

Xavier Morera
HELPING DEVELOPERS UNDERSTAND SEARCH & BIG DATA
@xmorera www.xaviermorera.com





Why Spark with Python & Cloudera?

- History
- What we will cover in this training
- What do you need for this course?





Getting an Environment & Data

- Cloudera cluster for Spark on Yarn or Spark Standalone
- StackOverflow / StackExchange





Python Fundamentals for This Course

- History, philosophy and paradigm
 - Functional programming
- Python interactive shell: REPL
- Covered the basics
- PySpark: Spark in the REPL





Understanding Spark: An overview

- Coarse grained transformations
- Parallelism and pipelining
- Narrow and wide transformations
- Lazy execution and lineage
- Libraries





Getting Technical with Spark

- Storage and file formats
- Spark APIs
- Performance, configuration & logging
- Visualizing your Spark application
- Navigating the documentation





Learning the Core of Spark: RDDs

- SparkContext
- Resilient Distributed Datasets
- Many ways of creating RDDs
- Return data to the Driver
- Save RDDs



Going Deeper into Spark Core

- Anonymous functions
- Transformations
- Actions
- Partition operations
- Set operations and aggregations
- Caching and data persistence
- Shared variables
- What's better?



Increasing Proficiency with Spark: DataFrames & Spark SQL

- "Everyone uses SQL"
- SparkSession
- Creating DataFrames
- Schemas
- Rows, columns, expressions & operators
- UDFs





Continuing the Journey with DataFrames & Spark SQL

- DataFrame DSL
- Saving DataFrames
- Spark SQL & Temporary Views
- Peristent Tables
- Hive support & external databases
- Aggregating, grouping & joining
- Catalog API





Understanding a Typed API: Datasets

- Works with Scala, not Python
- Typed API
- DataFrames = Dataset[Row]
- Scala is strongly typed



Continuing the Journey with Spark, Python & Cloudera: Next Steps

Getting a Hadoop Cluster







DirectorCloud



AltusPlatform-as-a-Service



Cloudera Altus



Platform-as-a-Service

Provision clusters in the cloud

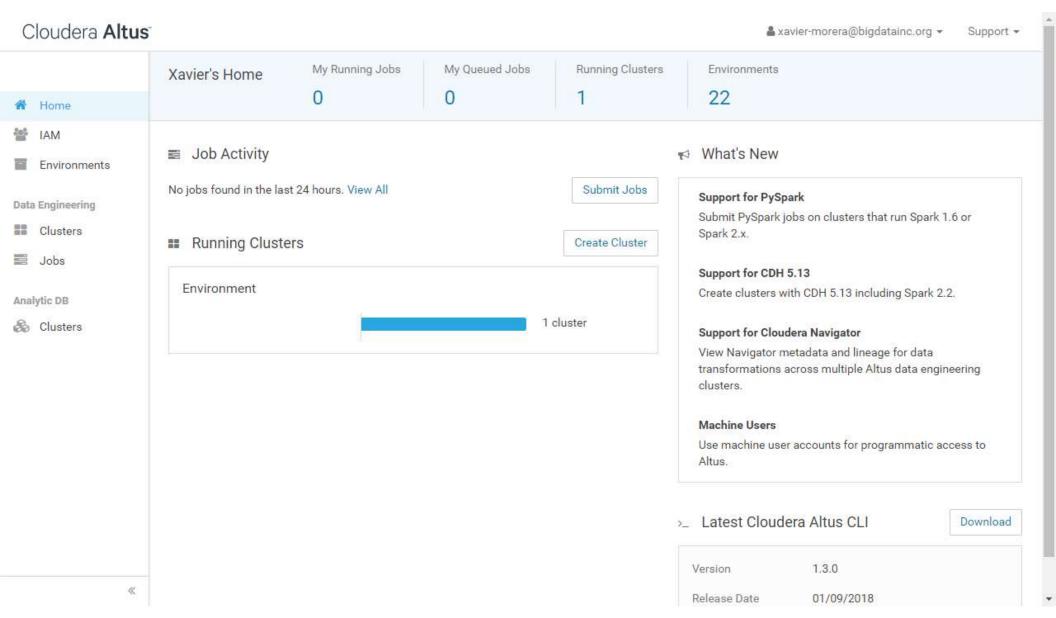
- AWS & Azure
- Managed cluster

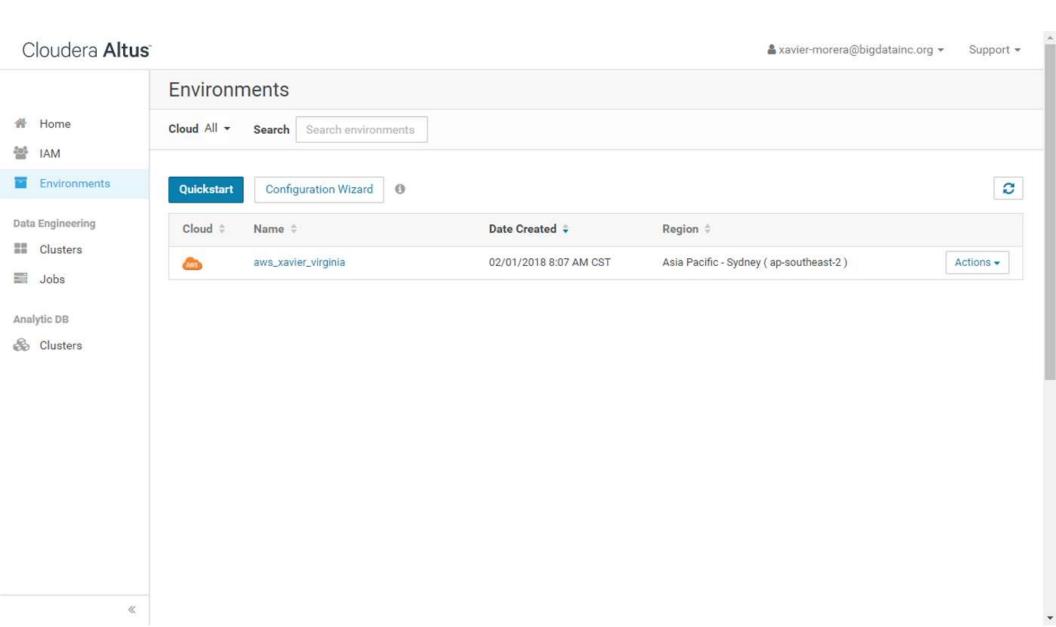
Data Engineering & Workload Analytics

Processing engines

- MR2, Hive on MR2, Spark & Hive on Spark



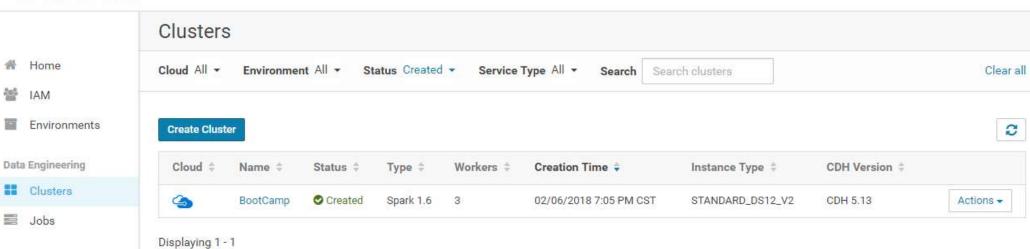




Cloudera Altus

Analytic DB

& Clusters



å xavier-morera@bigdatainc.org ▼

Support +

Submit Jobs

Job Settings

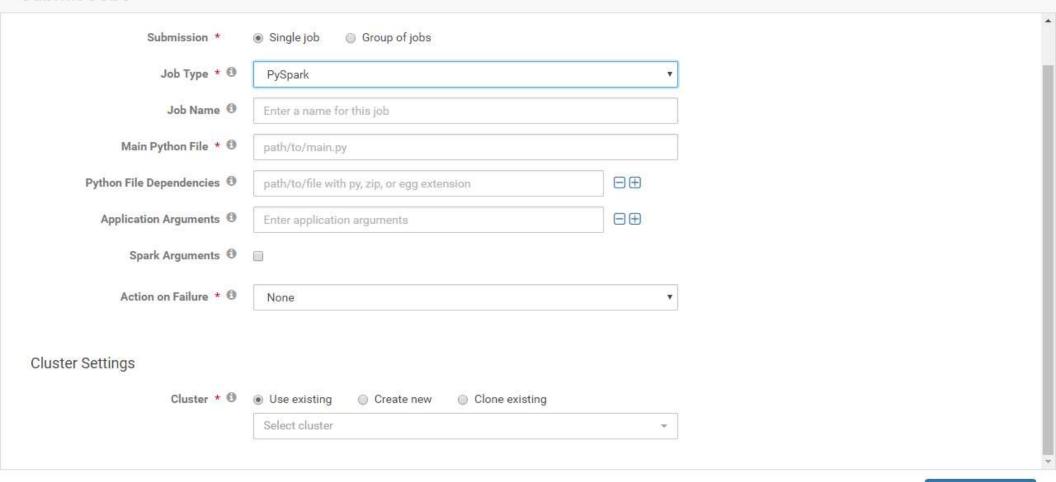


Cloudera Altus

▲ xavier-morera@bigdatainc.org ▼

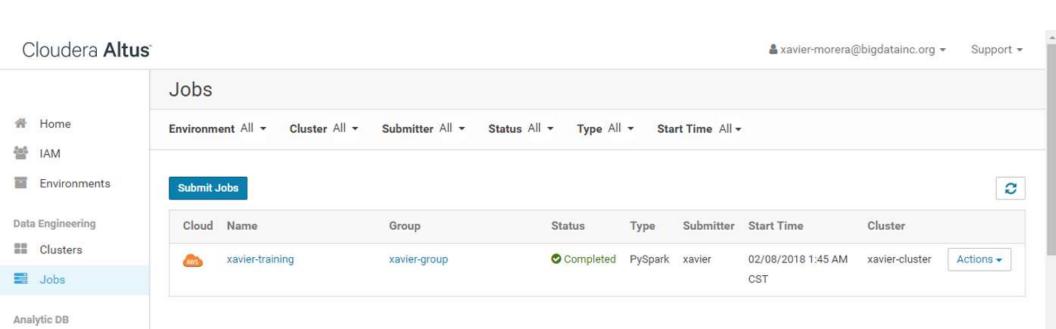
Support +

Submit Jobs



Cancel

Submit Jobs



«

& Clusters

Cloudera Data Science Workbench



Web application for working with Big Data

- Data Science
- From exploration to production

Work with R, Python & Scala

Self-Service Data Science for the Enterprise



View All Categories

About Cloudera Data Science Workbench

Release Notes

Known Issues and Limitations

Requirements and Supported Platforms

▼ Installation and Upgrade

Using Cloudera Manager

Using Packages

Getting Started

▶ User Guide

Managing the CDSW Service in Cloudera Manager

- Using Cloudera Distribution of Apache Spark 2
- ▶ Administration Guide
- ▶ Security Guide

Troubleshooting

Command Line Reference

Frequently Asked Questions (FAQs)

Installing and Upgrading Cloudera Data Science Workbench 1.3.x Using Cloudera Manager

This topic describes how to install and upgrade Cloudera Data Science Workbench using Cloudera Manager.

- Installing Cloudera Data Science Workbench 1.3.x Using Cloudera Manager
- Upgrading to the Latest Version of Cloudera Data Science Workbench 1.3.x

Installing Cloudera Data Science Workbench 1.3.x Using Cloudera Manager

Use the following steps to install Cloudera Data Science Workbench using Cloudera Manager.

- 1. Prerequisites
- 2. Install Cloudera Distribution of Apache Spark 2
- 3. Configure JAVA_HOME (Required for Spark 2.2)
- 4. Download and Install the Cloudera Data Science Workbench CSD
- 5. Install the Cloudera Data Science Workbench Parcel
- 6. Add the Cloudera Data Science Workbench Service
- 7. Create the Administrator Account
- 8. Next Steps

Prerequisites

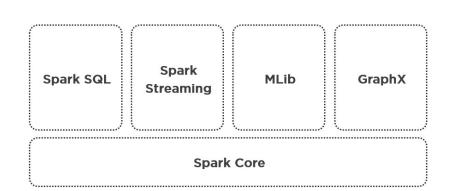
Before you begin installing Cloudera Data Science Workbench, make sure you have completed the steps to

http://tiny.bigdatainc.org/getcdsw

Streaming

Process live data streams in a scalable, high-throughput and fault tolerant way using



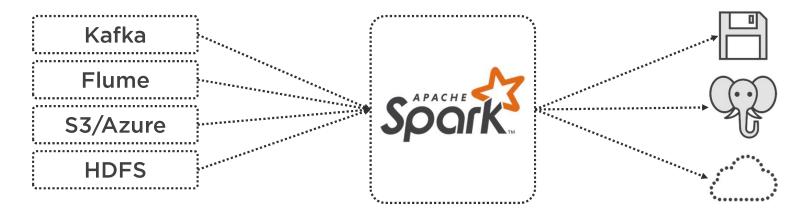


Spark Streaming

Process live data streams in a scalable, high-throughput and fault tolerant way using Apache Spark

Spark Streaming

Process live data streams in a scalable, high-throughput and fault tolerant way using Apache Spark



Structured Streaming

Process live data streams in a scalable, high-throughput and fault tolerant way using Apache Spark



TensorFlow



Open-source library

- Created by Google

Machine learning applications

- Deep learning

Use with Spark?



TensorFrames





TensorFlow on Spark DataFrames

Manipulate DataFrames

- TensorFlow programs

Machine Learning



Feedback Register a package Login Find a package Q

Tensorflow wrapper for DataFrames on Apache Spark

@databricks / **** (23)

TensorFrames (TensorFlow on Spark Dataframes) lets you manipulate Spark's DataFrames with TensorFlow programs.

This package provides Python bindings, a Scala DSL and a small runtime to express and run TensorFlow computation graphs.

Tags

2 tensorflow

How to [+]

Include this package in your Spark Applications using:

spark-shell, pyspark, or spark-submit

> \$SPARK_HOME/bin/spark-shell --packages databricks:tensorframes:0.2.9-s_2.11

Releases

Version: 0.2.9-s_2.11 (4be3b2 | zip | jar) / Date: 2017-09-05 / License: Apache-2.0

Version: 0.2.9-s_2.10 (4be3b2 | zip | jar) / Date: 2017-09-05 / License: Apache-2.0

Version: 0.2.9-rc3-s_2.11 (99e6d8 | zip | jar) / Date: 2017-08-31 / License: Apache-2.0

Version: 0.2.9-rc3-s_2.10 (99e6d8 | zip | jar) / Date: 2017-08-31 / License: Apache-2.0

Spark Packages is a community site hosting modules that are not part of Apache Spark. Your use of and access to this site is subject to the terms of use. Apache Spark and the Spark logo are trademarks of the Apache Software Foundation. This site is maintained as a community service by Databricks.

Thank You!



Xavier Morera
HELPING DEVELOPERS UNDERSTAND SEARCH & BIG DATA
@xmorera www.xaviermorera.com

