# cloudera®

# Spark Guide

**Cloudera, Inc.**
**1001 Page Mill Road, Bldg 3**
**Palo Alto, CA 94304**
**info@cloudera.com**
**US: 1-888-789-1488**
**Intl: 1-650-362-0488**
**www.cloudera.com**

**Release Information**

Version: CDH 5.6.x
Date: August 24, 2017

# Table of Contents