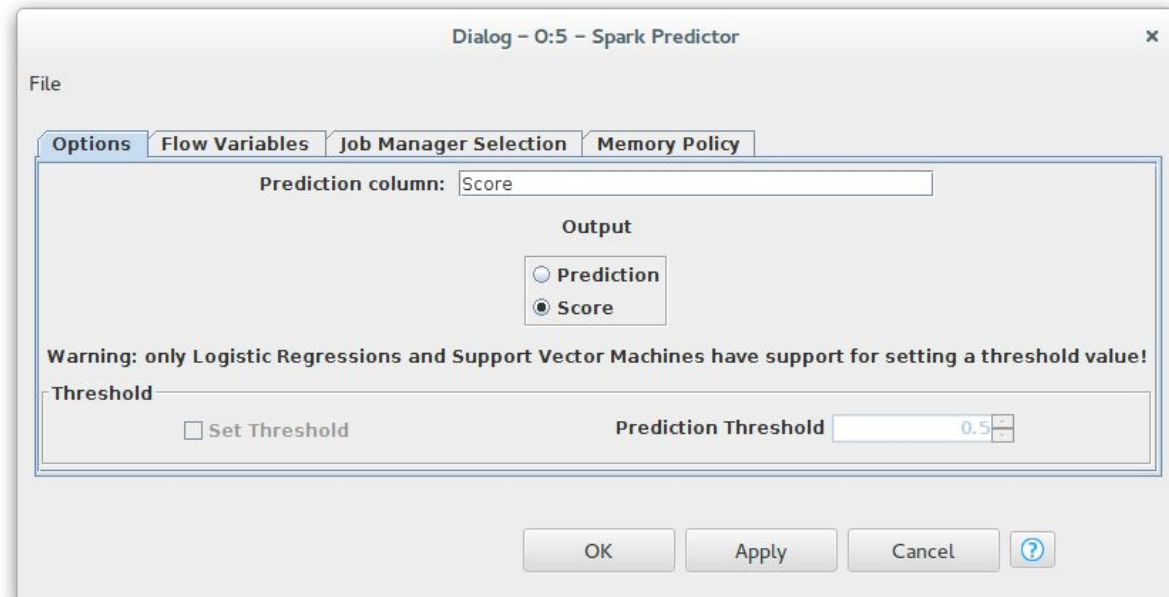


- Converts supported Spark MLlib models into PMML files.
- Supported model are (up to now):
 - Decision Tree
 - K-Means
 - Linear Regression
 - Logistic Regression
 - Naive Bayes
 - Random Forest
 - Support Vector Machines

Spark Predictor + Scoring

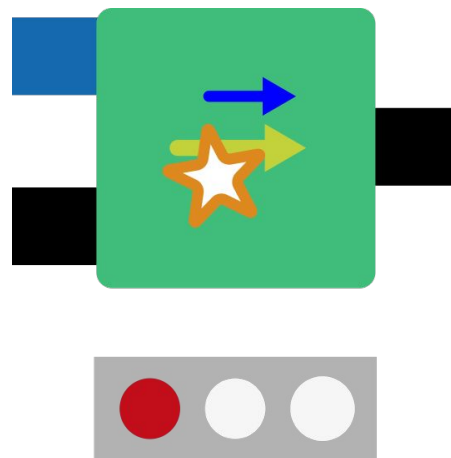


- Labels new data points using a learned Spark MLlib model
- Allows to get the score (the probability that an observation belongs to the positive class)
- Allows to set a threshold



Interlude: Spark JPMML Scorer

Once Upon a Time: Spark PMML Predictor



“Compiles the given PMML model into bytecode and runs it on the input data on Apache Spark”

Spark PMML Predictor: Issues

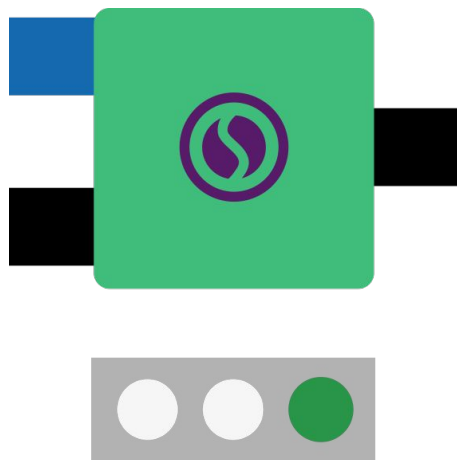
- Returns a prediction, not a score
- *“Compiles the given PMML model into bytecode”*
 - not testable
 - java.lang.ClassFormatError: Invalid method Code length
 - maximum code length is 65534 bytes!

Solution: JPMML

“Java API for Predictive Model Markup Language (PMML)”

- Adopted by the community
- Thoroughly tested (and testable!)
- Easily customizable
- Well documented

Spark JPMML Model Scorer



- Sends the PMML file to the Apache Spark cluster
- Takes advantage of the JPMML library to turn the PMML file into a predictive model
- Uses the JPMML library to give a score to the model's prediction

Other Activities

BugFixing: Spark to Hive & Hive to Spark

- Scala, Java and Hive all have different types
- E.g.:
 - Scala: Int
 - Java: int / Integer
 - Hive: INT
- All of these conversions were handled

Support: Database Connector + Impala

jdbc:impala://your.host:1234/auth=noSasl;UseNativeQuery=1

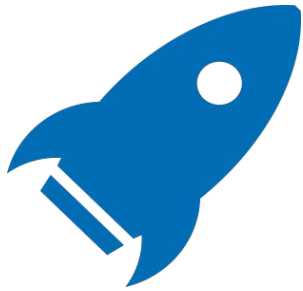
- Enabling the UseNativeQuery option, no transformation is needed to convert the queries into Impala SQL
 - If the application is Impala aware and already emits Impala SQL, enabling this feature avoids the extra overhead of query transformation
- Moreover, we noticed that this solves concurrency issues in Impala

And now...

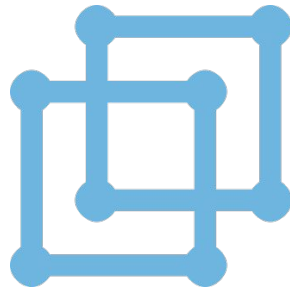
Section 5

Conclusion

Conclusion



- Apache Spark can run programs up to 100x faster than Hadoop MapReduce in memory, or 10x faster on disk



- KNIME gives a data scientist-friendly interface to Apache Spark



- Yet, when dealing with Spark, even from KNIME, a (basic!) understanding of its inner workings is required

Thanks!

Special Thanks

- Sara Aceto
- Stefano Baghino
- Emanuele Bezzi
- Nicola Breda
- Laura Fogliatto
- Simone Grandi
- Tobias Koetter
- Gaetano Mauro
- Thorsten Meinl
- Fabio Oberto
- Luigi Pomante
- Simone Robutti
- Stefano Rocco
- Riccardo Sakakini
- Enrico Scopelliti
- Rosaria Silipo
- Lorenzo Sommaruga
- Marco Tosini
- Marco Veronese
- Giuseppe Zavattoni

A photograph of four surfers standing on a beach at sunset. They are silhouetted against the bright orange and yellow sky. Each surfer is holding a surfboard. The ocean is visible in the background, and the wet sand reflects the figures and the sky. The overall mood is serene and adventurous.

Hey, We Are Hiring!

Send Us Your CV!
hr@databiz.it

Hey, We Are Hiring!



Send Us Your CV!

hr@databiz.it