# 基于Spark集群在Windows 10下使用IntelliJ IDEA开发Spark应用程序

# WINDOWS10软件安装

● 前提条件：
  1. 基于**Linux**的**Spark**集群已经搭建完成（考虑到稳定性，推荐使用**CentOS**）
     https://github.com/QiqiDuan257/parallel-pso-spark/blob/master/How-to-Install-Spark-on-CentOS7-Chinese.md
     **Linux**用于开发、测试，同时也是用于生产的标准（唯一）平台
     **Mac OS X**只用于开发、测试
     **WINDOWS**只用于开发、测试
  2. 采用**WINDOWS10**的个人电脑需要与**Spark**集群处在同一个局域网内

● 需要在WINDOWS10下安装的软件：
  1. java version 1.8.0_131
     （通过CMD命令"*java --version*"验证JAVA是否安装成功）
  2. Scala version 2.11.11
     （通过CMD命令"*scala --version*"验证Scala是否安装成功）
  3. IntelliJ IDEA（Community版本）
     https://www.jetbrains.com/idea/download/#section=windows
  特别注意：
     WINDOWS10下的Java、Scala版本需要与Spark集群下的版本<span style="color:red">保持一致</span>

# WINDOWS10系统环境变量配置

● 以管理员身份修改*hosts*配置文件，保存**Spark**集群中所有节点的**IP**地址与主机名之间的
一一映射关系：

　　　　*hosts*配置文件位置：*C:\Windows\System32\drivers\etc\hosts*

*hosts*配置文件示例：
\# for the Spark commodity cluster
10.20.51.154　　dc001.syhlab dc001
10.20.42.194　　dc002.syhlab dc002
10.20.42.177　　dc003.syhlab dc003
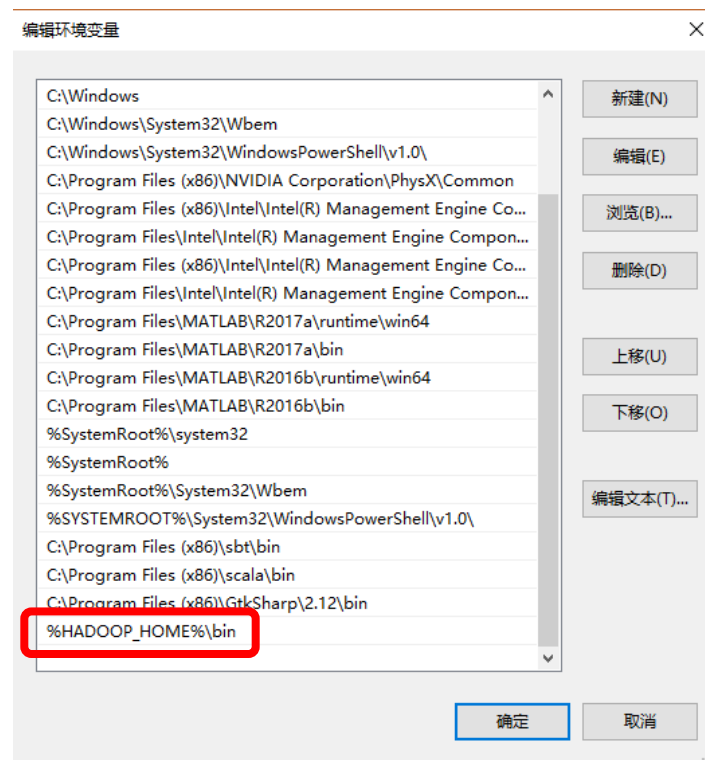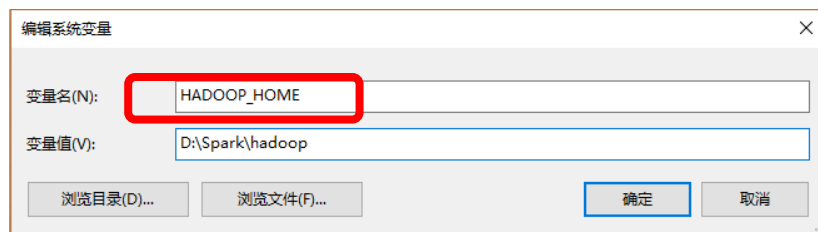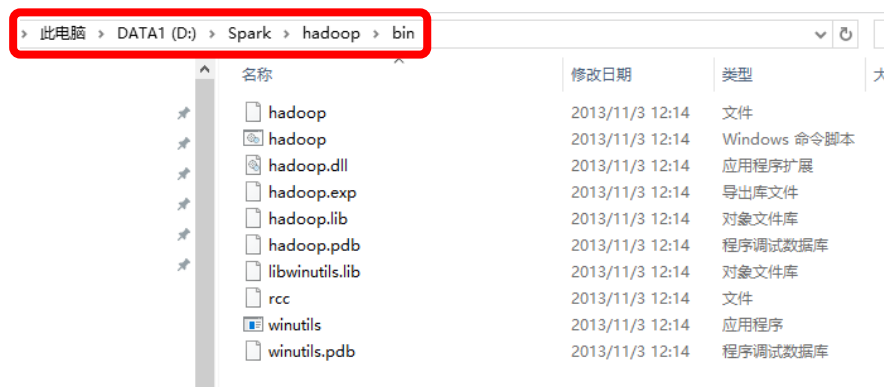10.20.42.175　　dc004.syhlab dc004
…　　　　　　…

# WINDOWS10系统环境变量配置

● 下载**Hadoop组件**到文件夹"**D:\Spark\hadoop**"中（此文件夹应只用于**Spark开发**）：
**https://github.com/srccodes/hadoop-common-2.2.0-bin**

● 配置环境变量*HADOOP_HOME*与*PATH*。否则，会抛出以下异常：
*java.io.IOException: Could not locate executable null\bin\winutils.exe in the Hadoop binaries.*

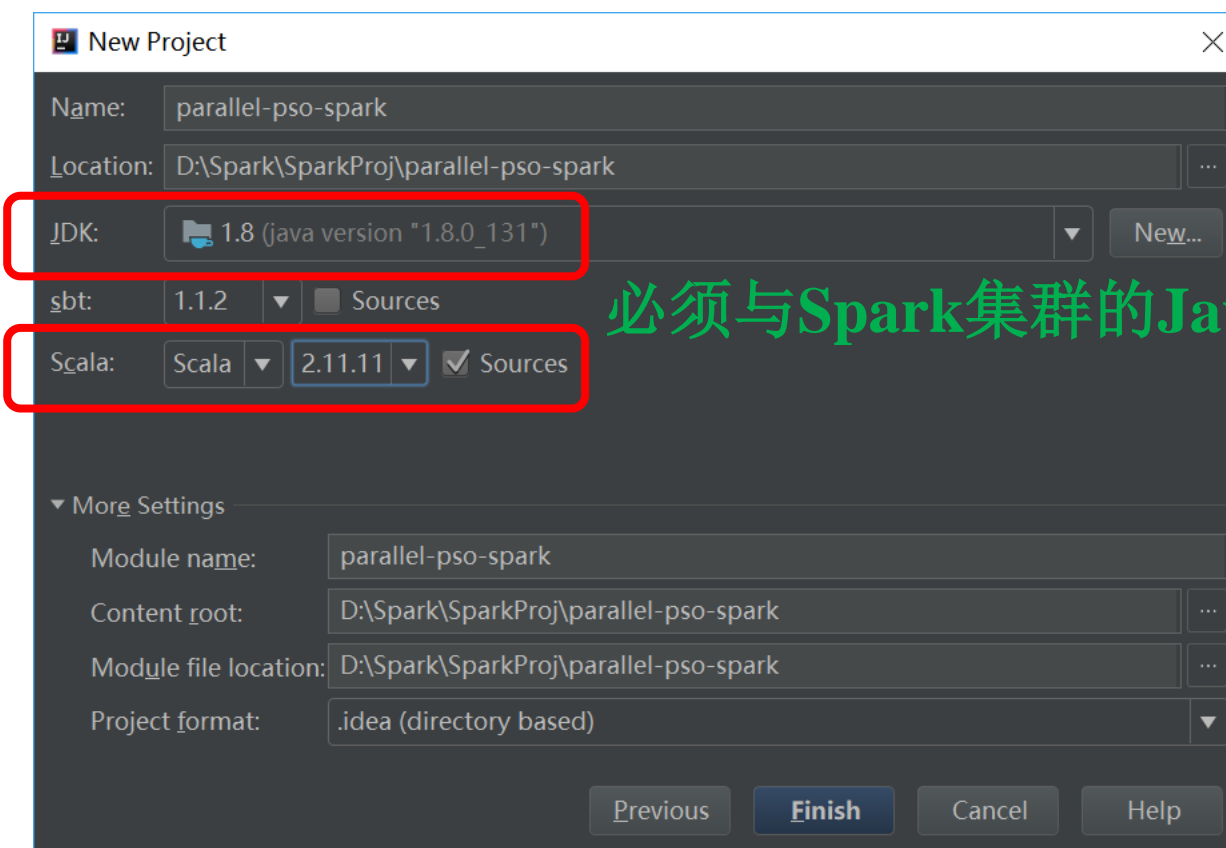# IntelliJ IDEA设置

● 创建基于**Spark**的应用程序：新建项目

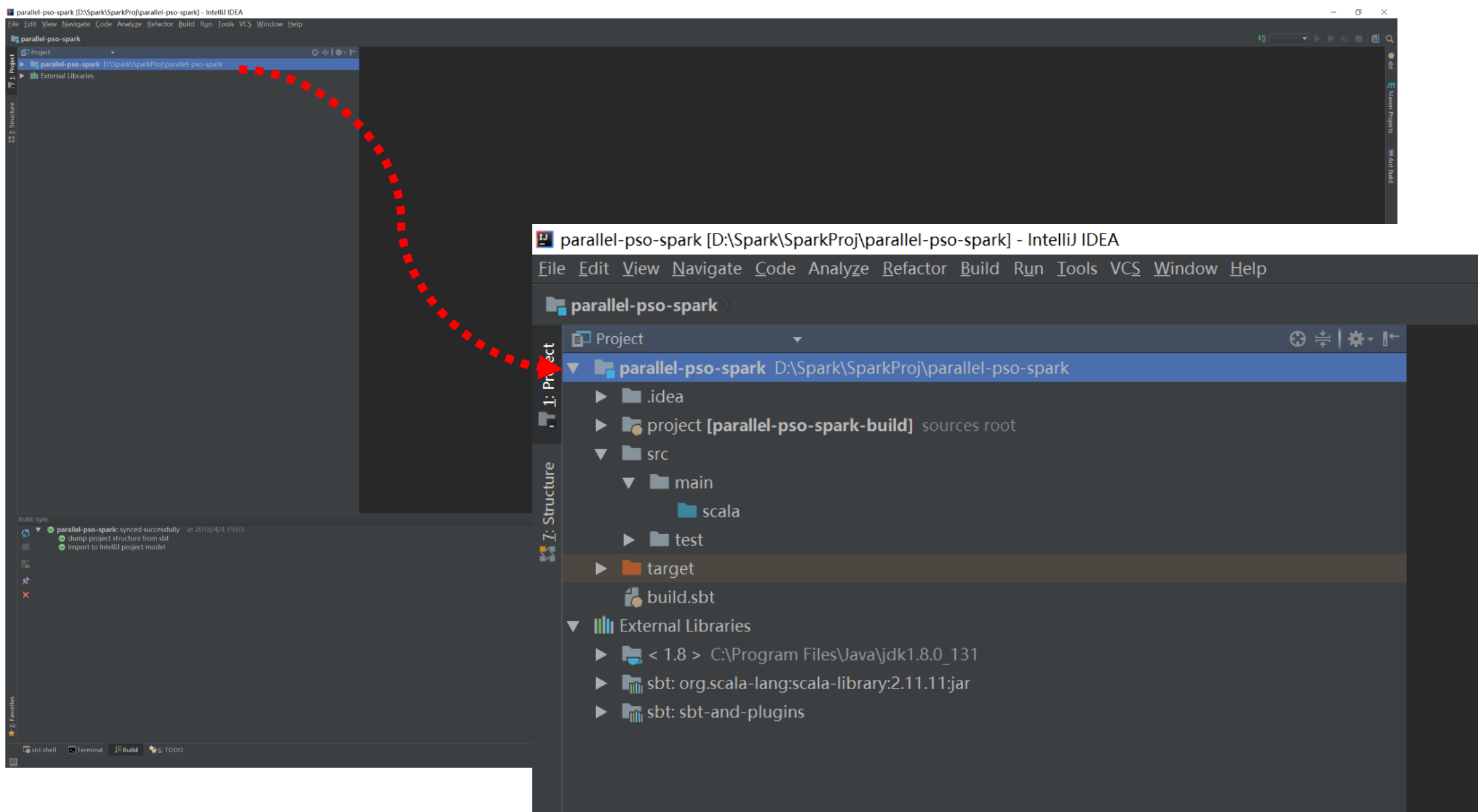# IntelliJ IDEA设置

● 创建基于**Spark**的应用程序：创建基于**Scala**语言的**sbt**项目

# IntelliJ IDEA设置

● 创建基于**Spark**的应用程序：配置**Java**、**sbt**、**Scala**版本

# IntelliJ IDEA设置

● 创建基于**Spark**的应用程序：项目初始化界面

# IntelliJ IDEA设置

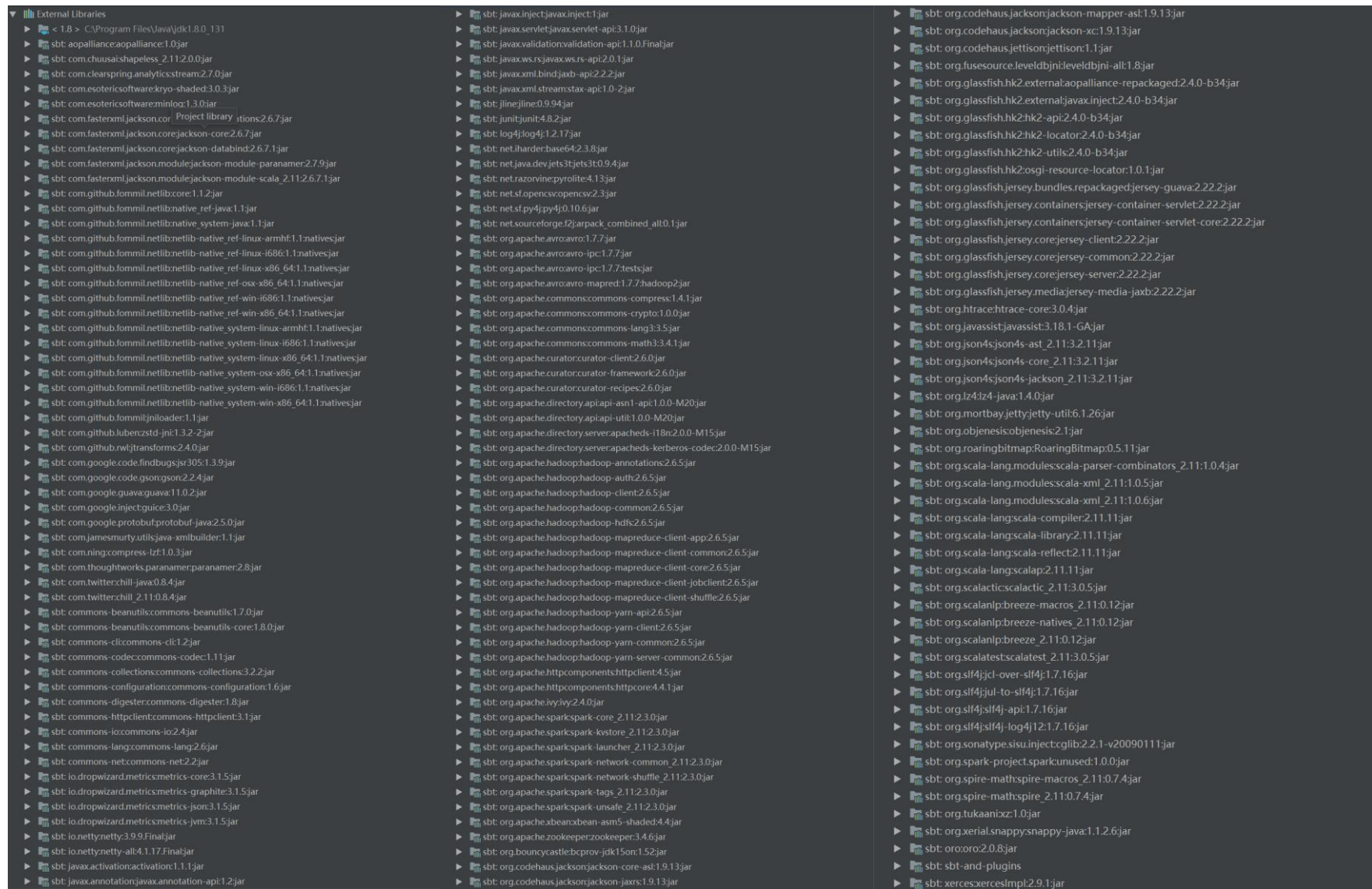● 从网站 **https://mvnrepository.com/** 中获取Spark项目的sbt依赖包

# IntelliJ IDEA设置

● 创建基于**Spark**的应用程序：更新*build.sbt*文件，自动载入所有的依赖包
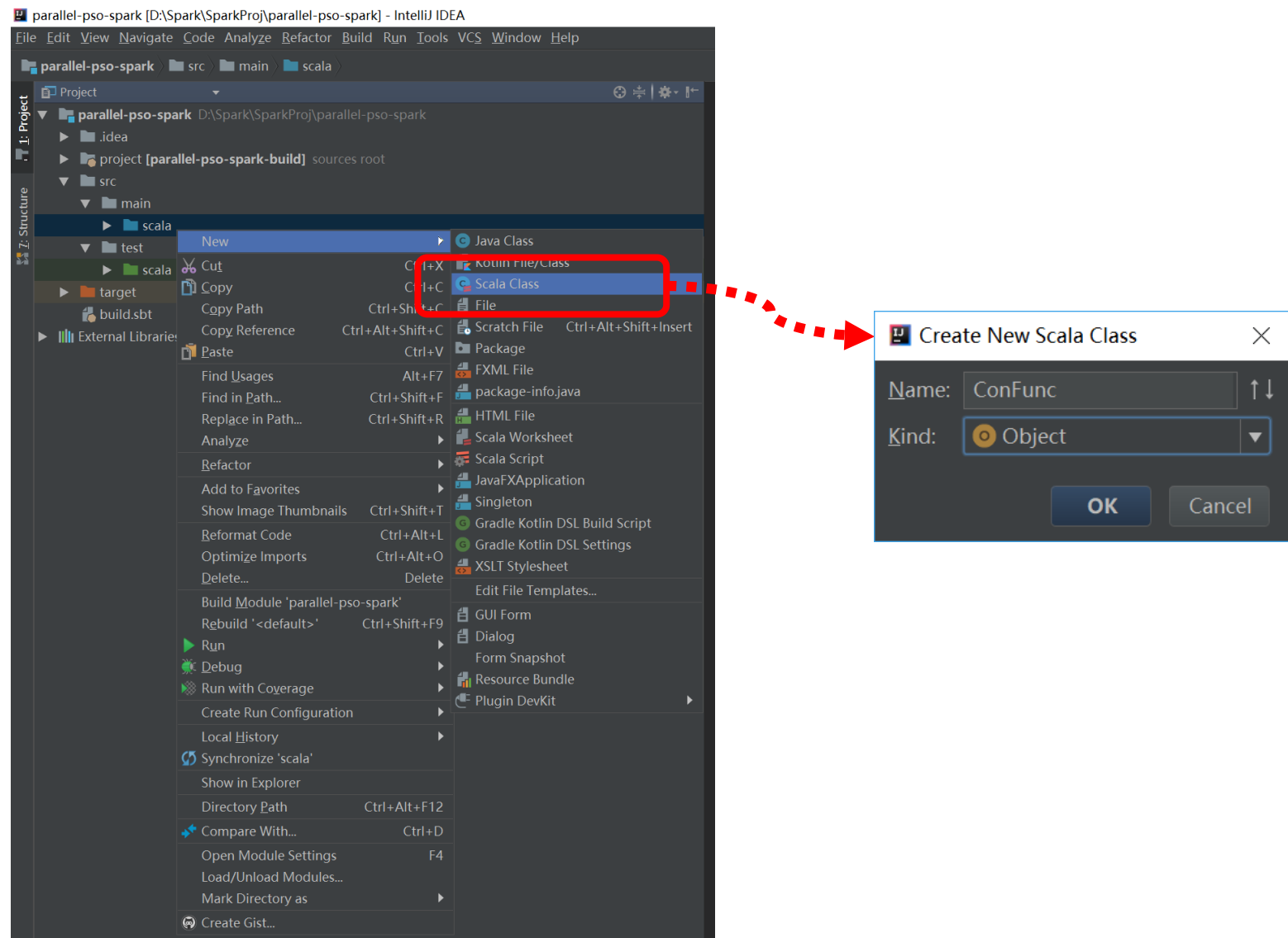　　可能等待较长时间（取决于网络速度）

# IntelliJ IDEA设置

● 创建基于**Spark**的应用程序：载入后的依赖包
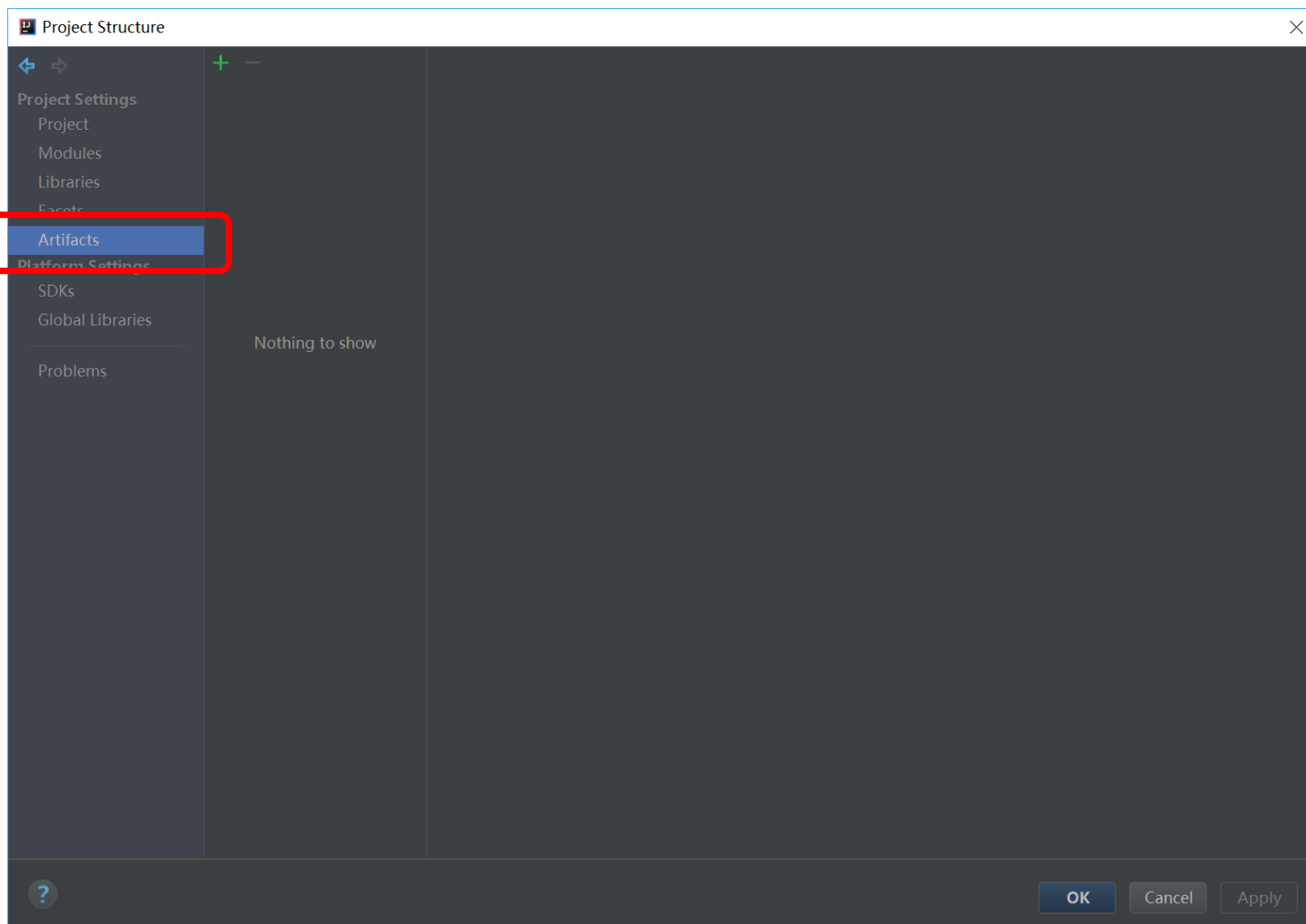
# IntelliJ IDEA设置

● 创建基于Spark的应用程序：新建object

# IntelliJ IDEA设置

● 创建基于Spark的应用程序：创建一个简单的示例，验证是否能够连接到Spark集群

```scala
package peas

import org.apache.spark.{SparkConf, SparkContext}

/**
  * Demo to Connect the Spark Commodity Cluster.
  */
object DemoConnectSparkCluster extends App {
  // configure and start the connection to the Spark commodity cluster
  val sparkConf = new SparkConf()
    .setAppName("DemoConnectSparkCluster")
    .setMaster("spark://dc001.syhlab:7077")
    .setJars(List("D:\\Spark\\SparkProj\\parallel-pso-spark\\out\\artifacts\\parallel_pso_spark_jar\\parallel-pso-spark.jar"))
  val sc = new SparkContext(sparkConf)

  // calculate the value of PI
  val numSamples = 1000000000
  val count = sc.parallelize(1 to numSamples, 100)
    .filter({_ =>
      val x = math.random
      val y = math.random
      x * x + y * y < 1
    }).count()
  println(f"PI =~= ${4.0 * count / numSamples}%9.7f")

  // close the connection to the Spark commodity cluster
  sc.stop()
}
```
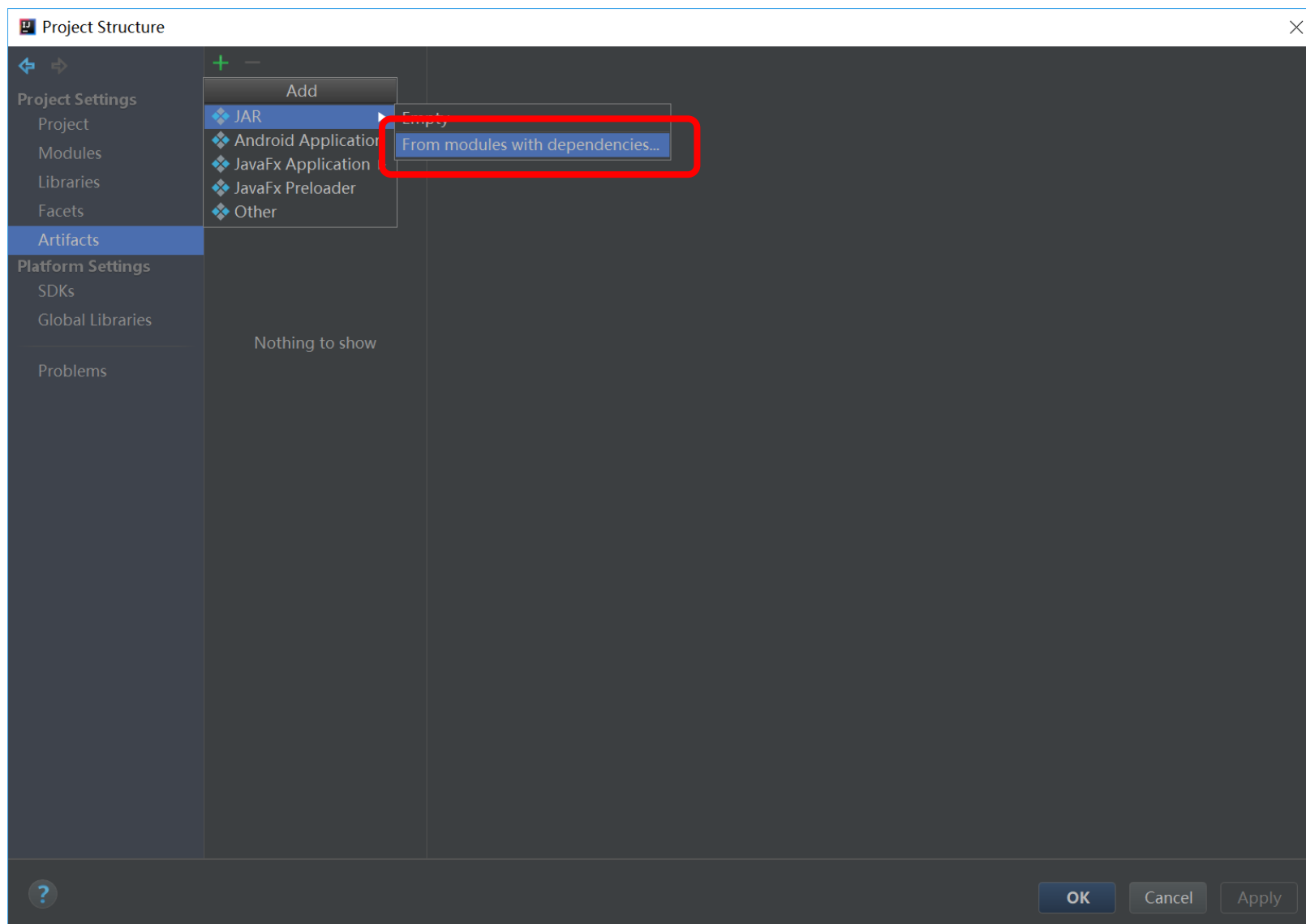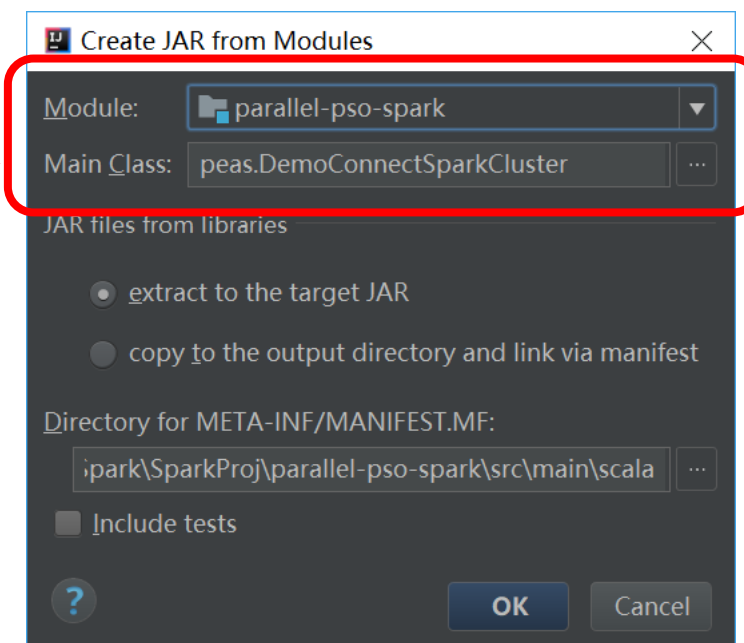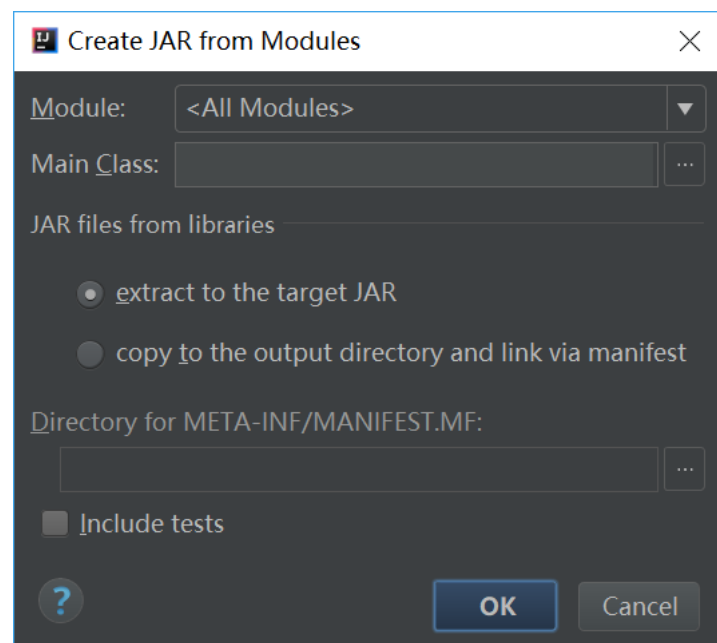
# IntelliJ IDEA设置

● 创建基于Spark的应用程序：配置项目输出

# IntelliJ IDEA设置

● 创建基于Spark的应用程序：配置项目输出

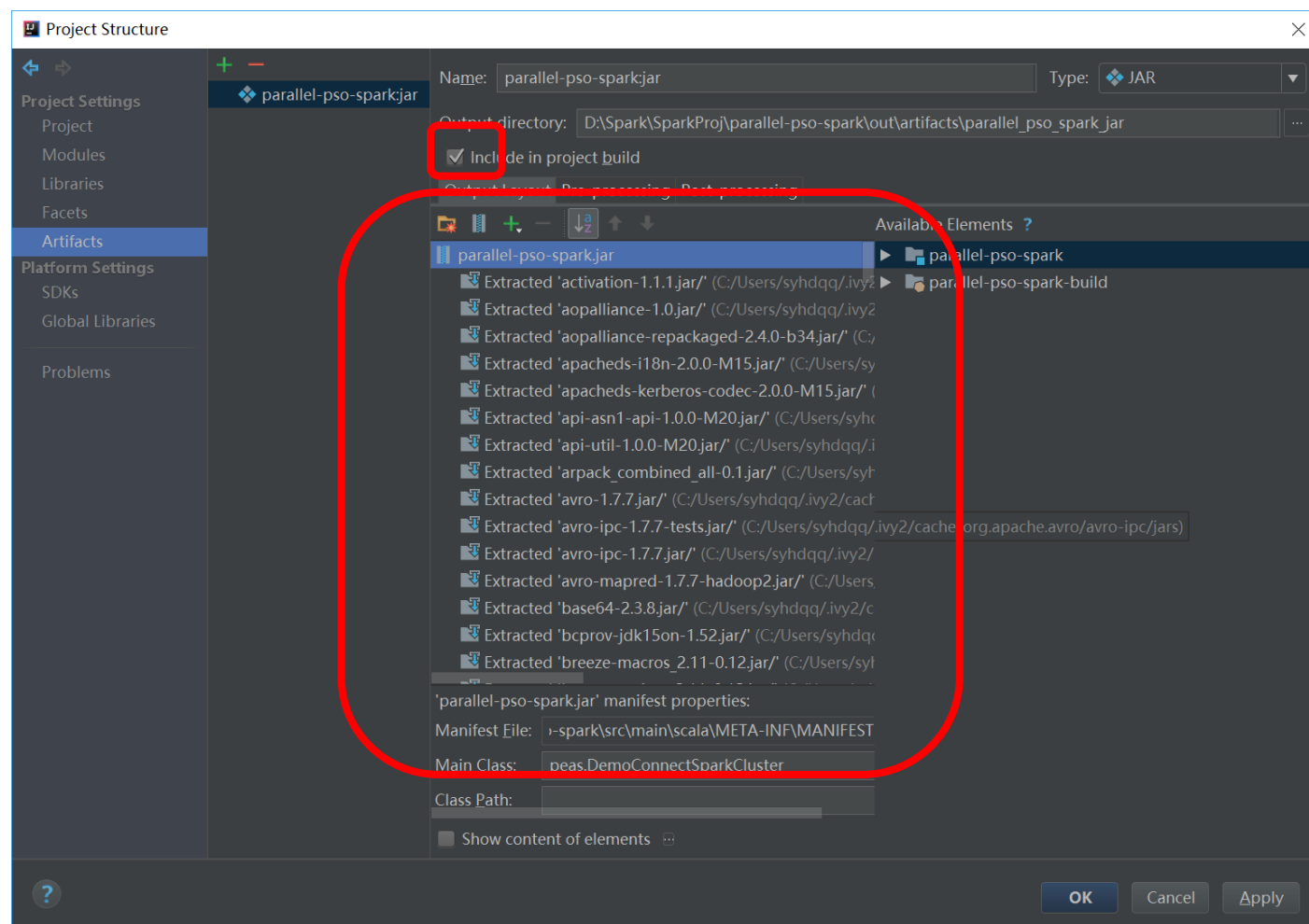# IntelliJ IDEA设置

● 创建基于**Spark**的应用程序：配置项目输出

# IntelliJ IDEA设置

● 创建基于**Spark**的应用程序：配置项目输出
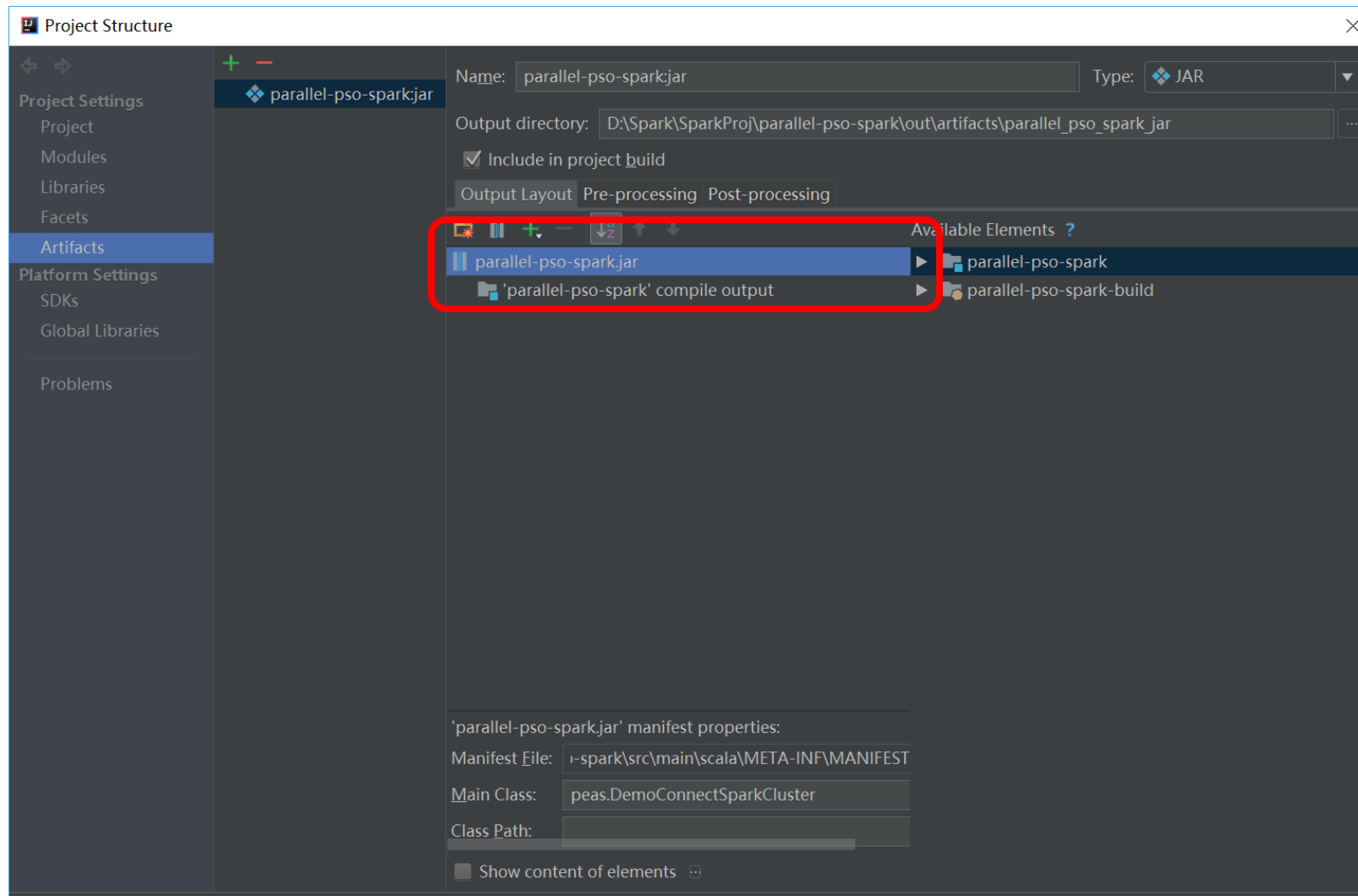
如果以下项未能配置正确，将出现以下错误提示：

*ERROR SparkContext: Failed to add \*.jar to Spark environment java.lang.ClassNotFoundException*
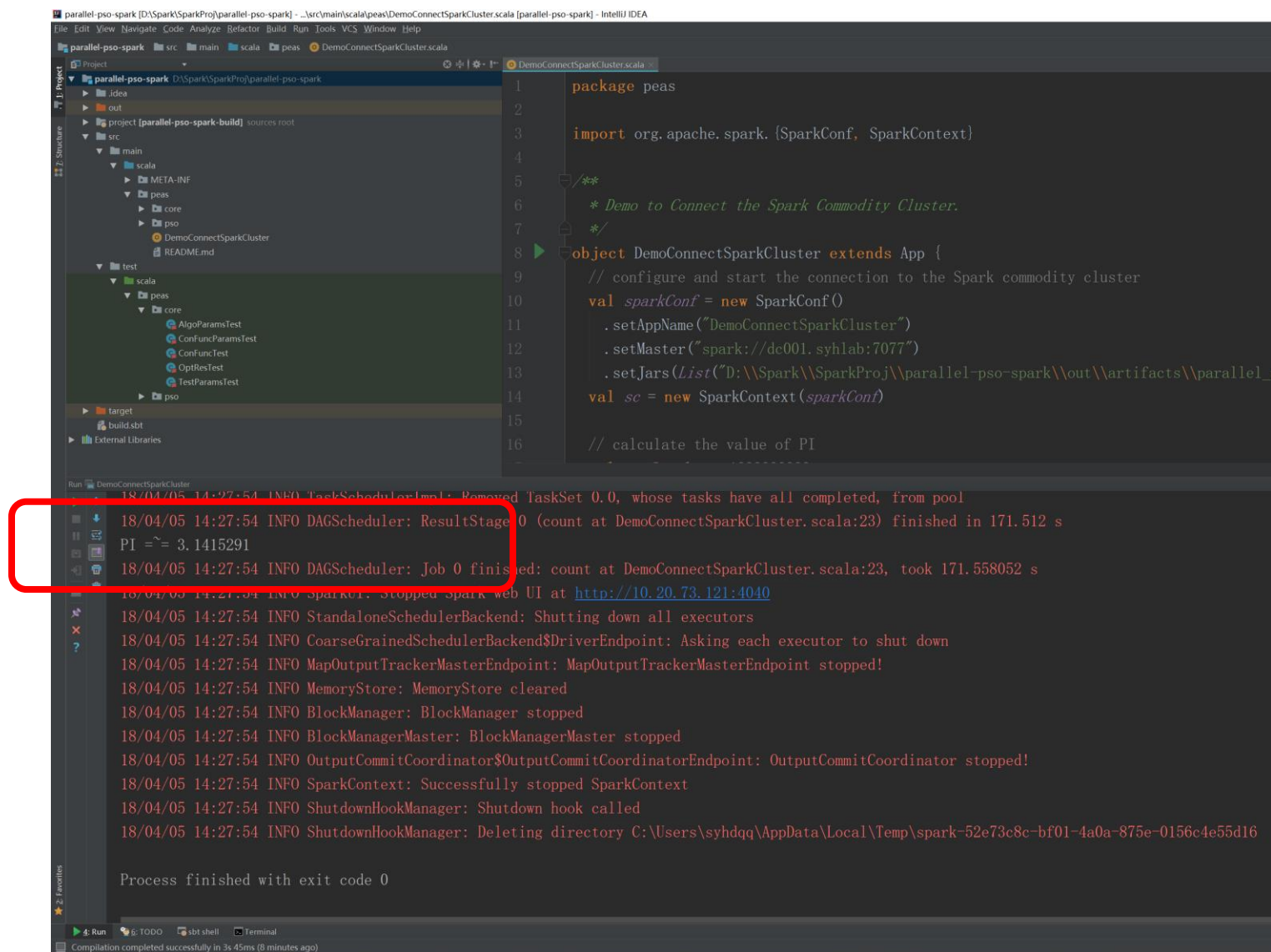
# IntelliJ IDEA设置

● 创建基于**Spark**的应用程序：配置项目输出

　　如果以下项未能配置正确，将出现以下错误提示：

*ERROR SparkContext: Failed to add \*.jar to Spark environment java.lang.ClassNotFoundException*

# IntelliJ IDEA设置

● 创建基于Spark的应用程序：程序输出（部分）

# IntelliJ IDEA设置

● 创建基于**Spark**的应用程序：程序监控（通过**WEB UI**）

# 常见问题汇总

● **TaskSchedulerImpl: Initial job has not accepted any resources; check your cluster UI to ensure that workers are registered and have sufficient resources.**

解决方案之一（造成此问题的原因可能有多种，这里只给出我个人遇到的情况）：

确保**Windows**开发端与**Spark**集群处于同一级的局域网内：

特别是在**Windows**开发端使用路由器（例如**TP-LINK**）进行网络连接时，当**Spark**集群处于上一级**IP**地址（例如*10.20.2.5*）而**Windows**开发端处于下一级**IP**地址时（例如**192.168.7.9**），需要将路由器转化为交换机模式。

● **WARN TransportChannelHandler: Exception in connection from 10.20.2.5:54321 java.io.IOException: Connection reset by peer**

解决方案之一（造成此问题的原因可能有多种，这里只给出我个人遇到的情况）：

关闭**Windows**防火墙　　*可选*

# 常见问题汇总

*查看、分析log日志*是对**Spark**应用程序进行排错的重要手段。

```
[dis@dc001 spark-2.3.0-bin-hadoop2.7]$ ls
bin   conf   data   examples   jars   kubernetes   LICENSE   licenses   logs   NOTICE   python   R   README.md   RELEASE   sbin   work   yarn
```