

# Spark – Big Data Processing

## Aula 6



**Semantix<sup>®</sup>**

All about data

# Quem sou eu?



## Rodrigo Augusto Rebouças

Engenheiro de dados da Semantix  
Instrutor do Semantix Academy

### Contatos

[rodrigo.augusto@semantix.com.br](mailto:rodrigo.augusto@semantix.com.br)  
[linkedin.com/in/rodrigo-reboucas](https://www.linkedin.com/in/rodrigo-reboucas)



# Spark application

# Spark shell x Spark applications

- Spark shell
  - Exploração e manipulação dos dados
- Spark applications
  - Rodar os programas de forma independente
  - Jobs para ETL e streaming
  - Criação de objetos
    - SparkSession – spark (spark SQL)
    - SparkContext – sc

# Rodar uma Spark application

- `spark-submit --class NameList MyJarFile.jar people.json namelist/`
- Opções submit
  - master: local, yarn, mesos ou spark standalone
  - jars: adicionar arquivos jar
  - py-files: lista de arquivos em .py, .zip ou .egg
  - driver-java-options: parâmetros para o driver JVM
  - deploy-mode: client ou cluster
  - driver-memory: Memória alocada para o spark driver (1G)
  - executor-memory: Memória alocada para a aplicação
  - num-executors: Número de executores para iniciar com a aplicação
  - driver-cores: Número de cores alocados para o spark driver
  - queue: Rodar na fila do Yarn
  - help

# Instalação IDE



# Intalação

- IntelliJ
  - Plugins:
    - Scala – 2.11
    - SBT (Preferível instalação separada)
  - Java 8
    - <https://www.oracle.com/br/java/technologies/javase/javase8u211-later-archive-downloads.html>
- PyCharm
  - Python 3

# IntelliJ

- <https://www.jetbrains.com/pt-br/idea/>

## Baixar IntelliJ IDEA

Windows

macOS

Linux

### Ultimate

Para desenvolvimento Web e corporativo

Baixar

.exe



Avaliação gratuita por 30 dias

### Community

Para desenvolvimento JVM e Android

Baixar

.exe

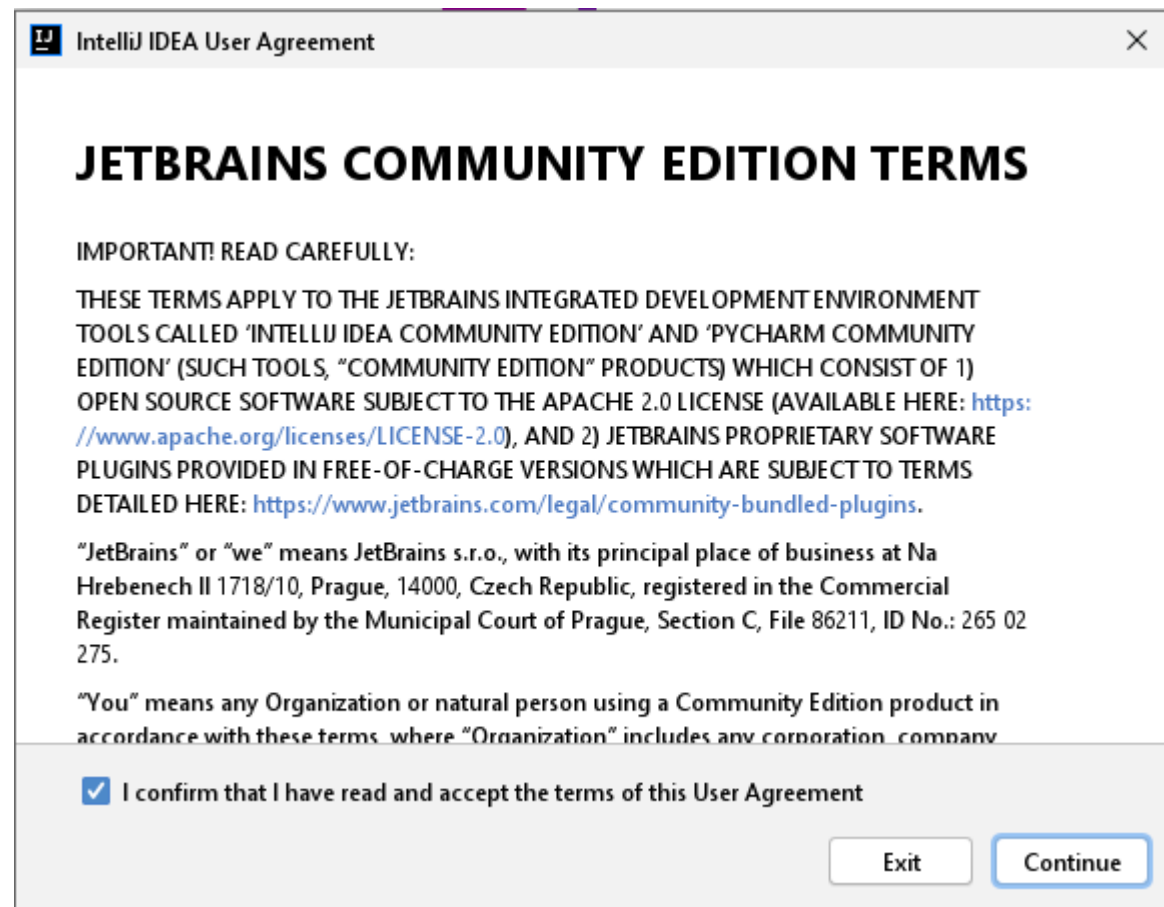
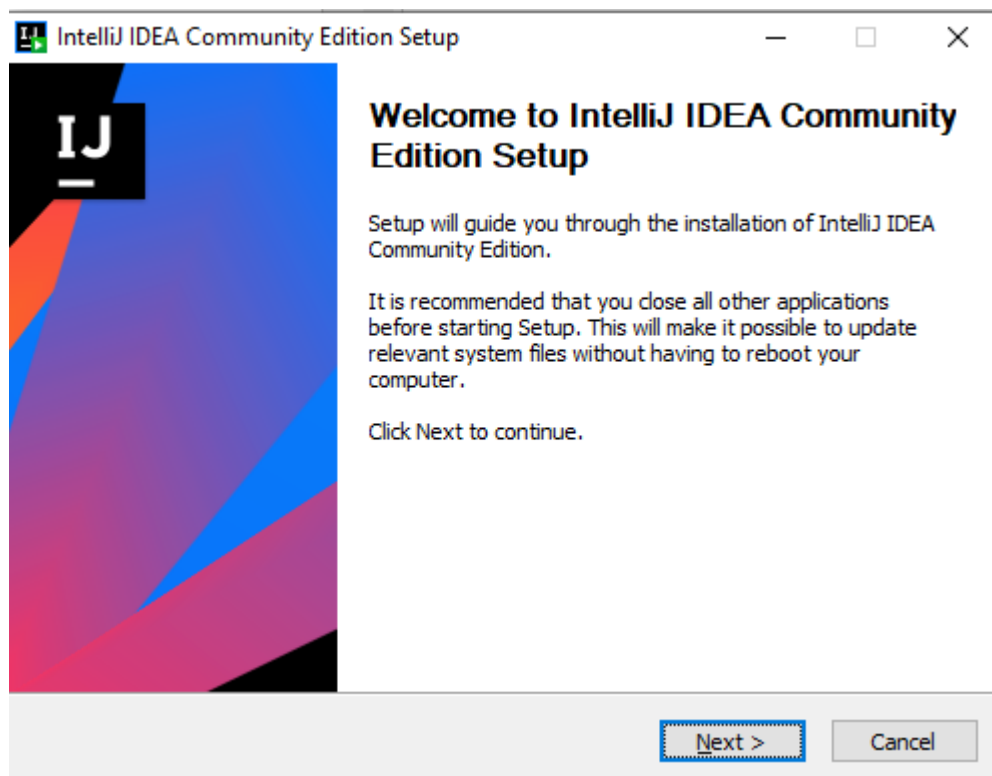


Gratuito, com base em open source



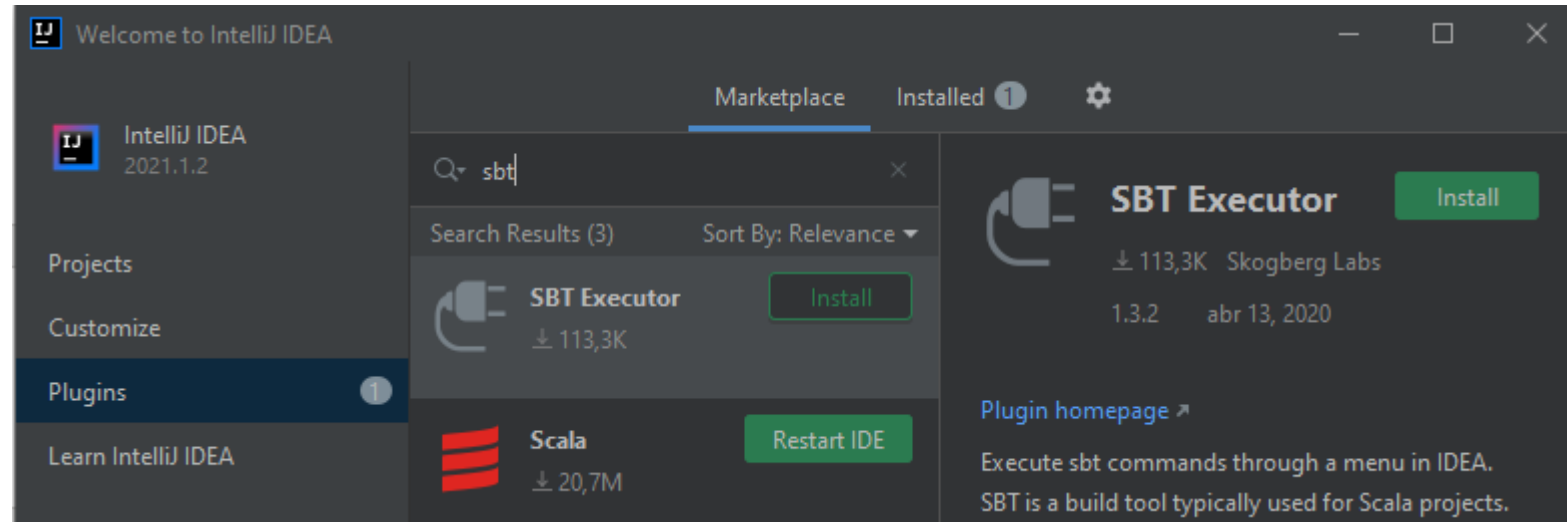
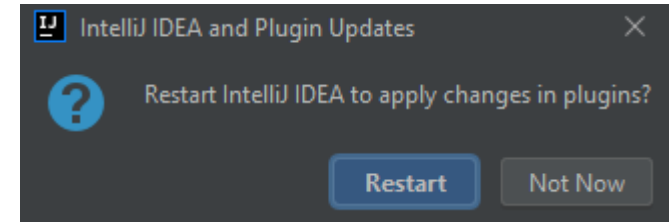
# IntelliJ

- Espaço em disco
  - 1,6 GB



# IntelliJ

- Plugins
  - Scala
  - SBT Executor



# Python

- Windows
  - `python --version`
- Linux
  - `python3 --version`
- Instalação
  - <https://www.python.org/>
- Anaconda (Opcional)



# PyCharm

- <https://www.jetbrains.com/pt-br/pycharm/>

## Baixar PyCharm

Windows

macOS

Linux

### Professional

Para desenvolvimento Web com Python e desenvolvimento científico. Com suporte para HTML, JS e SQL.

Baixar

Avaliação gratuita

### Community

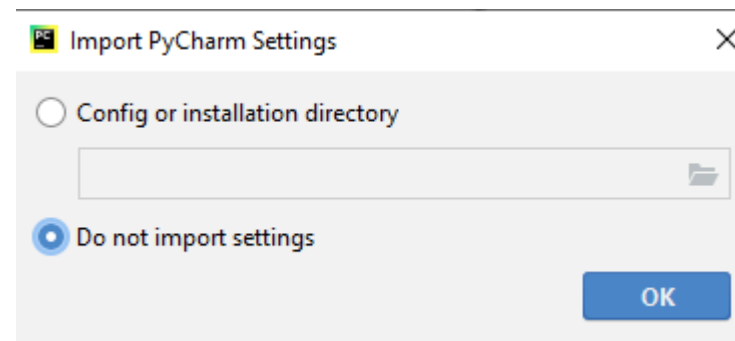
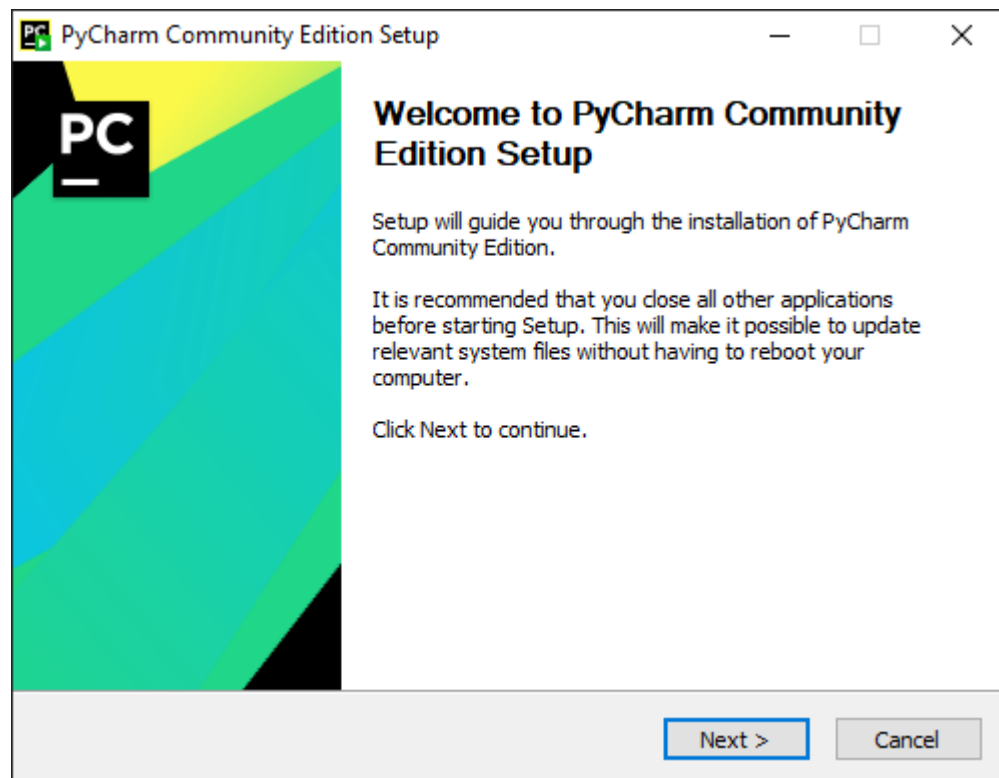
Para o autêntico desenvolvimento Python

Baixar

Open source gratuito

# PyCharm

- Espaço em disco
  - 1 GB



# Build Spark

# Build Spark

- IntelliJ
  - sbt package
    - <https://www.scala-sbt.org/sbt-native-packager/gettingstarted.html#setup>
- Pycharm
  - pip install
    - <https://packaging.python.org/tutorials/installing-packages/>
- <https://spark.apache.org/docs/latest/building-spark.html>



# Semantix<sup>®</sup>

All about data

[contato@semantix.com.br](mailto:contato@semantix.com.br)

[www.semantix.com.br](http://www.semantix.com.br)