



# 《Spark编程基础（Python版）》

教材官网：<http://dblab.xmu.edu.cn/post/spark-python/>

温馨提示：编辑幻灯片母版，可以修改每页PPT的厦大校徽和底部文字

## 第3章 Spark环境搭建和使用方法

（PPT版本号：2020年1月版）



扫一扫访问教材官网

林子雨

厦门大学计算机科学系

E-mail: [ziyulin@xmu.edu.cn](mailto:ziyulin@xmu.edu.cn) ▶▶

主页：<http://dblab.xmu.edu.cn/post/linziyu>

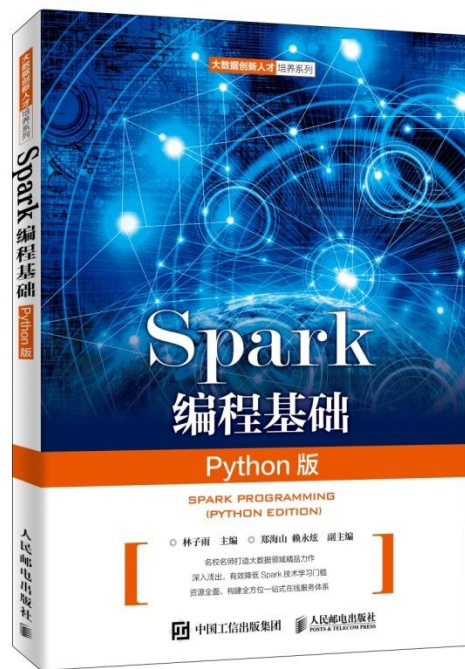




# 课程教材

## 林子雨，郑海山，赖永炫 编著 《Spark编程基础（Python版）》

教材官网：<http://dbllab.xmu.edu.cn/post/spark-python/>  
ISBN:978-7-115-52439-3 人民邮电出版社



本书以Python作为开发Spark应用程序的编程语言，系统介绍了Spark编程的基础知识。全书共8章，内容包括大数据技术概述、Spark的设计与运行原理、Spark环境搭建和使用方法、RDD编程、Spark SQL、Spark Streaming、Structured Streaming、Spark MLlib等。本书每个章节都安排了入门级的编程实践操作，以便读者更好地学习和掌握Spark编程方法。本书官网免费提供了全套的在线教学资源，包括讲义PPT、习题、源代码、软件、数据集、授课视频、上机实验指南等。



# 提纲

- 3.1 安装Spark
- 3.2 在pyspark中运行代码
- 3.3 开发Spark独立应用程序
- 3.4 Spark集群环境搭建
- 3.5 在集群上运行Spark应用程序



高校大数据课程

公共服务平台

百度搜索厦门大学数据库实验室网站访问平台





# 3.1 安装Spark

Spark的安装详细过程，请参考厦门大学数据库实验室建设的高校大数据课程公共服务平台上的技术博客：

《**Spark2.4.0**入门：**Spark**的安装和使用》

博客地址：<http://dblab.xmu.edu.cn/blog/1307-2/>

3.1.1 基础环境

3.1.2 下载安装文件

3.1.3 配置相关文件

3.1.4 Spark和Hadoop的交互



高校大数据课程

公共服务平台

平台每年访问量超过100万次



## 3.1.1 基础环境

- 安装Spark之前需要安装Linux系统、Java环境（Java8或JDK1.8以上版本）和Hadoop环境
- 如果没有安装Hadoop，请访问厦门大学数据库实验室建设的高校大数据课程公共服务平台，找到“Hadoop安装教程\_单机/伪分布式配置\_Hadoop2.6.0/Ubuntu14.04”（适用于Hadoop2.7.1/Ubuntu16.04），依照教程学习安装即可
- 注意，在这个Hadoop安装教程中，就包含了Java的安装，所以，按照这个教程，就可以完成JDK和Hadoop这二者的安装

Hadoop安装教程地址：<http://dblab.xmu.edu.cn/blog/install-hadoop/>





## 3.1.2 下载安装文件

- Spark安装包下载地址: <http://spark.apache.org>

进入下载页面后, 点击主页右侧的“Download Spark”按钮进入下载页面, 下载页面中提供了几个下载选项, 主要是Spark release及Package type的选择, 如下图所示。第1项Spark release一般默认选择最新的发行版本, 截至2019年1月份的最新版本为2.4.0(本教程采用2.4.0)。第2项package type则选择“Pre-build with user-provided Hadoop [can use with most Hadoop distributions]”, 可适用于多数Hadoop版本。选择好之后, 再点击第4项给出的链接就可以下载Spark了。

### Download Apache Spark™

1. Choose a Spark release:
2. Choose a package type:
3. Download Spark: [spark-2.4.0-bin-without-hadoop.tgz](#)
4. Verify this release using the [2.4.0 signatures and checksums](#) and [project release KEYS](#).

*Note: Starting version 2.0, Spark is built with Scala 2.11 by default. Scala 2.10 users should download the Spark source package and build with Scala 2.10 support.*



## 3.1.2 下载安装文件

- 解压安装包spark-2.4.0-bin-without-hadoop.tgz至路径 /usr/local:

```
$ sudo tar -zxf ~/下载/spark-2.4.0-bin-without-hadoop.tgz -C /usr/local/  
$ cd /usr/local  
$ sudo mv ./spark-2.4.0-bin-without-hadoop/ ./spark # 更改文件夹名  
$ sudo chown -R hadoop ./spark # 此处的 hadoop 为系统用户名
```



## 3.1.3 配置相关文件

- 配置Spark 的classpath

```
$ cd /usr/local/spark  
$ cp ./conf/spark-env.sh.template ./conf/spark-env.sh #拷贝配置文件
```

- 编辑该配置文件，在文件最后面加上如下一行内容：

```
export SPARK_DIST_CLASSPATH=$(/usr/local/hadoop/bin/hadoop classpath)
```

- 保存配置文件后，就可以启动、运行Spark了
- 若需要使用HDFS中的文件，则在使用Spark前需要启动Hadoop





## 3.1.4 Spark和Hadoop的交互

Spark部署模式包括：

- Local模式：单机模式
  - Standalone模式：使用Spark自带的简单集群管理器
  - YARN模式：使用YARN作为集群管理器
  - Mesos模式：使用Mesos作为集群管理器
- 经过上面的步骤以后，就在单台机器上按照“Hadoop（伪分布式）+Spark（Local模式）”这种方式完成了Hadoop和Spark组合环境的搭建。
- Hadoop和Spark可以相互协作，由Hadoop的HDFS、HBase等组件负责数据的存储和管理，由Spark负责数据的计算。



## 3.2 在pyspark中运行代码

- pyspark提供了简单的方式来学习Spark API
- pyspark可以以实时、交互的方式来分析数据
- pyspark提供了Python交互式执行环境



## 3.2 在pyspark中运行代码

pyspark命令及其常用的参数如下：

```
pyspark --master <master-url>
```

Spark的运行模式取决于传递给SparkContext的Master URL的值。Master URL可以是以下任一种形式：

- \* local 使用一个Worker线程本地化运行SPARK(完全不并行)
- \* local[\*] 使用逻辑CPU个数数量的线程来本地化运行Spark
- \* local[K] 使用K个Worker线程本地化运行Spark（理想情况下，K应该根据运行机器的CPU核数设定）
- \* spark://HOST:PORT 连接到指定的Spark standalone master。默认端口是7077
- \* yarn-client 以客户端模式连接YARN集群。集群的位置可以在HADOOP\_CONF\_DIR 环境变量中找到
- \* yarn-cluster 以集群模式连接YARN集群。集群的位置可以在HADOOP\_CONF\_DIR 环境变量中找到
- \* mesos://HOST:PORT 连接到指定的Mesos集群。默认接口是5050



## 3.2 在pyspark中运行代码

在Spark中采用本地模式启动pyspark的命令主要包含以下参数：

**--master:** 这个参数表示当前的pyspark要连接到哪个master，如果是local[\*]，就是使用本地模式启动pyspark，其中，中括号内的星号表示需要使用几个CPU核心(core)，也就是启动几个线程模拟Spark集群

**--jars:** 这个参数用于把相关的JAR包添加到CLASSPATH中；如果有多个jar包，可以使用逗号分隔符连接它们



## 3.2 在pyspark中运行代码

比如，要采用本地模式，在4个CPU核心上运行pyspark:

```
$ cd /usr/local/spark  
$ ./bin/pyspark --master local[4]
```

或者，可以在CLASSPATH中添加code.jar，命令如下:

```
$ cd /usr/local/spark  
$ ./bin/pyspark --master local[4] --jars code.jar
```

可以执行“pyspark --help”命令，获取完整的选项列表，具体如下:

```
$ cd /usr/local/spark  
$ ./bin/pyspark --help
```



## 3.2 在pyspark中运行代码

执行如下命令启动pyspark（默认是local模式）：

```
$ cd /usr/local/spark  
$ ./bin/pyspark
```

启动pyspark成功后在输出信息的末尾可以看到“>>>”的命令提示符

```
Welcome to  
  
Spark version 2.4.0  
  
Using Python version 3.4.3 (default, Nov 12 2018 22:25:49)  
SparkSession available as 'spark'.  
>>> |
```





## 3.2 在pyspark中运行代码

可以在里面输入scala代码进行调试:

```
>>> 8*2+5  
21
```

可以使用命令“exit()”退出pyspark:

```
>>> exit()
```



## 3.3 开发Spark独立应用程序

### 3.3.1 编写程序

### 3.3.2 通过spark-submit运行程序



## 3.3.1 安装编译打包工具

WordCount.py

```
1 from pyspark import SparkConf, SparkContext
2 conf = SparkConf().setMaster("local").setAppName("My App")
3 sc = SparkContext(conf = conf)
4 logFile = "file:///usr/local/spark/README.md"
5 logData = sc.textFile(logFile, 2).cache()
6 numAs = logData.filter(lambda line: 'a' in line).count()
7 numBs = logData.filter(lambda line: 'b' in line).count()
8 print('Lines with a: %s, Lines with b: %s' % (numAs, numBs))
```

对于这段Python代码，可以直接使用如下命令执行：

```
$ cd /usr/local/spark/mycode/python
$ python3 WordCount.py
```

执行该命令以后，可以得到如下结果：

```
Lines with a: 62, Lines with b: 30
```



## 3.3.2 通过spark-submit运行程序

可以通过spark-submit提交应用程序，该命令的格式如下：

```
spark-submit  
--master <master-url>  
--deploy-mode <deploy-mode> #部署模式  
... #其他参数  
<application-file> #Python代码文件  
[application-arguments] #传递给主类的主方法的参数
```

可以执行“spark-submit --help”命令，获取完整的选项列表

```
$ cd /usr/local/spark  
$ ./bin/spark-submit --help
```



## 3.3.2 通过spark-submit运行程序

Master URL可以是以下任一种形式：

- \* local 使用一个Worker线程本地化运行SPARK(完全不并行)
- \* local[\*] 使用逻辑CPU个数数量的线程来本地化运行Spark
- \* local[K] 使用K个Worker线程本地化运行Spark（理想情况下，K应该根据运行机器的CPU核数设定）
- \* spark://HOST:PORT 连接到指定的Spark standalone master。默认端口是7077.
- \* yarn-client 以客户端模式连接YARN集群。集群的位置可以在HADOOP\_CONF\_DIR 环境变量中找到。
- \* yarn-cluster 以集群模式连接YARN集群。集群的位置可以在HADOOP\_CONF\_DIR 环境变量中找到。
- \* mesos://HOST:PORT 连接到指定的Mesos集群。默认接口是5050。



## 3.3.2 通过spark-submit运行程序

以通过 spark-submit 提交到 Spark 中运行，命令如下：

```
$ /usr/local/spark/bin/spark-submit /usr/local/spark/mycode/python/WordCount.py
```

可以在命令中间使用“\”符号，把一行完整命令“人为断开成多行”进行输入，效果如下：

```
$ /usr/local/spark/bin/spark-submit \  
> /usr/local/spark/mycode/python/WordCount.py
```

上面命令的执行结果如下：

```
Lines with a: 62, Lines with b: 30
```

为了避免其他多余信息对运行结果的干扰，可以修改log4j的日志信息显示级别：

```
log4j.rootCategory=INFO, console
```

修改为

```
log4j.rootCategory=ERROR, console
```





## 3.4 Spark集群环境搭建

3.4.1 集群概况

3.4.2 准备工作：搭建Hadoop集群环境

3.4.3 安装Spark

3.4.4 配置环境变量

3.4.5 Spark配置

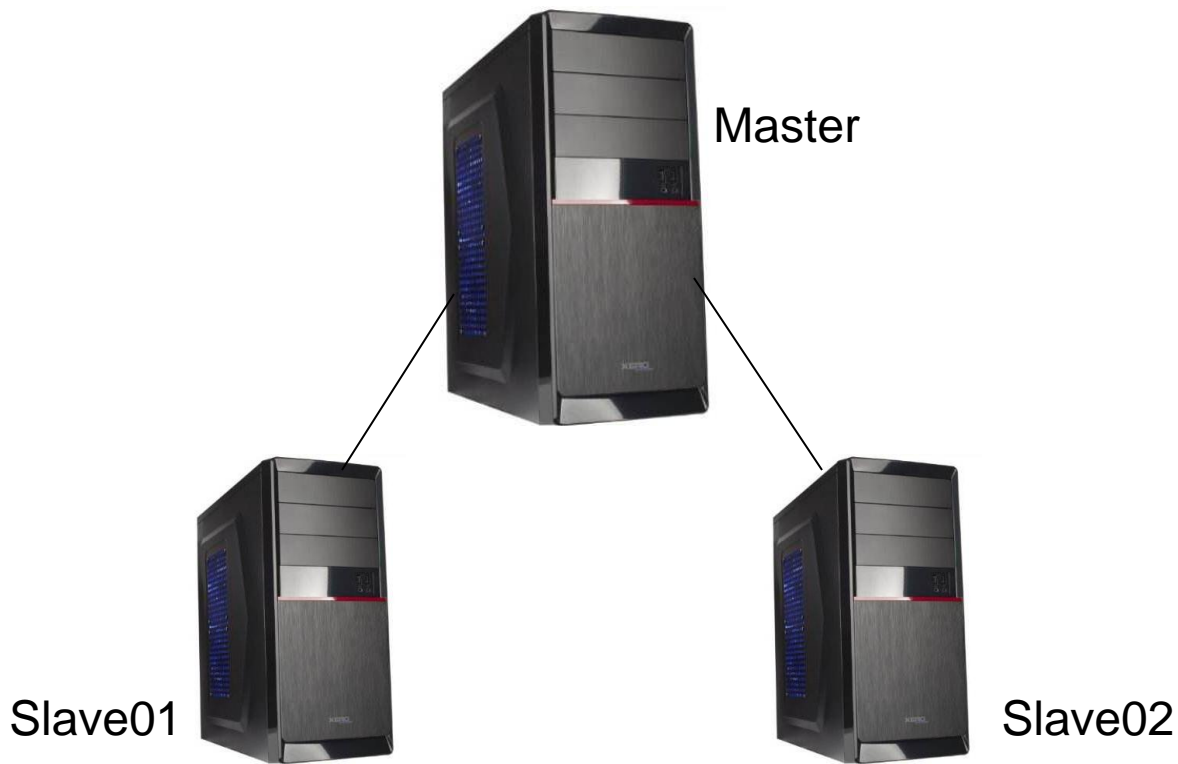
3.4.6 启动Spark集群

3.4.7 关闭Spark集群



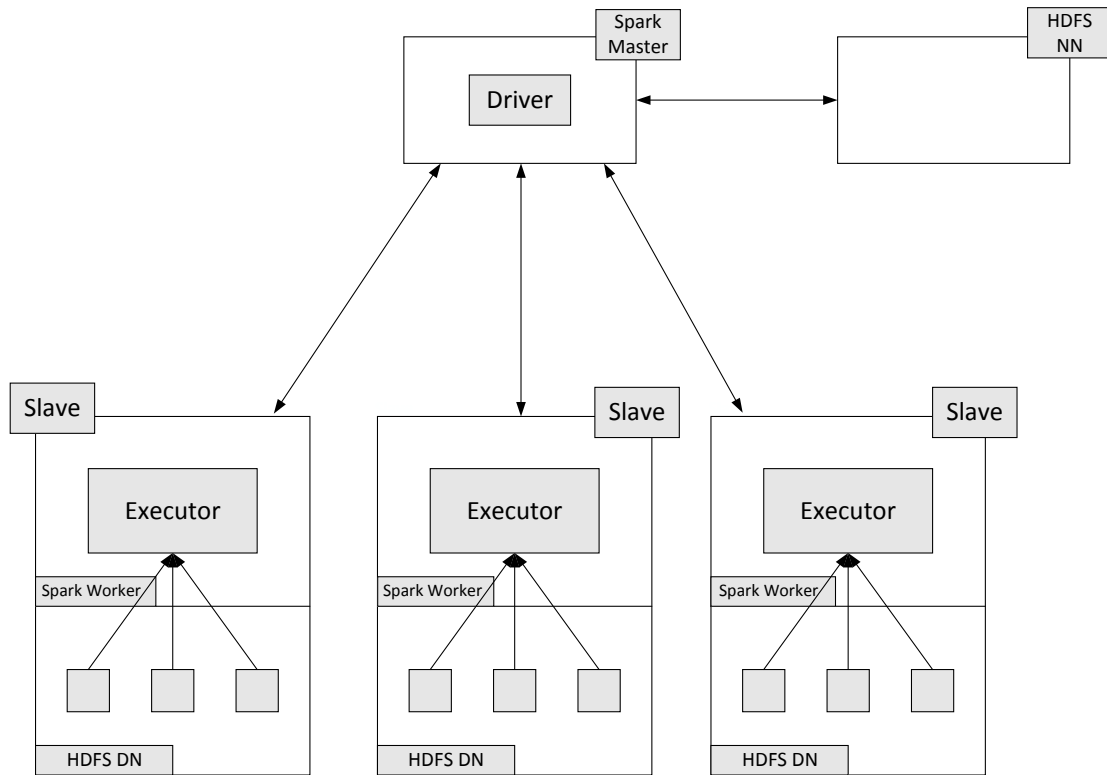
## 3.4.1 集群概况

- 采用3台机器（节点）作为实例来演示如何搭建Spark集群
- 其中1台机器（节点）作为Master节点
- 另外两台机器（节点）作为Slave节点（即作为Worker节点），主机名分别为Slave01和Slave02





## 3.4.2 准备工作：搭建Hadoop集群环境



Spark+HDFS运行架构

请参考厦门大学数据库实验室建设的“高校大数据课程公共服务平台”里面的技术博客：《**Hadoop 2.7**分布式集群环境搭建》  
文章地址：<http://dblab.xmu.edu.cn/blog/1177-2/>



## 3.4.3 安装Spark

在Master节点上，访问Spark官网下载Spark安装包

### Download Apache Spark™

1. Choose a Spark release:  ▾
2. Choose a package type:  ▾
3. Download Spark: [spark-2.4.0-bin-without-hadoop.tgz](#)
4. Verify this release using the [2.4.0 signatures and checksums](#) and [project release KEYS](#).

*Note: Starting version 2.0, Spark is built with Scala 2.11 by default. Scala 2.10 users should download the Spark source package and build [with Scala 2.10 support](#).*

```
sudo tar -zxf ~/下载/spark-2.4.0-bin-without-hadoop.tgz -C /usr/local/  
cd /usr/local  
sudo mv ./spark-2.4.0-bin-without-hadoop/ ./spark  
sudo chown -R hadoop ./spark
```



## 3.4.4 配置环境变量

在Master节点主机的终端中执行如下命令：

```
$ vim ~/.bashrc
```

在.bashrc添加如下配置：

```
export SPARK_HOME=/usr/local/spark  
export PATH=$PATH:$SPARK_HOME/bin:$SPARK_HOME/sbin
```

运行source命令使得配置立即生效：

```
$ source ~/.bashrc
```



## 3.4.5 Spark配置

### (1) 配置slaves文件

将 slaves.template 拷贝到 slaves

```
$ cd /usr/local/spark/  
$ cp ./conf/slaves.template ./conf/slaves
```

slaves文件设置Worker节点。编辑slaves内容,把默认内容localhost替换成如下内容:

```
Slave01  
Slave02
```





## 3.4.5 Spark配置

### (2) 配置spark-env.sh文件

将 spark-env.sh.template 拷贝到 spark-env.sh

```
$ cp ./conf/spark-env.sh.template ./conf/spark-env.sh
```

编辑spark-env.sh,添加如下内容:

```
export SPARK_DIST_CLASSPATH=$(/usr/local/hadoop/bin/hadoop classpath)
export HADOOP_CONF_DIR=/usr/local/hadoop/etc/hadoop
export SPARK_MASTER_IP=192.168.1.104
```



## 3.4.5 Spark配置

配置好后，将Master主机上的/usr/local/spark文件夹复制到各个节点上  
在Master主机上执行如下命令：

```
cd /usr/local/  
tar -zcf ~/spark.master.tar.gz ./spark  
cd ~  
scp ./spark.master.tar.gz slave01:/home/hadoop  
scp ./spark.master.tar.gz slave02:/home/hadoop
```

在slave01,slave02节点上分别执行下面同样的操作：

```
sudo rm -rf /usr/local/spark/  
sudo tar -zxf ~/spark.master.tar.gz -C /usr/local  
sudo chown -R hadoop /usr/local/spark
```



## 3.4.6 启动Spark集群

(1) 首先启动Hadoop集群。在Master节点主机上运行如下命令：

```
$ cd /usr/local/hadoop/  
$ sbin/start-all.sh
```

(2) 启动Master节点

在Master节点主机上运行如下命令：

```
$ cd /usr/local/spark/  
$ sbin/start-master.sh
```

(3) 启动所有Slave节点

在Master节点主机上运行如下命令：

```
$ sbin/start-slaves.sh
```



## 3.4.6 启动Spark集群

(4) 在浏览器上查看Spark独立集群管理器的集群信息

在Master主机上打开浏览器，访问<http://master:8080>,如下图:

### 2.4.0 Spark Master at spark://master:7077

URL: spark://master:7077  
REST URL: spark://master:6066 (cluster mode)  
Alive Workers: 2  
Cores in use: 2 Total, 0 Used  
Memory in use: 3.8 GB Total, 0.0 B Used  
Applications: 0 Running, 0 Completed  
Drivers: 0 Running, 0 Completed  
Status: ALIVE

#### Workers

Worker Id	Address	State
worker-20161205032642-192.168.1.108-35410	192.168.1.108:35410	ALIVE
worker-20161205032643-192.168.1.107-45533	192.168.1.107:45533	ALIVE



## 3.4.7 关闭Spark集群

在Master节点上执行下面命令

(1) 关闭Master节点

```
$ sbin/stop-master.sh
```

(2) 关闭Worker节点

```
$ sbin/stop-slaves.sh
```

(3) 关闭Hadoop集群

```
$ cd /usr/local/hadoop/  
$ sbin/stop-all.sh
```



## 3.5 在集群上运行Spark应用程序

3.5.1 启动Spark集群

3.5.2 采用独立集群管理器

3.5.3 采用Hadoop YARN管理器



## 3.5.1 启动Spark集群

请登录Linux系统，打开一个终端  
启动Hadoop集群

```
$ cd /usr/local/hadoop/  
$ sbin/start-all.sh
```

启动Spark的Master节点和所有slaves节点

```
$ cd /usr/local/spark/  
$ sbin/start-master.sh  
$ sbin/start-slaves.sh
```



## 3.5.2 采用独立集群管理器

### (1) 在集群中运行应用程序

- 向独立集群管理器提交应用，需要把spark://master:7077作为主节点参数递给spark-submit
- 可以运行Spark安装好以后自带的样例程序SparkPi，它的功能是计算得到pi的值（3.1415926）

```
$ cd /usr/local/spark/  
$ bin/spark-submit \  
> --master spark://master:7077 \  
> /usr/local/spark/examples/src/main/python/pi.py 2>&1 | grep "Pi is roughly"
```





## 3.5.2 采用独立集群管理器

### (2) 在集群中运行pyspark

也可以用pyspark连接到独立集群管理器上

```
$ cd /usr/local/spark/  
$ bin/pyspark --master spark://master:7077
```

```
>>> textFile = sc.textFile("hdfs://master:9000/README.md")  
>>> textFile.count()  
105  
>>> textFile.first()  
'# Apache Spark'
```



## 3.5.2 采用独立集群管理器

### (3) 查看集群信息

用户在独立集群管理Web界面查看应用的运行情况

<http://master:8080/>



Spark Master at spark://Master:7077

URL: spark://Master:7077

Alive Workers: 1

Cores in use: 2 Total, 0 Used

Memory in use: 2.8 GB Total, 0.0 B Used

Applications: 0 [Running](#), 4 [Completed](#)

Drivers: 0 Running, 0 Completed

Status: ALIVE

#### Workers (1)

Worker Id	Address	State	Cores	Memory
<a href="#">worker-20190103003706-192.168.1.110-42671</a>	192.168.1.110:42671	ALIVE	2 (0 Used)	2.8 GB (0.0 B Used)

#### Running Applications (0)

Application ID	Name	Cores	Memory per Executor	Submitted Time	User	State	Duration
----------------	------	-------	---------------------	----------------	------	-------	----------

#### Completed Applications (4)

Application ID	Name	Cores	Memory per Executor	Submitted Time	User	State	Duration
<a href="#">app-20190103011806-0003</a>	PythonPi	2	1024.0 MB	2019/01/03 01:18:06	hadoop	FINISHED	5 s
<a href="#">app-20190103010703-0002</a>	PythonPi	2	1024.0 MB	2019/01/03 01:07:03	hadoop	FINISHED	6 s
<a href="#">app-20190103010650-0001</a>	PythonPi	2	1024.0 MB	2019/01/03 01:06:50	hadoop	FINISHED	7 s
<a href="#">app-20190103003932-0000</a>	PythonPi	2	1024.0 MB	2019/01/03 00:39:32	hadoop	FINISHED	9 s



## 3.5.3 采用Hadoop YARN管理器

### (1) 在集群中运行应用程序

向Hadoop YARN集群管理器提交应用，需要把yarn-client或yarn-cluster作为主节点参数递给spark-submit

```
$ cd /usr/local/spark/  
$ bin/spark-submit \  
> --master yarn-client \  
> /usr/local/spark/examples/src/main/python/pi.py
```

运行后，根据在Shell中得到输出的结果地址查看，如下图：

```
2019-01-03 01:34:40 INFO Client:54 -  
client token: N/A  
diagnostics: N/A  
ApplicationMaster host: 192.168.1.110  
ApplicationMaster RPC port: -1  
queue: default  
start time: 1546508073776  
final status: UNDEFINED  
tracking URL: http://Master:8088/proxy/application_1546504350361_0006/  
user: hadoop
```



## 3.5.3 采用Hadoop YARN管理器

复制结果地址到浏览器，点击查看Logs，再点击stdout，即可查看结果，如下图：

Application Metrics					
<b>Total Resource Preempted:</b>	<memory:0, vCores:0>				
<b>Total Number of Non-AM Containers Preempted:</b>	0				
<b>Total Number of AM Containers Preempted:</b>	0				
<b>Resource Preempted from Current Attempt:</b>	<memory:0, vCores:0>				
<b>Number of Non-AM Containers Preempted from Current Attempt:</b>	0				
<b>Aggregate Resource Allocation:</b>	107090 MB-seconds, 50 vcore-seconds				
<b>Aggregate Preempted Resource Allocation:</b>	0 MB-seconds, 0 vcore-seconds				

Show 20 entries		Search: <input type="text"/>			
Attempt ID	Started	Node	Logs	Nodes blacklisted by the app	Nodes blacklisted by the system
<a href="#">appattempt_1546504350361_0006_000001</a>	Thu Jan 3 01:34:33 -0800 2019	<a href="http://Slave:8042">http://Slave:8042</a>	<a href="#">Logs</a>	0	0



## 3.5.3 采用Hadoop YARN管理器

### (2) 在集群中运行pyspark

也可以用pyspark连接到采用YARN作为集群管理器的集群上

```
$ bin/pyspark --master yarn
```

假设HDFS的根目录下已经存在一个文件README.md，下面在pyspark环境中执行相关语句：

```
>>> textFile = sc.textFile("hdfs://master:9000/README.md")
>>> textFile.count()
105
>>> textFile.first()
'# Apache Spark'
```



# 3.5.3 采用Hadoop YARN管理器

## (3) 查看集群信息

用户在Hadoop YARN集群管理Web界面查看所有应用的运行情况

http://master:8088/cluster



### All Applications

Cluster Metrics		Apps Submitted	Apps Pending	Apps Running	Apps Completed	Containers Running	Memory Used	Memory Total	Memory Reserved	Vcores Used	Vcores Total	Vcores Free						
		10	0	10	0	0 B	8 GB	0 B	0	8	0							
Cluster Nodes Metrics		Active Nodes	Decommissioning Nodes	Decommissioned Nodes	Lost Nodes	Unhealthy Nodes	Rebooted Nodes	Shutdown Nodes										
		1	0	0	0	0	0	0										
Scheduler Metrics		Scheduler Type	Scheduling Resource Type	Minimum Allocation	Maximum Allocation	Maximum Cluster Application Priority												
		Capacity Scheduler	[memory-mb (unit=M), vcores]	<memory:2048, vCores:1>	<memory:8192, vCores:4>	0												
Show 20 entries												Search:						
ID	User	Name	Application Type	Queue	Application Priority	StartTime	FinishTime	State	FinalStatus	Running Containers	Allocated CPU Vcores	Allocated Memory MB	Reserved CPU Vcores	Reserved Memory MB	% of Queue	% of Cluster	Progress	Tracking UI
application_1546504350361_0010	hadoop	PySparkShell	SPARK	default	0	Thu Jan 3 02:35:56 -0800 2019	Thu Jan 3 02:38:40 -0800 2019	FINISHED	SUCCEEDED	N/A	N/A	N/A	N/A	N/A	0.0	0.0	<div style="width: 100%;"></div>	History
application_1546504350361_0009	hadoop	PythonPi	SPARK	default	0	Thu Jan 3 02:33:52 -0800 2019	Thu Jan 3 02:34:13 -0800 2019	FINISHED	SUCCEEDED	N/A	N/A	N/A	N/A	N/A	0.0	0.0	<div style="width: 100%;"></div>	History
application_1546504350361_0008	hadoop	pi.py	SPARK	default	0	Thu Jan 3 02:32:13 -0800 2019	Thu Jan 3 02:32:44 -0800 2019	FINISHED	SUCCEEDED	N/A	N/A	N/A	N/A	N/A	0.0	0.0	<div style="width: 100%;"></div>	History
application_1546504350361_0007	hadoop	pi.py	SPARK	default	0	Thu Jan 3 02:31:06 -0800 2019	Thu Jan 3 02:31:32 -0800 2019	FINISHED	SUCCEEDED	N/A	N/A	N/A	N/A	N/A	0.0	0.0	<div style="width: 100%;"></div>	History
application_1546504350361_0006	hadoop	PythonPi	SPARK	default	0	Thu Jan 3 01:34:33 -0800 2019	Thu Jan 3 01:34:54 -0800 2019	FINISHED	SUCCEEDED	N/A	N/A	N/A	N/A	N/A	0.0	0.0	<div style="width: 100%;"></div>	History
application_1546504350361_0005	hadoop	PythonPi	SPARK	default	0	Thu Jan 3 01:20:49 -0800 2019	Thu Jan 3 01:21:08 -0800 2019	FINISHED	SUCCEEDED	N/A	N/A	N/A	N/A	N/A	0.0	0.0	<div style="width: 100%;"></div>	History



# 附录A：主讲教师林子雨简介



## 主讲教师：林子雨

单位：厦门大学计算机科学系

E-mail: ziyulin@xmu.edu.cn

个人网页: <http://dblab.xmu.edu.cn/post/linziyu>

数据库实验室网站: <http://dblab.xmu.edu.cn>



扫一扫访问个人主页

林子雨，男，1978年出生，博士（毕业于北京大学），现为厦门大学计算机科学系助理教授（讲师），曾任厦门大学信息科学与技术学院院长助理、晋江市发展和改革委员会副局长。中国计算机学会数据库专业委员会委员，中国计算机学会信息系统专业委员会委员。国内高校首个“数字教师”提出者和建设者，厦门大学数据库实验室负责人，厦门大学云计算与大数据研究中心主要建设者和骨干成员，2013年度和2017年度厦门大学教学类奖教金获得者，荣获2017年福建省精品在线开放课程、2018年厦门大学高等教育成果特等奖、2018年福建省高等教育教学成果二等奖、2018年国家精品在线开放课程。主要研究方向为数据库、数据仓库、数据挖掘、大数据、云计算和物联网，并以第一作者身份在《软件学报》《计算机学报》和《计算机研究与发展》等国家重点期刊以及国际学术会议上发表多篇学术论文。作为项目负责人主持的科研项目包括1项国家自然科学基金青年基金项目(No.61303004)、1项福建省自然科学基金项目(No.2013J05099)和1项中央高校基本科研业务费项目(No.2011121049)，主持的教改课题包括1项2016年福建省教改课题和1项2016年教育部产学协作育人项目，同时，作为课题负责人完成了国家发改委城市信息化重大课题、国家物联网重大应用示范工程区域试点泉州市工作方案、2015泉州市互联网经济调研等课题。中国高校首个“数字教师”提出者和建设者，2009年至今，“数字教师”大平台累计向网络免费发布超过500万字高价值的研究和教学资料，累计网络访问量超过500万次。打造了中国高校大数据教学知名品牌，编著出版了中国高校第一本系统介绍大数据知识的专业教材《大数据技术原理与应用》，并成为京东、当当网等网店畅销书籍；建设了国内高校首个大数据课程公共服务平台，为教师教学和学生学习大数据课程提供全方位、一站式服务，年访问量超过200万次。





# 附录B：大数据学习路线图



大数据学习路线图访问地址：<http://dmlab.xmu.edu.cn/post/10164/>





# 附录C：《大数据技术原理与应用》教材

《大数据技术原理与应用——概念、存储、处理、分析与应用（第2版）》，由厦门大学计算机科学系林子雨博士编著，是国内高校第一本系统介绍大数据知识的专业教材。人民邮电出版社 ISBN:978-7-115-44330-4 定价：49.80元



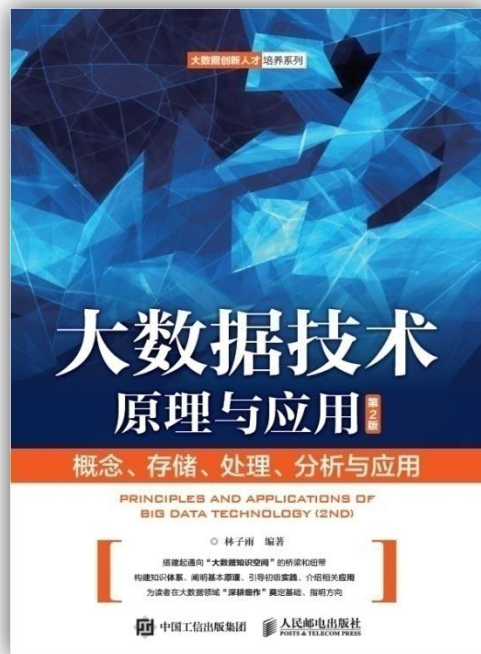
扫一扫访问教材官网

全书共有15章，系统地论述了大数据的基本概念、大数据处理架构Hadoop、分布式文件系统HDFS、分布式数据库HBase、NoSQL数据库、云数据库、分布式并行编程模型MapReduce、Spark、流计算、图计算、数据可视化以及大数据在互联网、生物医学和物流等各个领域的应用。在Hadoop、HDFS、HBase和MapReduce等重要章节，安排了入门级的实践操作，让读者更好地学习和掌握大数据关键技术。

本书可以作为高等院校计算机专业、信息管理等相关专业的大数据课程教材，也可供相关技术人员参考、学习、培训之用。

欢迎访问《大数据技术原理与应用——概念、存储、处理、分析与应用》教材官方网站：

<http://dbl原因.xmu.edu.cn/post/bigdata>





# 附录D：《大数据基础编程、实验和案例教程》

本书是与《大数据技术原理与应用（第2版）》教材配套的唯一指定实验指导书

大数据教材



1+1黄金组合  
厦门大学林子雨编著

配套实验指导书



- 步步引导，循序渐进，详尽的安装指南为顺利搭建大数据实验环境铺平道路
- 深入浅出，去粗取精，丰富的代码实例帮助快速掌握大数据基础编程方法
- 精心设计，巧妙融合，五套大数据实验题目促进理论与编程知识的消化和吸收
- 结合理论，联系实际，大数据课程综合实验案例精彩呈现大数据分析全流程

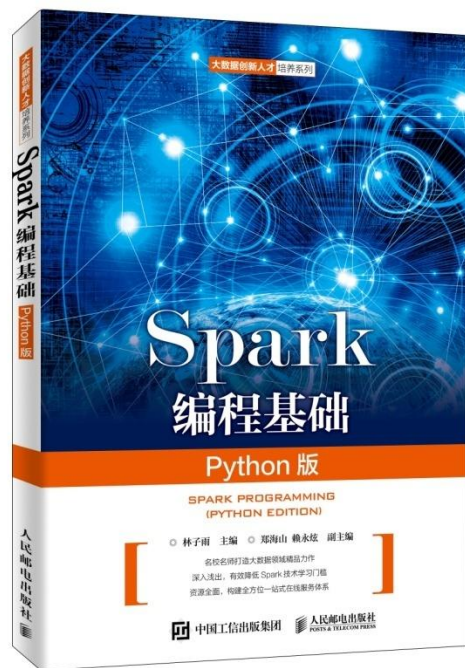
清华大学出版社 ISBN:978-7-302-47209-4 定价：59元



# 附录E：《Spark编程基础（Python版）》

林子雨，郑海山，赖永炫 编著 《Spark编程基础（Python版）》

教材官网：<http://dbllab.xmu.edu.cn/post/spark-python/>  
ISBN:978-7-115-52439-3 人民邮电出版社



本书以Python作为开发Spark应用程序的编程语言，系统介绍了Spark编程的基础知识。全书共8章，内容包括大数据技术概述、Spark的设计与运行原理、Spark环境搭建和使用方法、RDD编程、Spark SQL、Spark Streaming、Structured Streaming、Spark MLlib等。本书每个章节都安排了入门级的编程实践操作，以便读者更好地学习和掌握Spark编程方法。本书官网免费提供了全套的在线教学资源，包括讲义PPT、习题、源代码、软件、数据集、授课视频、上机实验指南等。





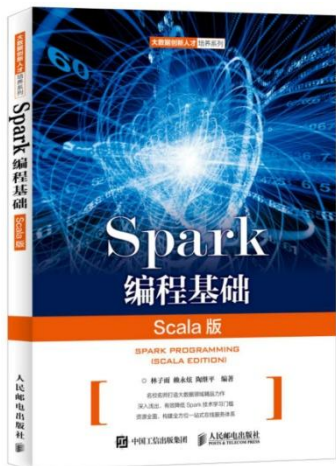
# 附录F：《Spark编程基础（Scala版）》

## 《Spark编程基础（Scala版）》

厦门大学 林子雨，赖永炫，陶继平 编著

披荆斩棘，在大数据丛林中开辟学习捷径  
填沟削坎，为快速学习Spark技术铺平道路  
深入浅出，有效降低Spark技术学习门槛  
资源全面，构建全方位一站式在线服务体系

人民邮电出版社出版发行，ISBN:978-7-115-48816-9  
教材官网：<http://dmlab.xmu.edu.cn/post/spark/>



本书以Scala作为开发Spark应用程序的编程语言，系统介绍了Spark编程的基础知识。全书共8章，内容包括大数据技术概述、Scala语言基础、Spark的设计与运行原理、Spark环境搭建和使用方法、RDD编程、Spark SQL、Spark Streaming、Spark MLlib等。本书每个章节都安排了入门级的编程实践操作，以便读者更好地学习和掌握Spark编程方法。本书官网免费提供了全套的在线教学资源，包括讲义PPT、习题、源代码、软件、数据集、授课视频、上机实验指南等。



# 附录G：高校大数据课程公共服务平台



## 高校大数据课程

公 共 服 务 平 台

<http://dbllab.xmu.edu.cn/post/bigdata-teaching-platform/>



扫一扫访问平台主页



扫一扫观看3分钟FLASH动画宣传片

The background of the slide features a blue gradient with several faint, light-blue silhouettes of people. At the top, there are two groups of people standing and holding hands. On the right side, a person is shown in profile, looking towards the center. On the left side, two people are shown in profile, one appearing to be speaking or gesturing towards the other. The overall theme is one of community and collaboration.

**Thank You!**

**Department of Computer Science, Xiamen University, 2020**