

ILLINOIS DATA SCIENCE INITIATIVE

TECHNICAL REPORTS

PySpark on Python 3: Configuration and Package Mangement Guide

Author:
Benjamin Congdon

March 1, 2017

PySpark on Python 3: Configuration and Package Management Guide

BENJAMIN CONGDON¹ AND THE BIG DOG²

¹*National Center For Supercomputing Applications (NCSA)*

²*Laboratory for Computation, Data, and Machine Learning*

Compiled March 1, 2017

PySpark requires Python 2.6 or later, and does not officially support Python 3. However, with only minor configuration changes, it is possible to successfully run a Spark cluster using Python 3 for both PySpark drivers and workers.

<https://github.com/lcdm-uiuc>

INTRODUCTION

ASSUMPTIONS

This technical report will make the assumption that you have already set up a Spark Cluster with a CentOS worker base image.

EXAMPLES OF ARTICLE COMPONENTS

The sections below show examples of different article components.

FUNDING INFORMATION

ACKNOWLEDGMENTS

SUPPLEMENTAL DOCUMENTS

REFERENCES