in 3 Minutes

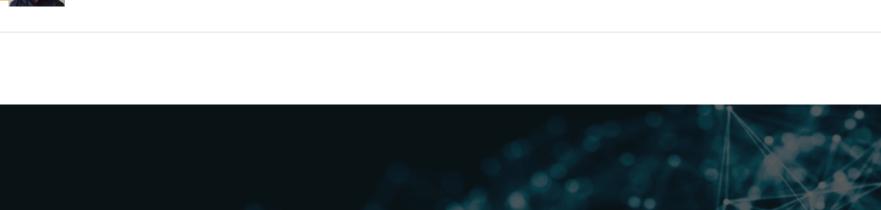
JANUARY 19, 2020

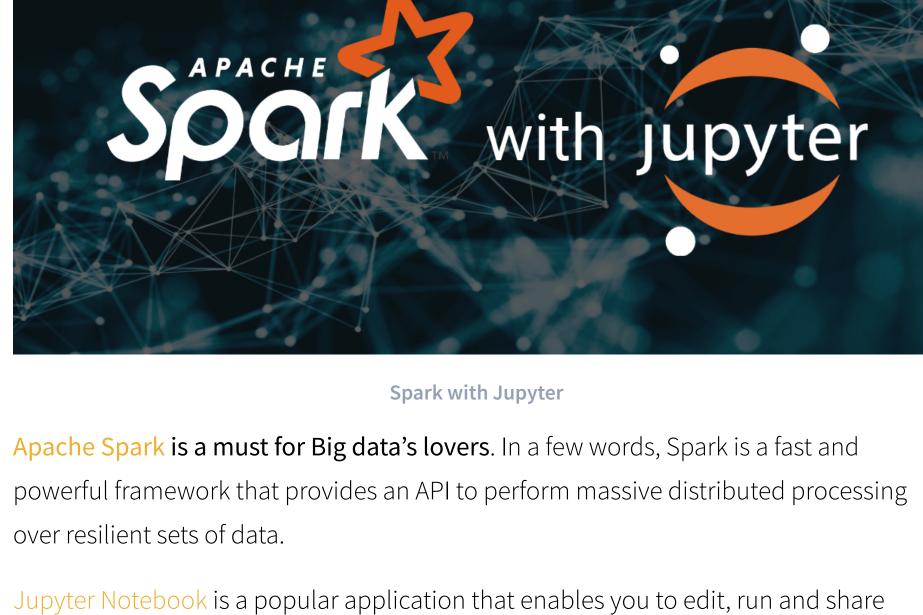
## Spark is a fast and powerful framework. **Charles Bochet**

5 min read

How to install PySpark and Jupyter Notebook

Data Scientist



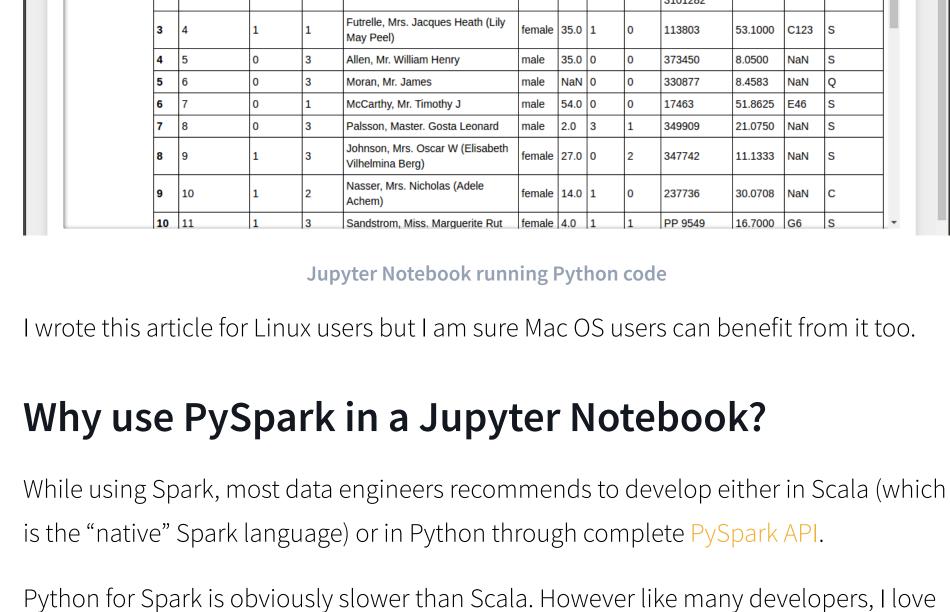


Python code into a web view. It allows you to modify and re-execute parts of your code

in a very flexible way. That's why Jupyter is a great tool to test and prototype programs.

① localhost:8888/notebooks/dev/sicara/titanic/jupyter/main.py.ipynb

Python [conda root] O In [1]: import pandas as pd import numpy as np # create data frame containing your data, each column can be accessed # by df['column data = pd.read csv('../data/train.csv') In [2]: data STON/O2 Heikkinen, Miss, Laina female 26.0 0 7.9250 NaN 3101282



favorites libraries. In my opinion, Python is the perfect language for prototyping in Big Data/Machine Learning fields.

To learn more about Python vs. Scala pro and cons for Spark context, please refer to this interesting article: Scala vs. Python for Apache Spark.

If you prefer to develop in Scala, you will find many alternatives on the following

Python because it's flexible, robust, easy to learn, and benefits from all my

github repository: alexarchambault/jupyter-scala

Install pySpark Before installing pySpark, you must have Python and Spark installed. I am using Python 3 in the following examples but you can easily adapt them to Python 2. Go to

Then, visit the Spark downloads page. Select the latest Spark release, a prebuilt

Create a symbolic link:

export SPARK\_HOME=/opt/spark

export PATH=\\$SPARK\_HOME/bin:\\$PATH

package for Hadoop, and download it directly.

Unzip it and move it to your /opt folder:

Finally, tell your bash (or zsh, etc.) where to find Spark. To do so, configure your \$PATH variables by adding the following lines in your ~/.bashrc (or ~/.zshrc) file:

This way, you will be able to download and use multiple Spark versions.

\$ tar -xzf spark-1.2.0-bin-hadoop2.4.tgz

\$ mv spark-1.2.0-bin-hadoop2.4 /opt/spark-1.2.0

Install Jupyter Notebook

You can run a regular jupyter notebook by typing:

## \$ pysparkWelcome to

SparkSession available as 'spark'.

(I bet you understand what it does!)

pi = 4 \* count / num samples print(pi) sc.stop() The output will probably be around 3.14. **PySpark in Jupyter** 

export PYSPARK\_DRIVER\_PYTHON=jupyter export PYSPARK\_DRIVER\_PYTHON OPTS='notebook' Restart your terminal and launch PySpark again:

Copy and paste our Pi calculation script and run it by pressing Shift + Enter.

Jupyter Notebook: Pi Calculation script

Method 2 — FindSpark package There is another and more generalized way to use PySpark in a Jupyter Notebook:

You are now able to run PySpark in a Jupyter Notebook:)

findspark.init()

sc = pyspark.SparkContext(appName="Pi")

x, y = random.random(), random.random()

I hope this 3-minutes guide will help you easily getting started with Python and Spark.

Thanks to Pierre-Henri Cumenge, Antoine Toubhans, Adil Baaj, Vincent Quagliaro, and Adrien Lina.

And if you want to tackle some bigger challenges, don't miss out the more evolved

JupyterLab environnement or the PyCharm integration of jupyter notebooks.

the Python official website to install it. I also encourage you to set up a virtualenv To install Spark, make sure you have Java 8 or higher installed on your computer.

Now, let's get started.

\$ ln -s /opt/spark-1.2.0 /opt/spark

Install Jupyter notebook: \$ pip install jupyter

Let's check if PySpark is properly installed without using Jupyter Notebook first. You may need to restart your terminal to be able to run PySpark. Run:

Your first Python program on Spark

\$ jupyter notebook

17:53:06)

import random

def inside(p):

 $num_samples = 100000000$ 

return x\*x + y\*y < 1

count = sc.parallelize(range(0,

num samples)).filter(inside).count()

>>> It seems to be a good start! Run the following program:

x, y = random.random(), random.random()

Using  $\overline{Py}$ thon version 3.5.2 (default, Jul 2 2016

There are two ways to get PySpark available in a Jupyter Notebook: Configure PySpark driver to use Jupyter Notebook: running pyspark will automatically open a Jupyter Notebook Load a regular Jupyter Notebook and load PySpark using findSpark package First option is quicker but specific to Jupyter Notebook, second option is a broader approach to get PySpark available in your favorite IDE. Method 1 — Configure PySpark driver Update PySpark driver environment variables: add these lines to

\$ pyspark Now, this command should start a Jupyter Notebook in your web browser. Create a new notebook by clicking on 'New' > 'Notebooks Python [default]'.

Jupyter Notebook: Pi Calculation script

Done!

To install findspark:

\$ pip install findspark

Launch a regular Jupyter Notebook:

\$ jupyter notebook

import pyspark import random

def inside(p):

print(pi)

sc.stop()

num samples = 100000000

return x\*x + y\*y < 1

count = sc.parallelize(range(0,

pi = 4 \* count / num samples

num\_samples)).filter(inside).count()

your ~/.bashrc (or ~/.zshrc) file.

use findSpark package to make a Spark Context available in your code. findSpark package is not specific to Jupyter Notebook, you can use this trick in your favorite IDE too.

Create a new Python [default] notebook and write the following script: import findspark

The output should be: Jupyter Notebook: Pi calculation Jupyter Notebook: Pi calculation Here are a few resources if you want to go the extra mile: https://www.dezyre.com/article/scala-vs-python-for-apache-spark/213 http://queirozf.com/entries/comparing-interactive-solutions-for-runningapache-spark-zeppelin-spark-notebook-and-jupyter-scala http://spark.apache.org/docs/latest/api/python/index.html https://github.com/jadianes/spark-py-notebooksl

how to build a successful ai poc thief bridge

Arnault 9 min read Data Scientist

How To Build A Successful AI

Turn Your Artificial Intelligence Ideas Into

PoC

**Working Software** 

Antoine 7 min read Data Scientist

How to Perform Fraud

detection with Graph Analysis.

Page Rank

Legal infos

Terms

Privacy

**Detection with Personalized** 

This article shows how to perform fraud

9 min read

y in 🖸

Image Registration: From SIFT

How the field has evolved from OpenCV

to Deep Learning

to Neural Networks.

Sicara © 2016-2019 All rights reserved

**SICARA** 

Contact

Sicara SAS

75017 Paris

contact@sicara.ai

48 Bld des Batignolles

**Contact us** 

Your email