# APACHE SPARK FOR DATA SCIENTISTS

# Session 1: Introduction to Apache Spark

ie
SCHOOL OF
HUMAN SCIENCES
& TECHNOLOGY

# Recommended Prerequisite Knowledge in Python Language

- <u>Primitive Data Types</u> in Python
- <u>Flow Control</u> in Python (for, while  loops, if...else etc.)
- <u>Standard Functions</u> & <u>Lambda Expressions</u>
- Python <u>Collections</u> and optionally <u>Classes</u>
- <u>Pandas</u>, <u>Numpy</u>, <u>Scipy</u>
- <u>String Operations</u> & <u>Regular Expressions</u>
- Data Visualization in Python (matplotlib, bokeh, plotly etc.)
- <u>Basic Functional Methods in Python (map, filter, reduce,</u> <u>itertools</u> etc.)
- <u>Scikit</u>-learn (good to know it resembles syntactically to <u>Spark</u> MLib)

<u>https://www.codecademy.com/learn/learn-python</u>

https://lectures.quantecon.org/py/

https://runestone.academy/runestone/static/pythonds/index.html

ie

SCHOOL OF
HUMAN SCIENCES
& TECHNOLOGY

I.   Course Overview

II.  Spark: Background & Position in Big Data Analytics

III. Core Concepts & Challenges of Distributed Computing

IV.  Overview of Spark Components

V.   Conceptual Introduction to Spark Application

VI.  Appendix

ie
SCHOOL OF
HUMAN SCIENCES
& TECHNOLOGY

# Apache Spark Overview

*Apache Spark:*

- is a fast and general purpose big data processing engine.

- is an open source project incubated by Apache Software Foundation

https://spark.apache.org/

https://github.com/apache/spark

- is an unified engine with built-in modules for SQL, streaming, machine learning, graph processing & third-party packages.



| Spark SQL + DataFrames | Streaming | MLlib Machine Learning | GraphX Graph Computation |

**Apache Spark Core API**

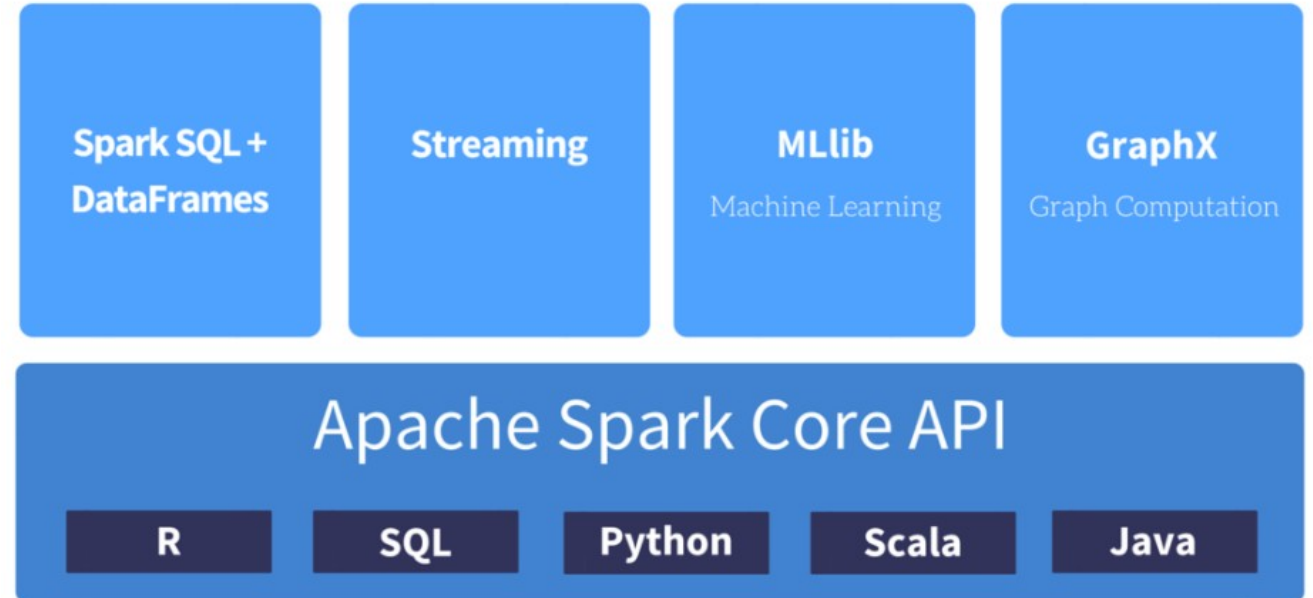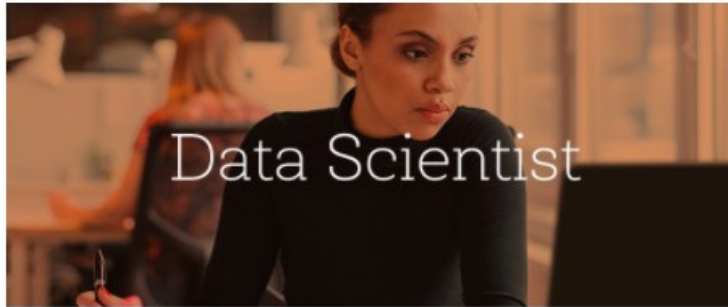| R | SQL | Python | Scala | Java |

Figure Ref: http://syntelli.com/data-capability/apache-spark/

# Common Use Cases of Apache Spark for Big Data

I. Extract-Transform-Load (ETL) operations

II. Predictive analytics and machine learning

III. Data access operations, such as interactive SQL queries and visualizations

IV. Text mining and processing

V. Real-time / Near real-time event processing

VI. Graph analysis applications

VII. Pattern recognition & Deep Learning

VIII. Recommendation engines

IX. and so on…

# Apache Spark for Data Scientists & Data Engineers

## Data Scientist

You have questions
and you need answers quickly.

- **Test hypotheses iteratively to converge on a solution.**
- **Leverage machine learning or graph processing algorithms to aid your investigation.**
- **Visually explore data and diagnose issues.**
- **Document your thinking and publish findings.**

## Data Engineer

Your team counts on you to prepare data
and deploy applications to production.

- **Build, operate, and manage infrastructure.**
- **Convert your existing codebase to a distributed setting.**
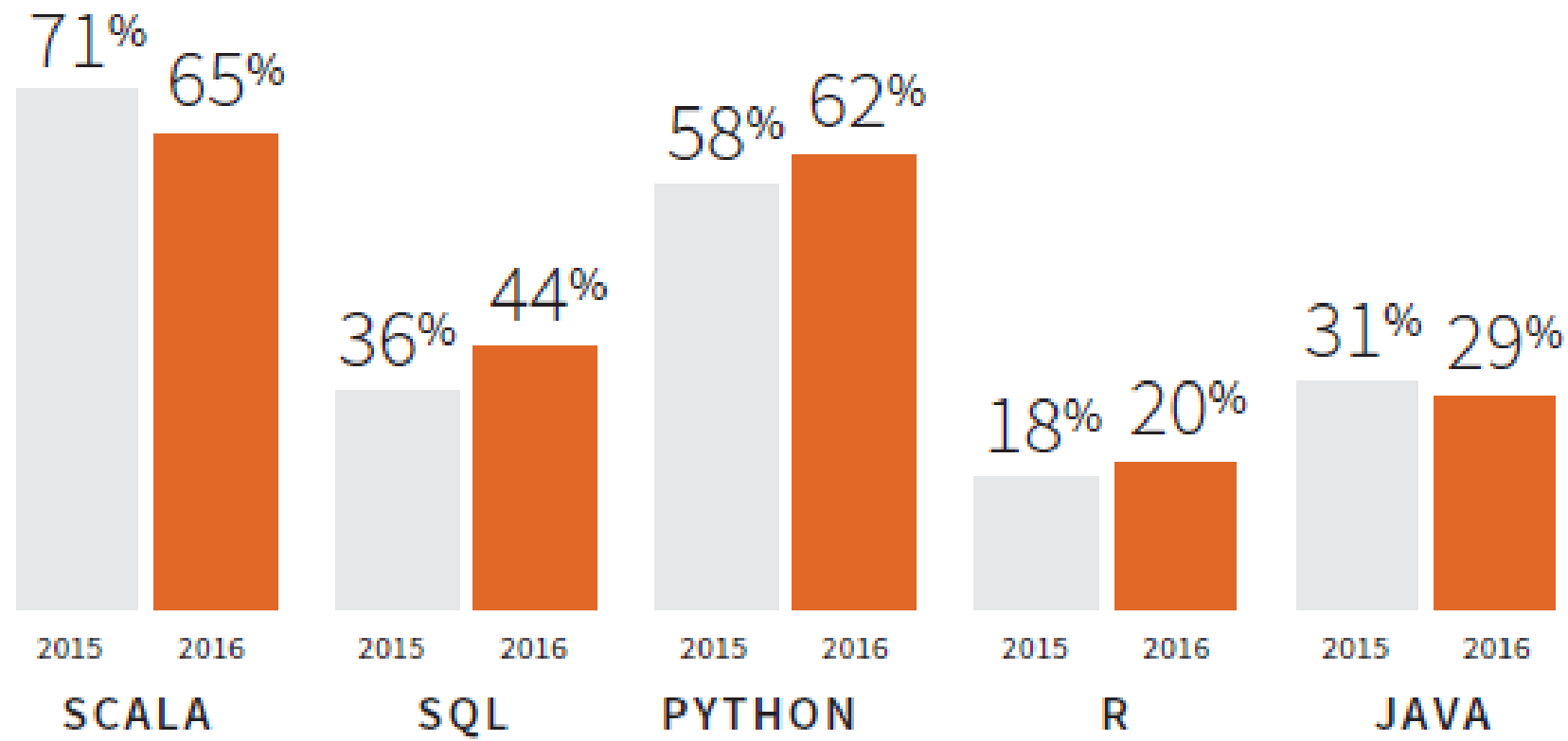- **Deploy production-quality data pipelines.**

python™    R    Spark SQL    Scala    Java    python™

Ref: https://databricks.com/solutions/by-role

SCHOOL OF
HUMAN SCIENCES
& TECHNOLOGY

# Spark Programming Languages

## LANGUAGES USED IN APACHE SPARK

*Respondents were allowed to select more than one language.*

**SCALA**
- 2015: 71%
- 2016: 65%

**SQL**
- 2015: 36%
- 2016: 44%

**PYTHON**
- 2015: 58%
- 2016: 62%

**R**
- 2015: 18%
- 2016: 20%

**JAVA**
- 2015: 31%
- 2016: 29%

March 2017 Ref:
https://www.slideshare.net/JohanPicard/a-short-introduction-to-spark-and-its-benefits

Master Business Analytics & Big Data

**ie**
SCHOOL OF
HUMAN SCIENCES
& TECHNOLOGY

# How Technologies Are Connected

Ref:https://insights.stackoverflow.com/survey/2018/#correlated-technologies

# Python vs. Scala for Spark Development (1/2)

Python is an <u>interpreted scripting language</u> and Pyspark API is a set of wrappers on Scala-Spark methods/functions, therefore there is always additional computational overhead.

It involves many widely used data science libraries like numpy, pandas, scipy, pyarrow, matplotlib, tensorflow, pixiedust, numba, spacy etc. which can be <u>mixed-in</u> <u>with Pyspark</u>.

Apache Spark is built with Scala, therefore it is:

✓ 10x faster than Pyspark

✓ New features are always initially available in Scala API

✓ More robust & extensive integration with other ecosyste tools like Kafka, Akka, Hadoop etc.

Scala & Java has less commonly used data science libraries like breeze, vegas, standford NLP, deeplearning4j etc.

# Python vs. Scala for Spark Development (2/2)

python™                                    Scala

It is dynamically typed, which limits type-safe Spark features and hardens debugging.

It is a statically typed & functional, complex functional & type-safe implementations are feasible and majority of bug can be caught in compile time.

Suits better to interactive data exploration, model training & evaluation and prototyping projects.

Suits better for production (live system) implementations & potentially a more robust option in any real-time analytics

Limited streaming features more relevant to prototyping rather than production applications & there are limitations of Pyspark Connectors to other Big Data Tools.
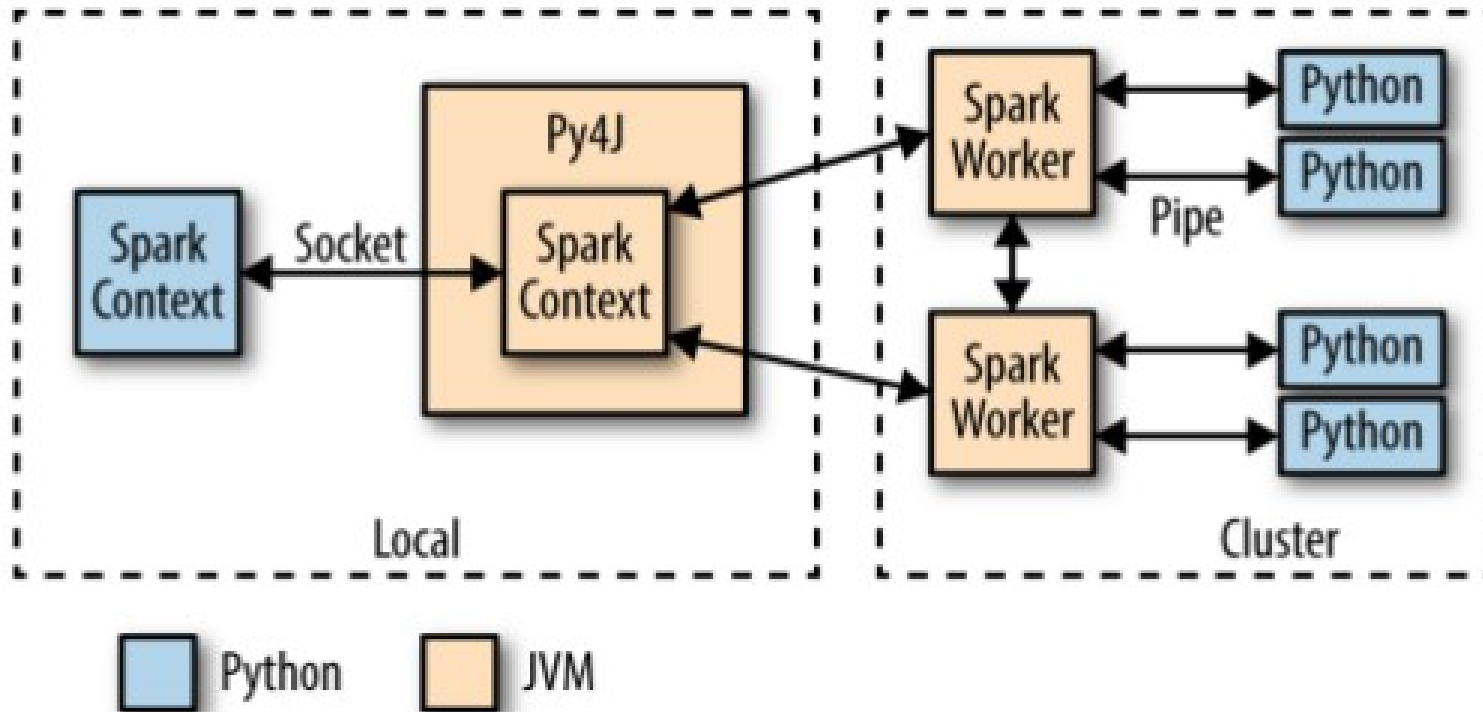
Advanced streaming & data product development capabilities, more variety of Scala-Spark Connectors to other Big Data Tools which are also developed in Scala or Java.

ie
SCHOOL OF
HUMAN SCIENCES
& TECHNOLOGY

# About Pyspark (Python API for Spark)



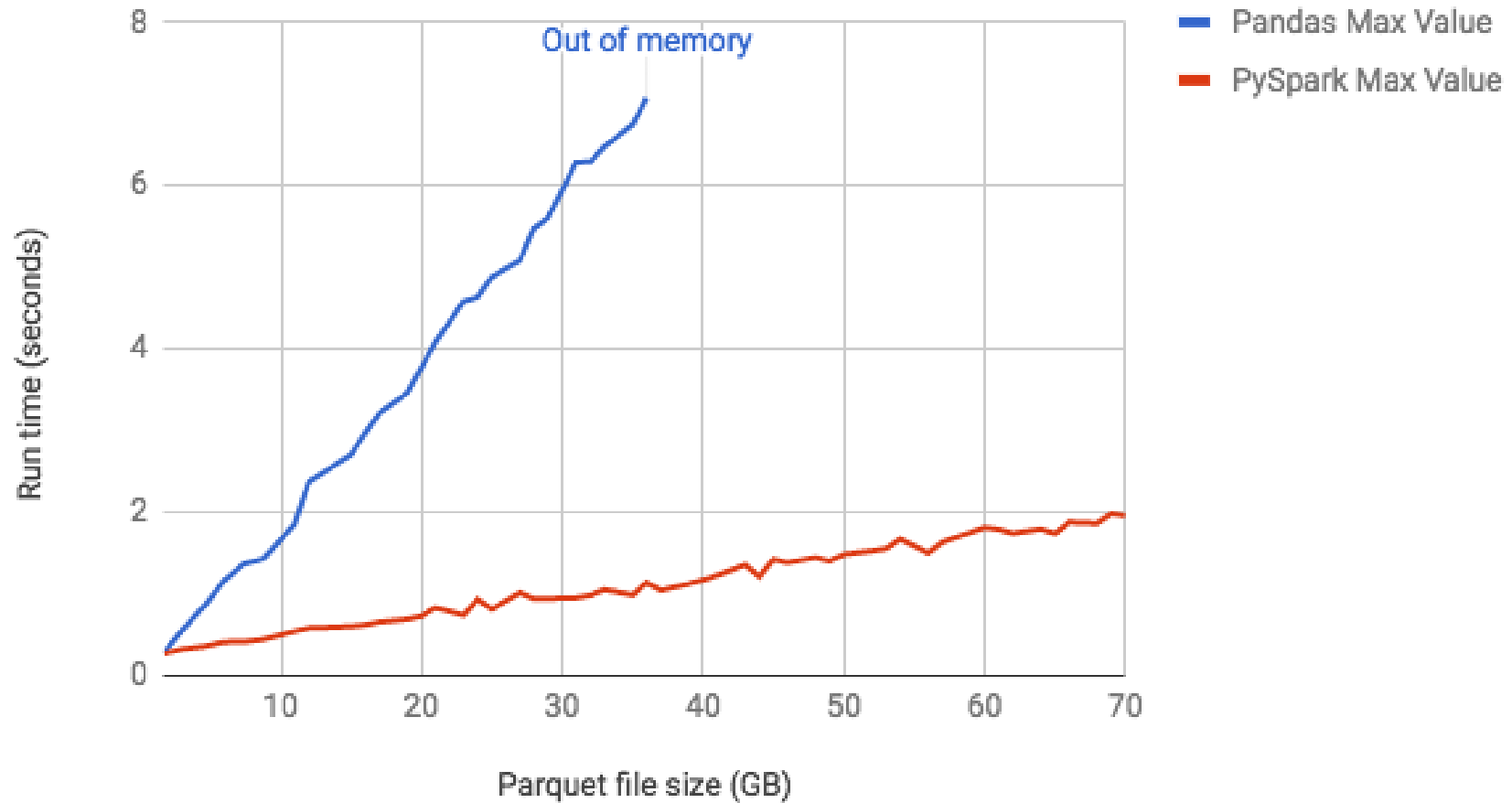Spark is internally developed in Scala. Scala Apps run on JVM (Java Virtual Machine) Threads similar to Java. List of JVM Languages

IMPORTANT: Pyspark API Documentation: https://spark.apache.org/docs/latest/api/python/index.html

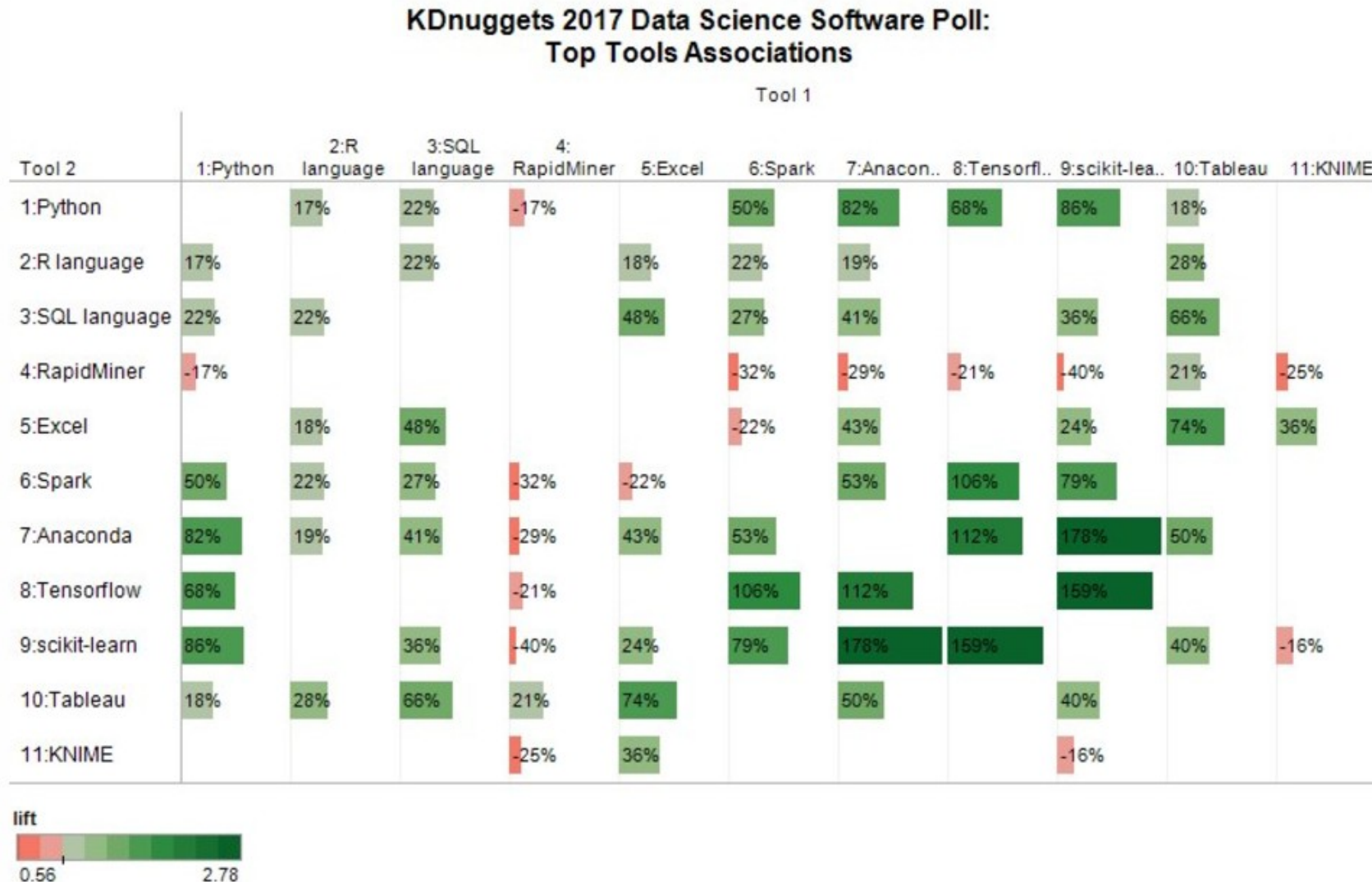# Benchmarking Pandas Vs Pyspark Dataframes on a Single Node Machine



Pandas VS PySpark: max value

Ref:
https://databricks.com/blog/2018/05/03/benchmarking-apache-spark-on-a-single-node-ma

**KDnuggets 2017 Data Science Software Poll:**
**Top Tools Associations**

Tool 1

| Tool 2 | 1:Python | 2:R language | 3:SQL language | 4:RapidMiner | 5:Excel | 6:Spark | 7:Anacon.. | 8:Tensorfl.. | 9:scikit-lea.. | 10:Tableau | 11:KNIME |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1:Python | | 17% | 22% | -17% | | 50% | 82% | 68% | 86% | 18% | |
| 2:R language | 17% | | 22% | | 18% | 22% | 19% | | | 28% | |
| 3:SQL language | 22% | 22% | | | 48% | 27% | 41% | | 36% | 66% | |
| 4:RapidMiner | -17% | | | | | -32% | -29% | -21% | -40% | 21% | -25% |
| 5:Excel | | 18% | 48% | | | -22% | 43% | | 24% | 74% | 36% |
| 6:Spark | 50% | 22% | 27% | -32% | -22% | | 53% | 106% | 79% | | |
| 7:Anaconda | 82% | 19% | 41% | -29% | 43% | 53% | | 112% | 178% | 50% | |
| 8:Tensorflow | 68% | | | -21% | | 106% | 112% | | 159% | | |
| 9:scikit-learn | 86% | | 36% | -40% | 24% | 79% | 178% | 159% | | 40% | -16% |
| 10:Tableau | 18% | 28% | 66% | 21% | 74% | | 50% | | 40% | | |
| 11:KNIME | | | | -25% | 36% | | | | -16% | | |

lift

0.56   2.78

Ref:

http://www.kdnuggets.com/2017/06/ecosystem-data-science-machine-learning-software.html

ie
SCHOOL OF
HUMAN SCIENCES
& TECHNOLOGY

# Top Data Science Tools Associations (2017 KDnuggets Survey)



Bar length is Bias_Py_R as defined above, bar height is the popularity of the tool.

Ref:
https://www.kdnuggets.com/2018/05/poll-tools-analytics-data-science-machine-learning-results.html

# Deep Learning vs Spark/Hadoop affinity for top Data Science Tools(2017 KDnuggets)



Circle size corresponds to tool share of use, and color to Python (blue) vs R (Orange) bias.
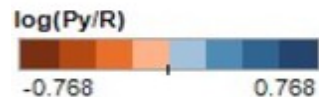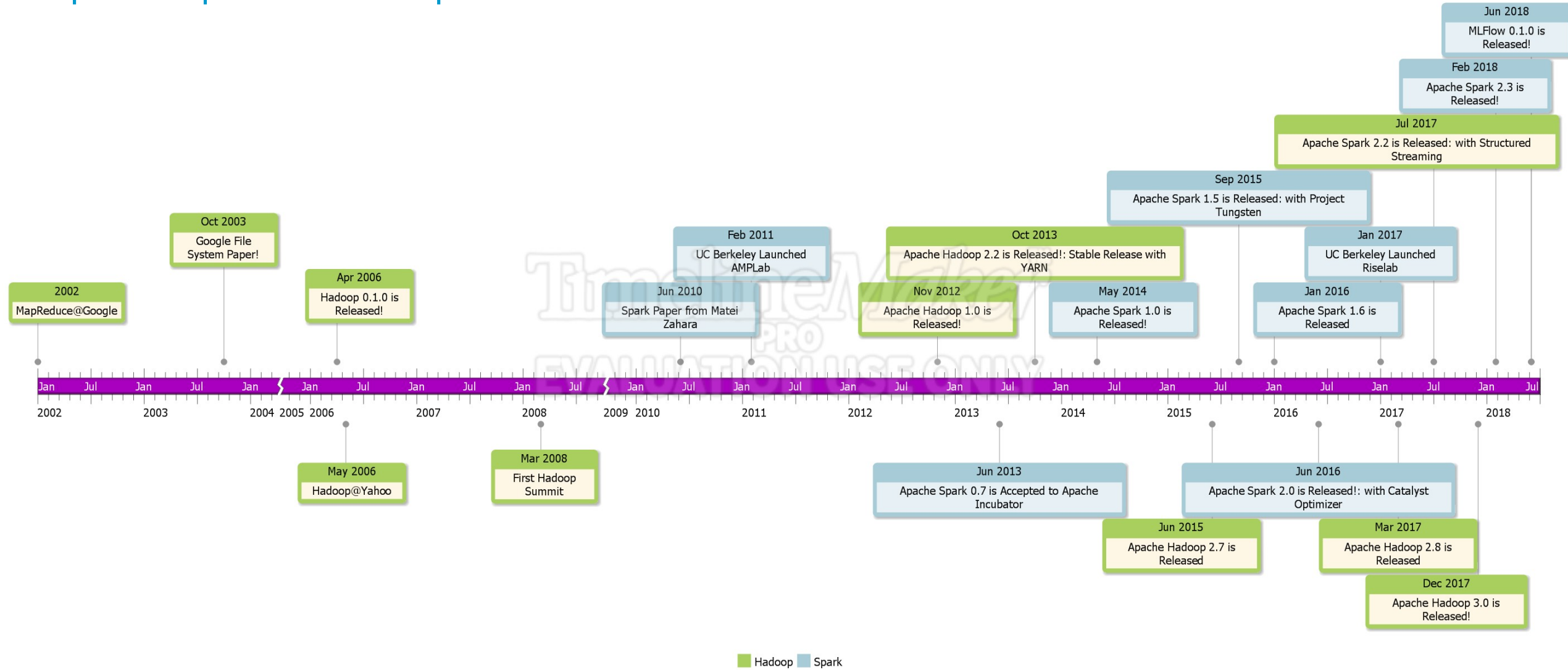
Ref:
http://www.kdnuggets.com/2017/06/ecosystem-data-science-machine-learning-software.html/2

# Apache Spark & Hadoop Timeline



Jun 2018
MLFlow 0.1.0 is Released!

Feb 2018
Apache Spark 2.3 is Released!

Jul 2017
Apache Spark 2.2 is Released: with Structured Streaming

Sep 2015
Apache Spark 1.5 is Released: with Project Tungsten

Oct 2003
Google File System Paper!

Feb 2011
UC Berkeley Launched AMPLab

Oct 2013
Apache Hadoop 2.2 is Released!: Stable Release with YARN

Jan 2017
UC Berkeley Launched Riselab

Apr 2006
Hadoop 0.1.0 is Released!

Jun 2010
Spark Paper from Matei Zahara

Nov 2012
Apache Hadoop 1.0 is Released!

May 2014
Apache Spark 1.0 is Released!

Jan 2016
Apache Spark 1.6 is Released

2002
MapReduce@Google

2002 2003 2004 2005 2006 2007 2008 2009 2010 2011 2012 2013 2014 2015 2016 2017 2018

May 2006
Hadoop@Yahoo

Mar 2008
First Hadoop Summit

Jun 2013
Apache Spark 0.7 is Accepted to Apache Incubator

Jun 2016
Apache Spark 2.0 is Released!: with Catalyst Optimizer

Jun 2015
Apache Hadoop 2.7 is Released

Mar 2017
Apache Hadoop 2.8 is Released

Dec 2017
Apache Hadoop 3.0 is Released!

Hadoop  Spark

Created with Timeline Maker Pro v4. Produced on Aug 11 2018.

ie
SCHOOL OF HUMAN SCIENCES & TECHNOLOGY