# PySpark Cheat Sheet: Spark in Python
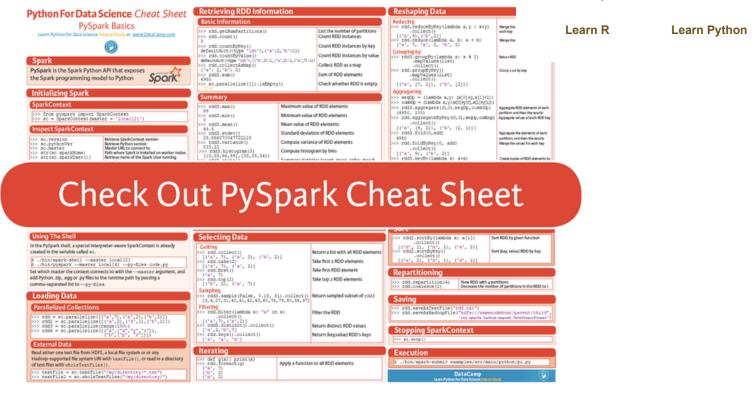
March 21st, 2017 in Python

Karlijn Willems

Apache Spark is generally known as a fast, general and open-source engine for big data processing, with built-in modules for streaming, SQL, machine learning and graph processing. It allows you to speed analytic applications up to 100 times faster compared to technologies on the market today. You can interface Spark with Python through "PySpark". This is the Spark Python API exposes the Spark programming model to Python.

Even though working with Spark will remind you in many ways of working with Pandas DataFrames, you'll also see that it can be tough getting familiar with all the functions that you can use to query, transform, inspect, ... your data. What's more, if you've never worked with any other programming language or if you're new to the field, it might be hard to distinguish between RDD operations.

Even though the documentation is very elaborate, it never hurts to have a cheat sheet by your side, especially when you're just getting into it.

This PySpark cheat sheet covers the basics, from initializing Spark and loading your data, to retrieving RDD information, sorting, filtering and sampling your data. But that's not all. You'll also see that topics such as repartitioning, iterating, merging, saving your data and stopping the SparkContext are included in the cheat sheet.

Are you hungry for more? Don't miss our other Python cheat sheets for data science that cover top
[basics](#), [Numpy](#), [Pandas](#), [Pandas Data Wrangling](#) and much more!

What do you think?

🏷    Python

# Up Next



## New Course! Supervised Learning in R: Classification

R Programming                                                                          79 views

This beginner-level introduction to machine learning covers four of the most common classification algorithms. You will come away with a b…

September 27th, 2017 in Blog

## New Course: Case Studies in Statistical Thinking!

Python

1,139 views

Hone your applied data science skills in Python by doing a variety of case studies across multiple disciplines. Use data science to solve …

September 20th, 2017 in Blog

## Jupyter Notebook Cheat Sheet

Python

19,148 views

This Jupyter Notebook cheat sheet will help you to find your way around the well-known Jupyter Notebook App, a subproject of Project Jupyt…

September 19th, 2017 in Blog

# Comments

**vaidas-armonas**

Hi Karlijn, thanks for sharing this. But *why* an RDD cheat sheet? RDDs should be left for the most exotic of uses while most DataSet APIs where we get Tungsten performance benefits and more familiar interface. I would love to hear the rationale and to share with the analysts in my company :)

06/28/17 6:43 AM |

    **karlijn**

    Hi vaidas-armonas! Thanks for your message! The idea behind this is that when you start out with Spark, you usually first cover the basic building blocks, which are still the RDDs, even though you might not use them often in daily practice. I hope this reasoning makes somewhat sense to you; You might also have seen that I recently added a cheat sheet to work with Spark SQL (DataFrames) on the community. You can find it here: https://www.datacamp.com/community/blog/pyspark-sql-cheat-sheet :)

    07/26/17 7:45 AM |

**vaidas-armonas**

Hi Karlijn, thanks for sharing this. But *why* an RDD cheat sheet? RDDs should be left for the most exotic of uses while most use cases should be covered by DataFrame / DataSet APIs where we get Tungsten performance benefits and more familiar interface. I would love to hear the rationale and maybe a cheat sheet for PySpark DataFrames to share with the analysts in my company :)

06/28/17 6:42 AM |

    **karlijn**

    Hi vaidas-armonas, I replied to your message above :)

    07/26/17 7:47 AM |

**alfredo-g-marquez**

This is great! Can't wait to see the cheat sheet for SparkR.

03/24/17 4:57 AM |

**benmainye**

I have never used PySpark in my work. But, it seems awesome. Thanks for the cheatsheet.

03/22/17 4:23 AM |

    **karlijn**

    Hi there! Thanks for writing. Spark is really wonderful when you're working with big data. Definitely give PySpark a go and let us know what you thought of it!

    03/22/17 7:49 AM |

**mannurulz**