# 零基础学习Spark 1.x应用开发系列课程

## 如何使用IDEA开发程序
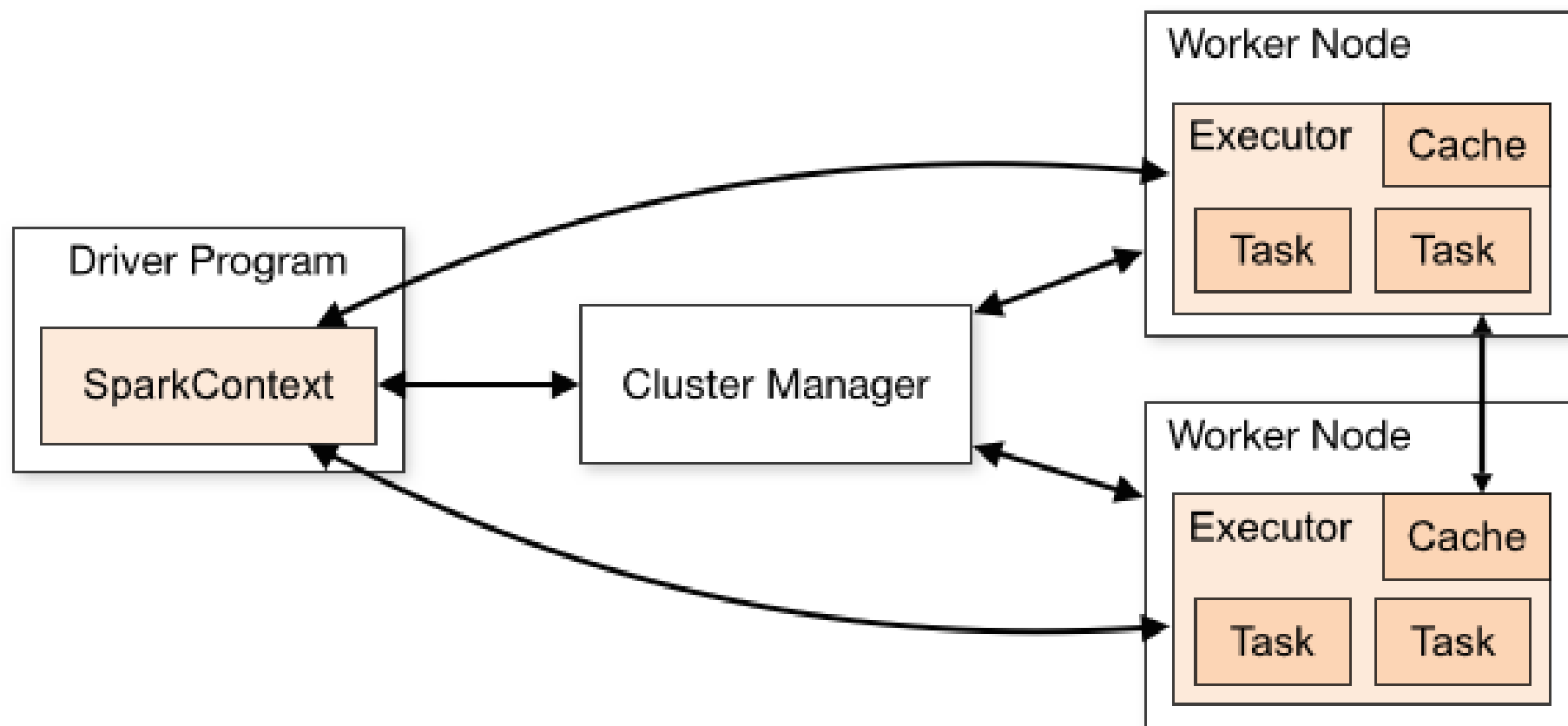
讲师-梦琪

【声明】本视频和讲义等均为云帆大数据网络课程的教学资料，所有资料只能在课程内使用，不允许在课程以外范围散播，违者将可能被追究法律和经济责任。

课程详情访问云帆大数据官方培训网站

**http://www.cloudyhadoop.com**

**1、构建Spark Application运行环境；**

在Driver Program中新建SparkContext（包含sparkcontext的程序称为Driver Program）；

**Spark Application运行的表现方式为：在集群上运行着一组独立的executor进程，这些进程由sparkcontext来协调；**

**2、SparkContext向资源管理器申请运行Executor资源，并启动StandaloneExecutorBackend，executor向sparkcontent申请task；**

集群通过SparkContext连接到不同的cluster manager(standalone、yarn、mesos)，cluster manager为运行应用的Executor分配资源；一旦连接建立之后，Spark每个Application就会获得各个节点上的Executor（进程）；每个Application都有自己独立的executor进程；Executor才是真正运行在WorkNode上的工作进程，它们为应用来计算或者存储数据；

**3、SparkContext获取到executor之后，Application的应用代码将会被发送到各个executor；**

**4、SparkContext构建RDD DAG图，将RDD DAG图分解成Stage DAG图，将Stage提交给TaskScheduler，最后由TaskScheduler将Task发送给Executor运行；**

**5、Task在Executor上运行，运行完毕后释放所有资源；**

| Term | Meaning |
|---|---|
| Application | User program built on Spark. Consists of a *driver program* and *executors* on the cluster. |
| Application jar | A jar containing the user's Spark application. In some cases users will want to create an "uber jar" containing their application along with its dependencies. The user's jar should never include Hadoop or Spark libraries, however, these will be added at runtime. |
| Driver program | The process running the main() function of the application and creating the SparkContext |
| Cluster manager | An external service for acquiring resources on the cluster (e.g. standalone manager, Mesos, YARN) |
| Deploy mode | Distinguishes where the driver process runs. In "cluster" mode, the framework launches the driver inside of the cluster. In "client" mode, the submitter launches the driver outside of the cluster. |

| Worker node | Any node that can run application code in the cluster |
| --- | --- |
| Executor | A process launched for an application on a worker node, that runs tasks and keeps data in memory or disk storage across them. Each application has its own executors. |
| Task | A unit of work that will be sent to one executor |
| Job | A parallel computation consisting of multiple tasks that gets spawned in response to a Spark action (e.g. save, collect); you'll see this term used in the driver's logs. |
| Stage | Each job gets divided into smaller sets of tasks called *stages* that depend on each other (similar to the map and reduce stages in MapReduce); you'll see this term used in the driver's logs. |

```
./bin/spark-submit \
  --class <main-class>
  --master <master-url> \
  --deploy-mode <deploy-mode> \
  --conf <key>=<value> \
  ... # other options
  <application-jar> \
  [application-arguments]
```

Some of the commonly used options are:

- `--class`: The entry point for your application (e.g. `org.apache.spark.examples.SparkPi`)
- `--master`: The master URL for the cluster (e.g. `spark://23.195.26.187:7077`)
- `--deploy-mode`: Whether to deploy your driver on the worker nodes (`cluster`) or locally as an external client (`client`) (default: `client`) †
- `--conf`: Arbitrary Spark configuration property in key=value format. For values that contain spaces wrap "key=value" in quotes (as shown).
- `application-jar`: Path to a bundled jar including your application and all dependencies. The URL must be globally visible inside of your cluster, for instance, an `hdfs://` path or a `file://` path that is present on all nodes.
- `application-arguments`: Arguments passed to the main method of your main class, if any

# Master URLs

| Master URL | Meaning |
| --- | --- |
| local | Run Spark locally with one worker thread (i.e. no parallelism at all). |
| local[K] | Run Spark locally with K worker threads (ideally, set this to the number of cores on your machine). |
| local[*] | Run Spark locally with as many worker threads as logical cores on your machine. |
| spark://HOST:PORT | Connect to the given Spark standalone cluster master. The port must be whichever one your master is configured to use, which is 7077 by default. |
| mesos://HOST:PORT | Connect to the given Mesos cluster. The port must be whichever one your is configured to use, which is 5050 by default. Or, for a Mesos cluster using ZooKeeper, use `mesos://zk://....` |
| yarn-client | Connect to a YARN cluster in client mode. The cluster location will be found based on the HADOOP_CONF_DIR variable. |
| yarn-cluster | Connect to a YARN cluster in cluster mode. The cluster location will be found based on HADOOP_CONF_DIR. |

# 零基础学习Spark 1.x应用

# 开发系列课程

## 安装IDEA及Scala插件

讲师-梦琪

# 选择IDEA UI Theme

# 安装Scala 插件

# 零基础学习Spark 1.x应用开发系列课程

## 导入Spark源码

讲师-梦琪

# 导入Spark源码

# 导入Spark源码
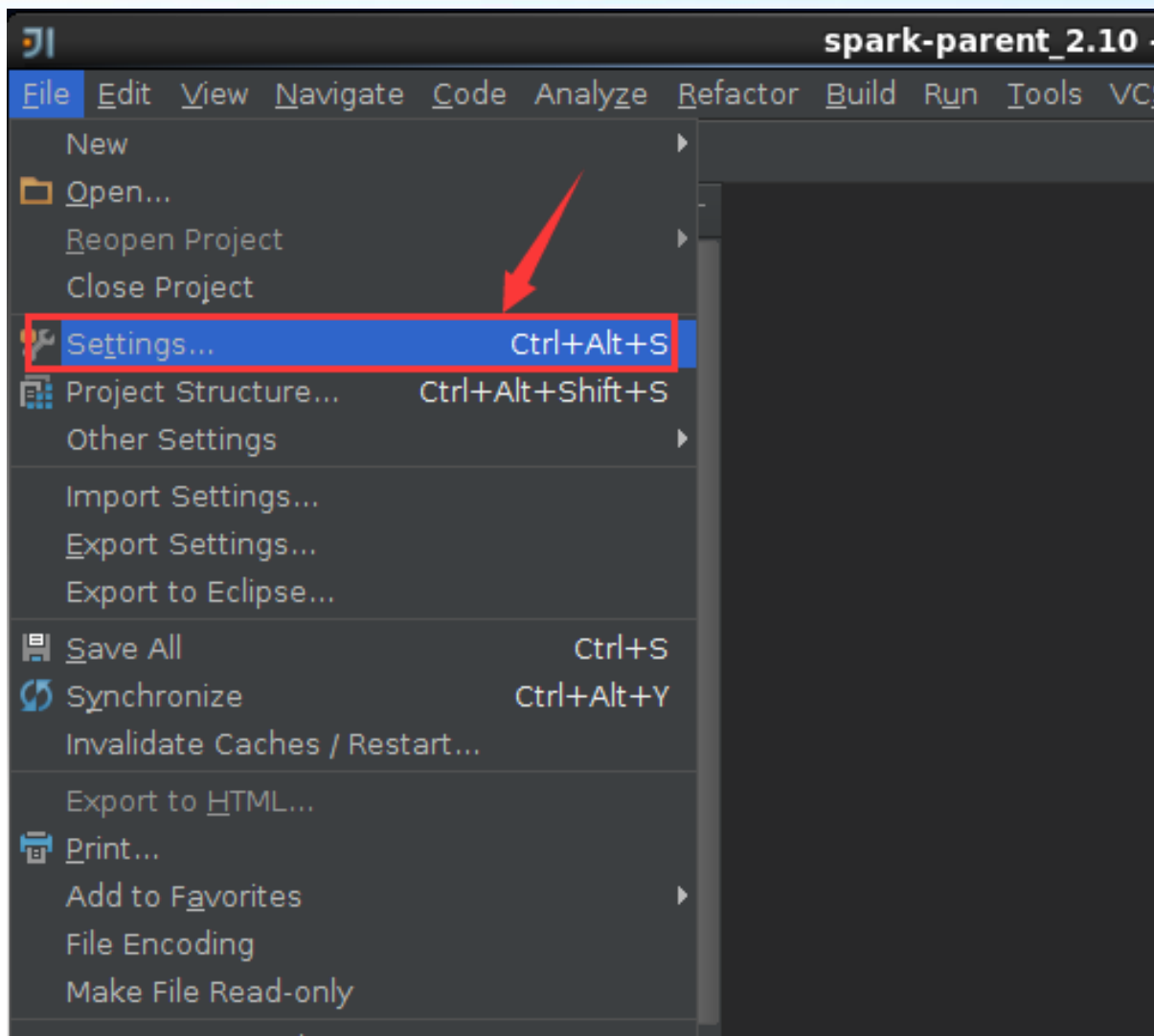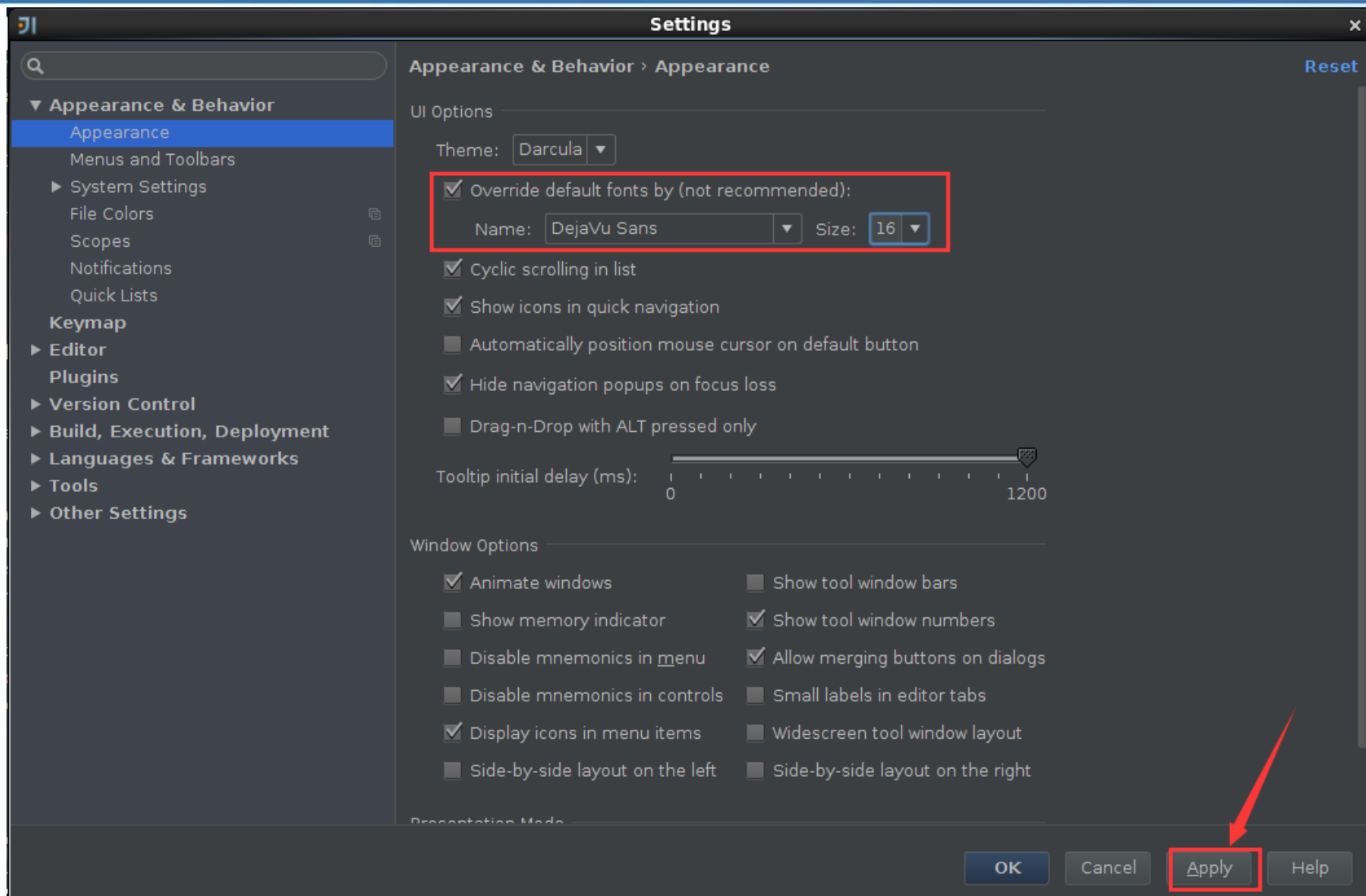


《零基础学习Spark 1.x应用开发系列课程》 讲师：梦琪

26

# 零基础学习Spark 1.x应用开发系列课程

## IDEA基本配置

讲师-梦琪

# IDEA基本配置

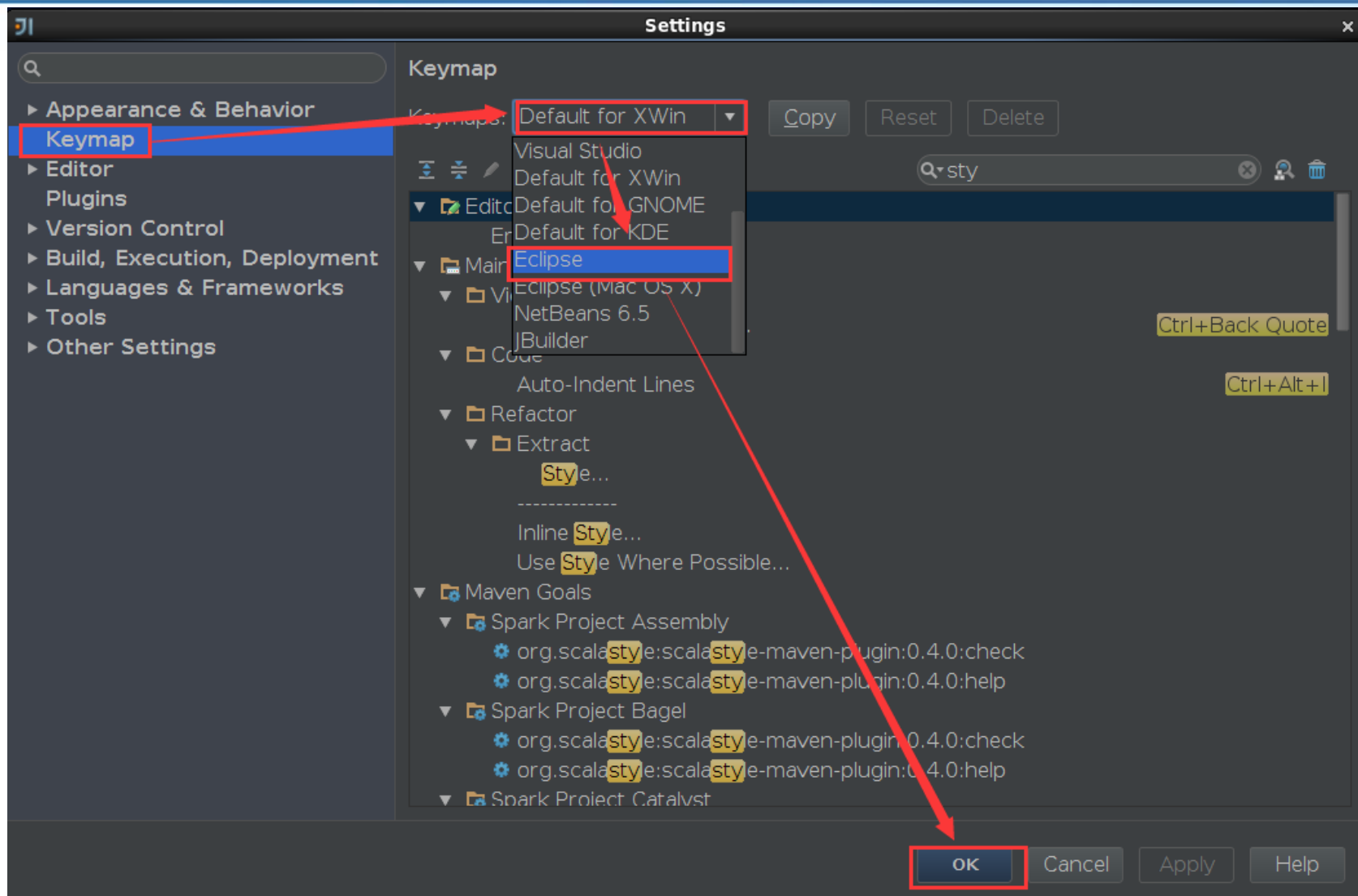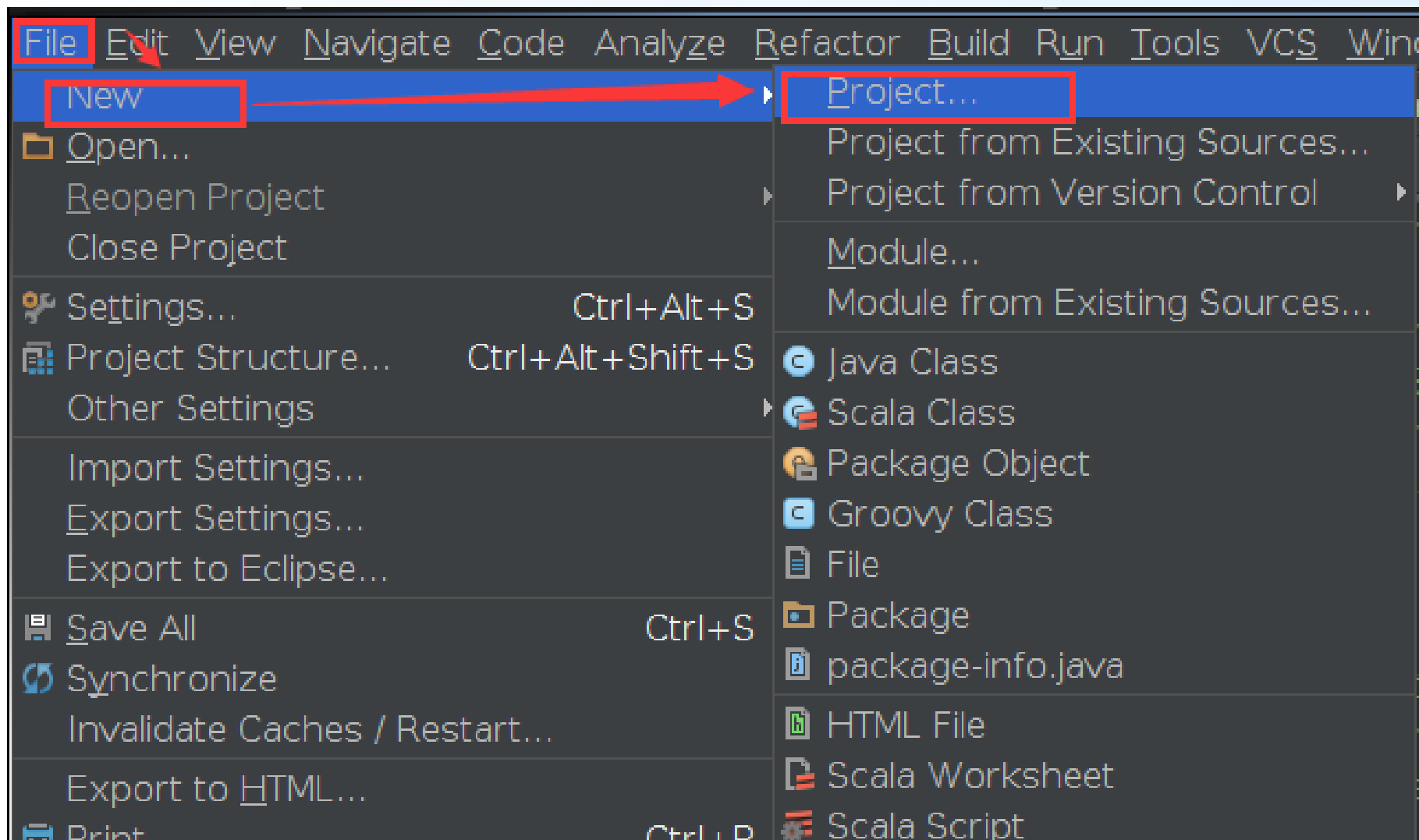# IDEA基本配置

# 零基础学习Spark 1.x应用

# 开发系列课程

## 创建Scala Project

讲师-梦琪

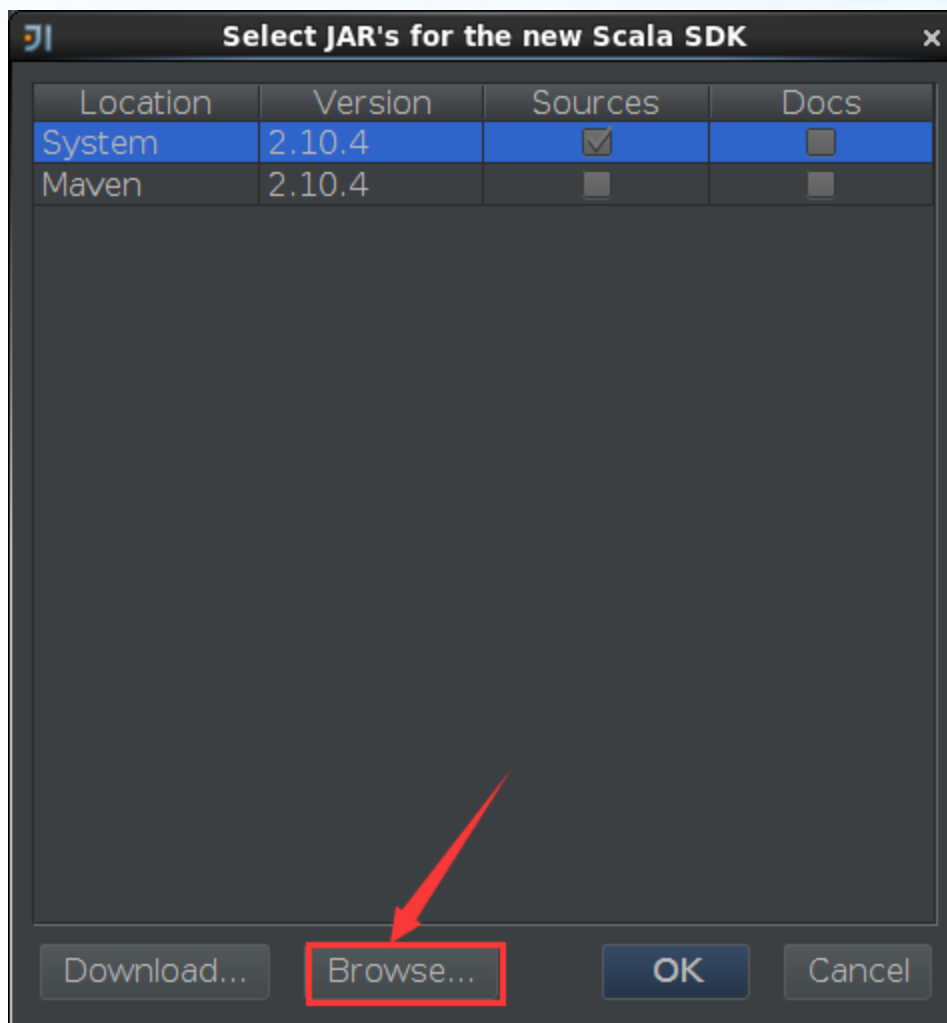# Create Scala Project

| File | Edit | View | Navigate | Code | Analyze | Refactor | Build | Run | Tools | VCS | Wind |

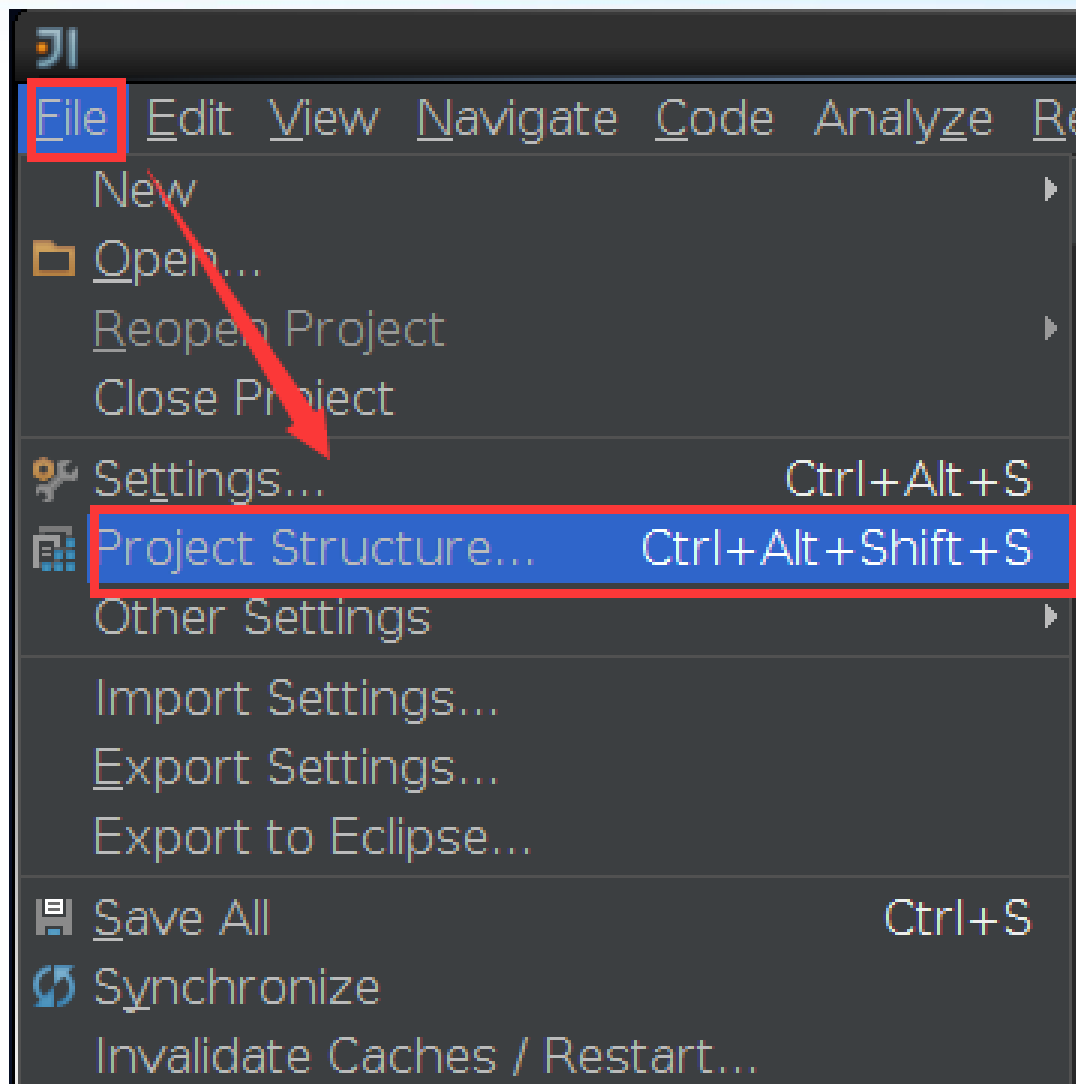| New | | | | | | ▶ | Project... |
| 📁 Open... | | | | | | | Project from Existing Sources... |
| Reopen Project | | | ▶ | | | | Project from Version Control | ▶ |
| Close Project | | | | | | | Module... |
| 🔧 Settings... | Ctrl+Alt+S | | | | | | Module from Existing Sources... |
| 📇 Project Structure... | Ctrl+Alt+Shift+S | | | | | | Ⓒ Java Class |
| Other Settings | | | ▶ | | | | 🅲 Scala Class |
| Import Settings... | | | | | | | 🅟 Package Object |
| Export Settings... | | | | | | | 🅖 Groovy Class |
| Export to Eclipse... | | | | | | | 🗎 File |
| 💾 Save All | Ctrl+S | | | | | | 📁 Package |
| 🔄 Synchronize | | | | | | | 📄 package-info.java |
| Invalidate Caches / Restart... | | | | | | | 🅗 HTML File |
| Export to HTML... | | | | | | | 🅢 Scala Worksheet |
| 🖨 Print | Ctrl+P | | | | | | 🅢 Scala Script |

# Create Scala Project

# Set Scala SDK

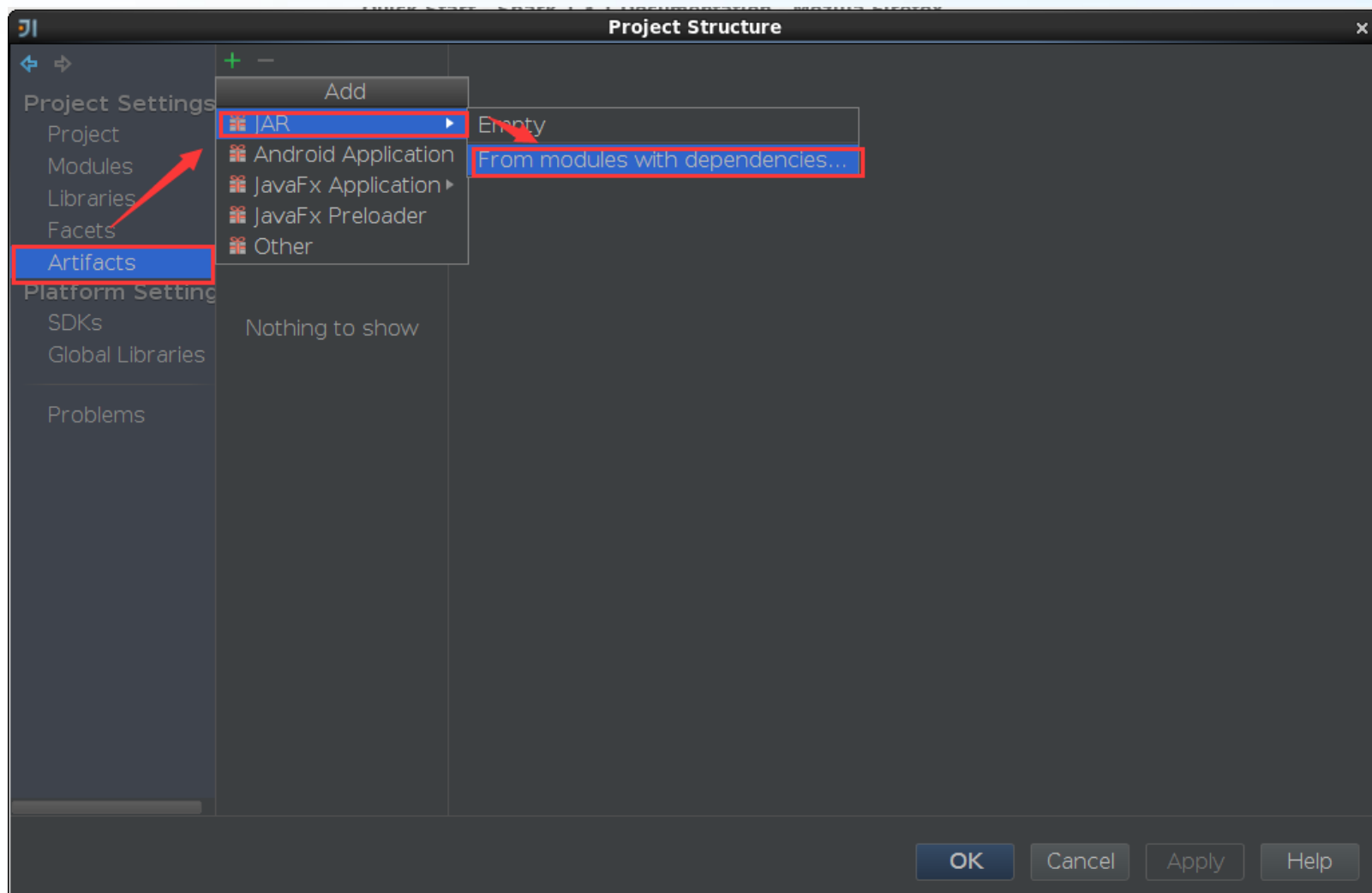# Create Scala Project

# Add Spark JARS
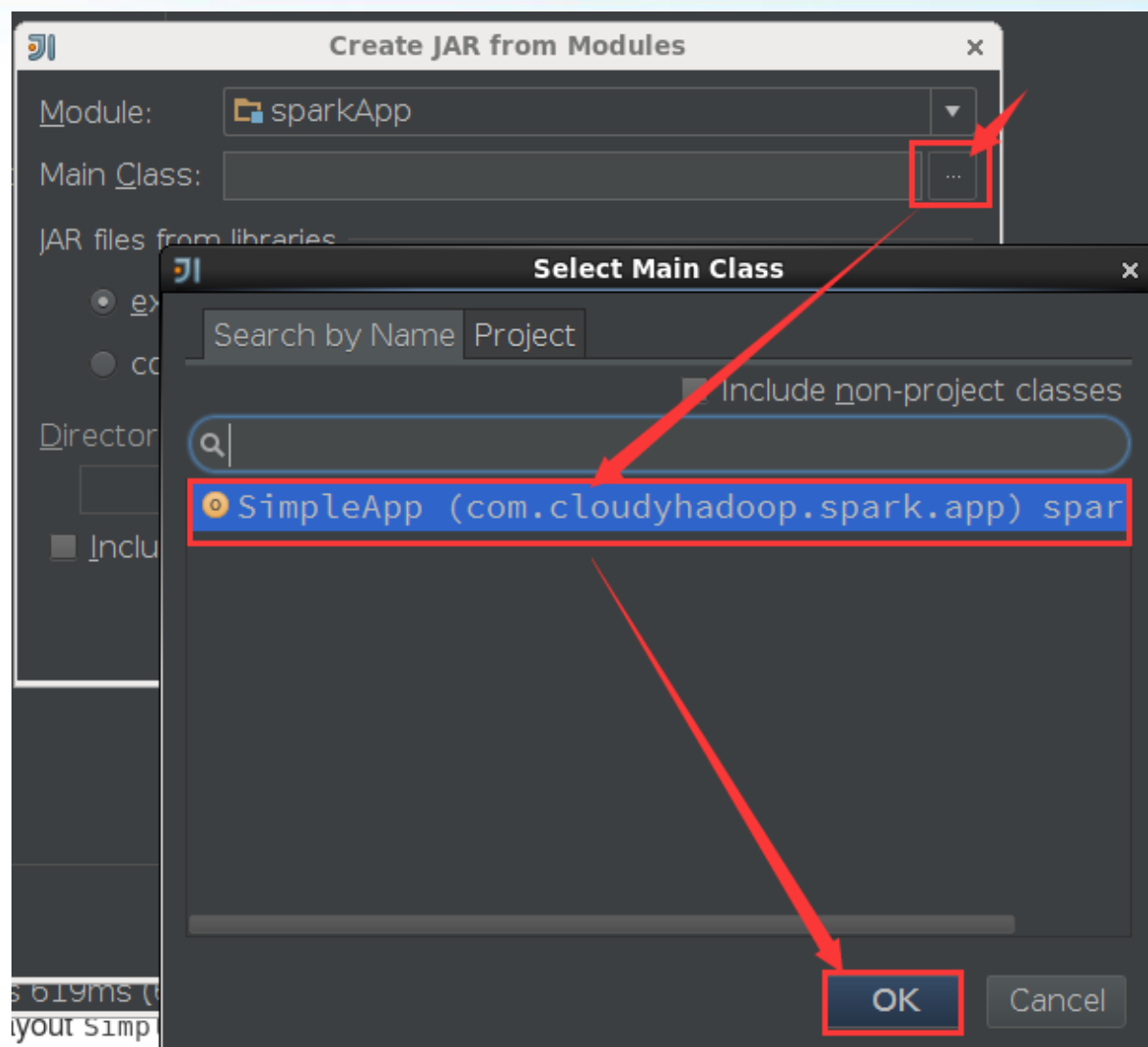
# 零基础学习Spark 1.x应用开发系列课程

## IDEA打包Spark Application

讲师-梦琪

```scala
package com.cloudyhadoop.spark.app
import org.apache.spark.SparkContext
import org.apache.spark.SparkConf
/**
 * Created by cyhp on 1/18/15.
 */
object SimpleApp {
  def main(args: Array[String]) {
    val logFile = "hdfs://hadoop-yarn.cloudyhadoop.com:8020/user/cyhp/spark/wc.input"
    val conf = new SparkConf()//
                  .setAppName("Simple Application")//
    //                .setMaster("local")
                  .setMaster("spark://hadoop-yarn.cloudyhadoop.com:7077")
    val sc = new SparkContext(conf)
    val logData = sc.textFile(logFile)
    val numAs = logData.filter(line => line.contains("a")).count()
    val numBs = logData.filter(line => line.contains("b")).count()
    println("Lines with a: %s, Lines with b: %s".format(numAs, numBs))
    sc.stop()
  }
}
```
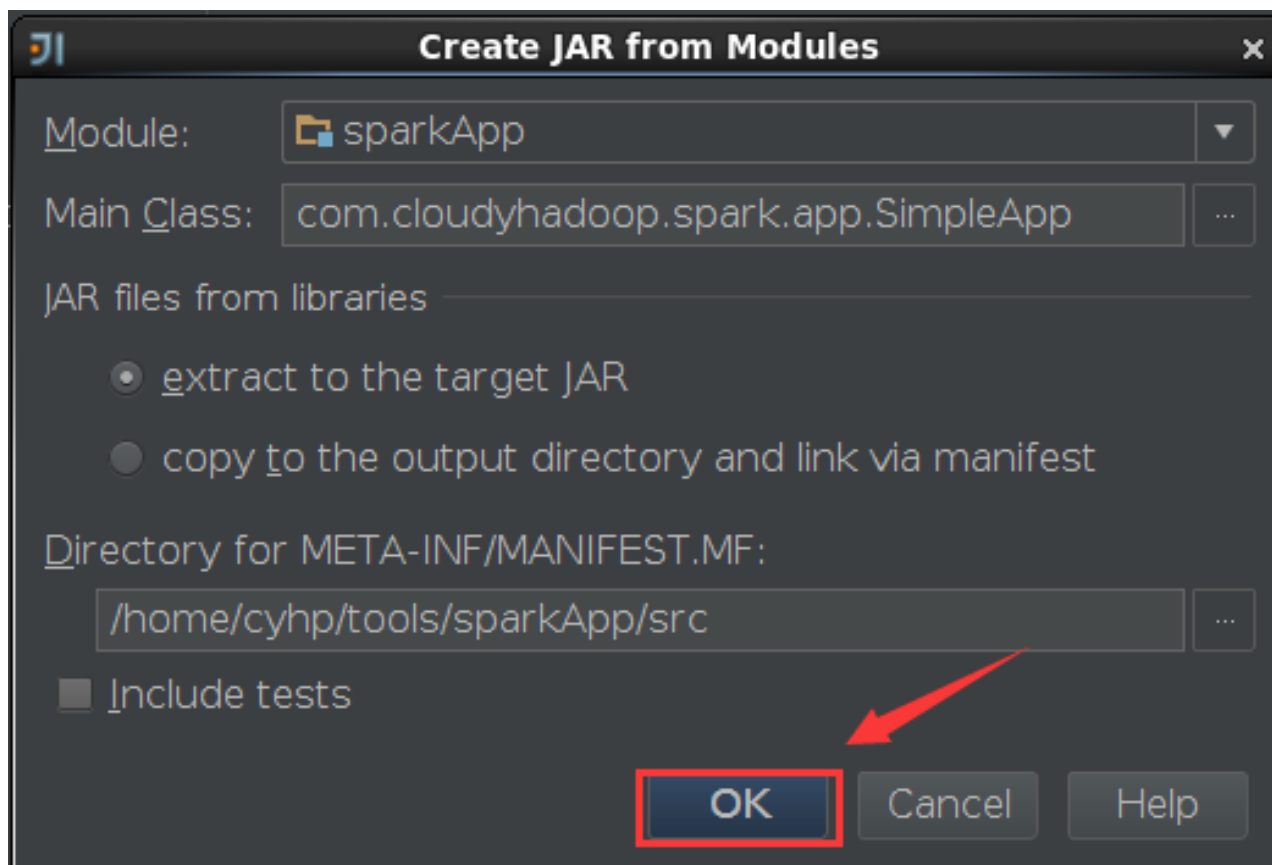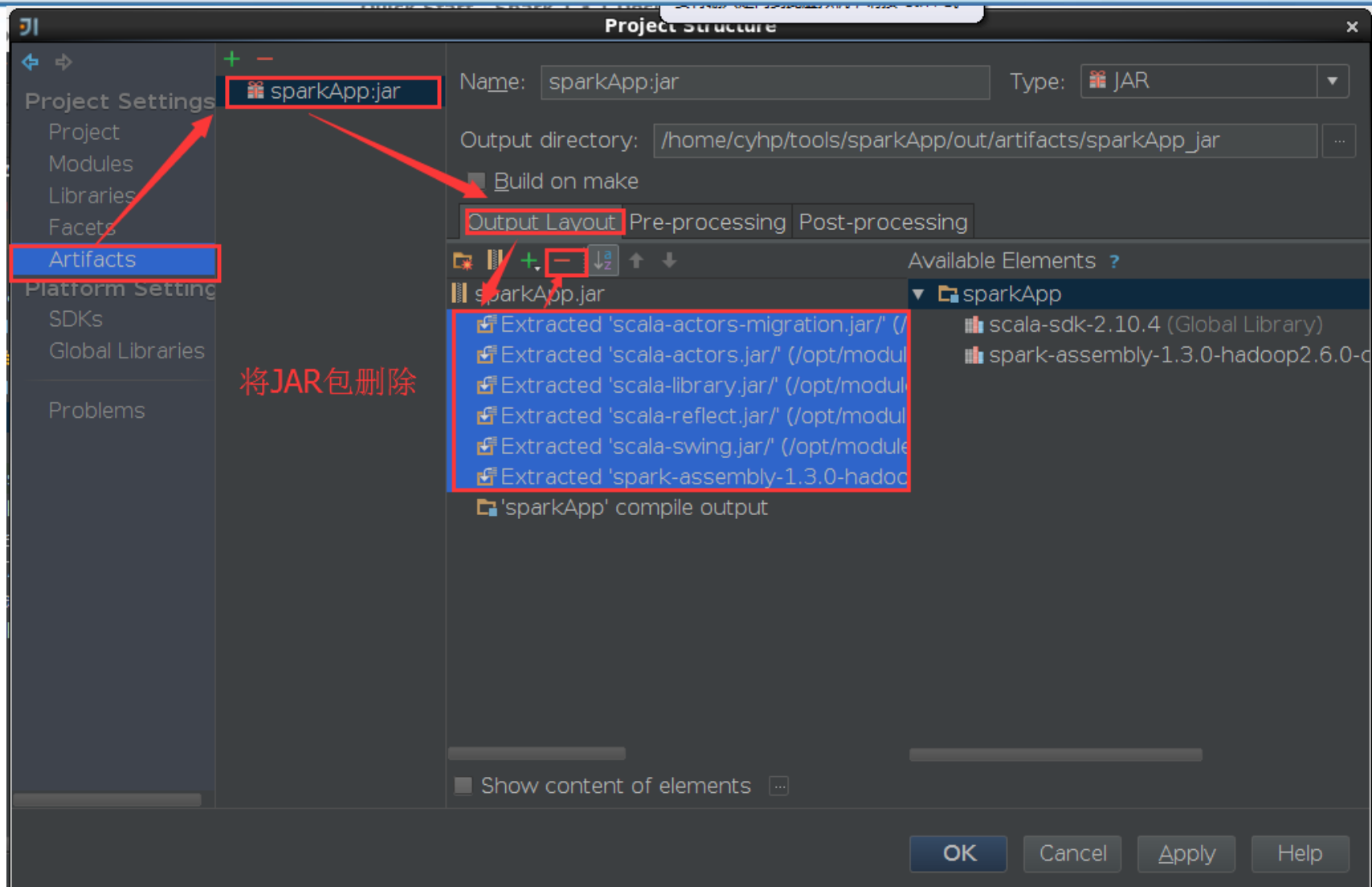
# Set Output

# Build Artifacts
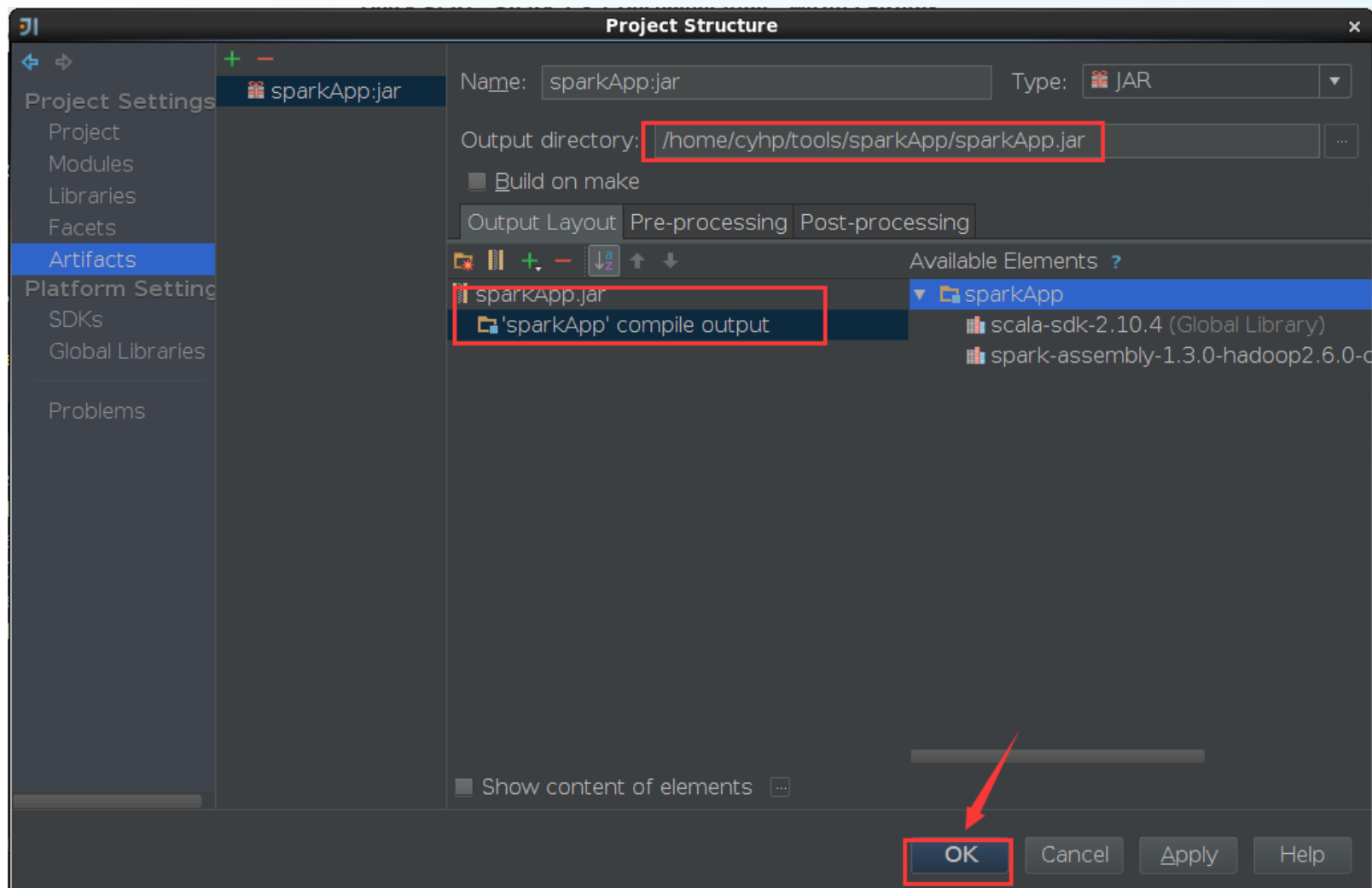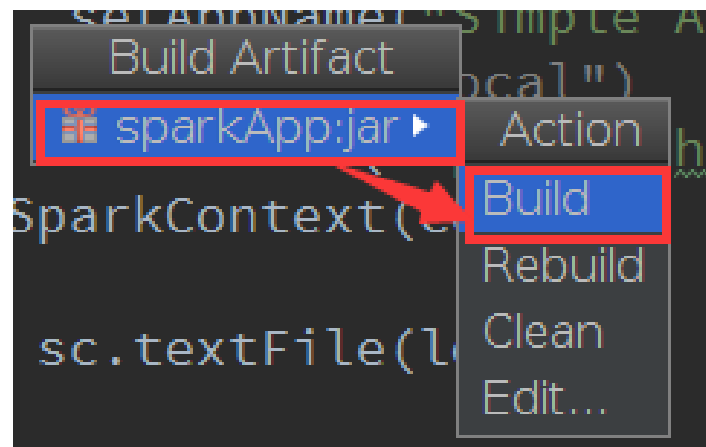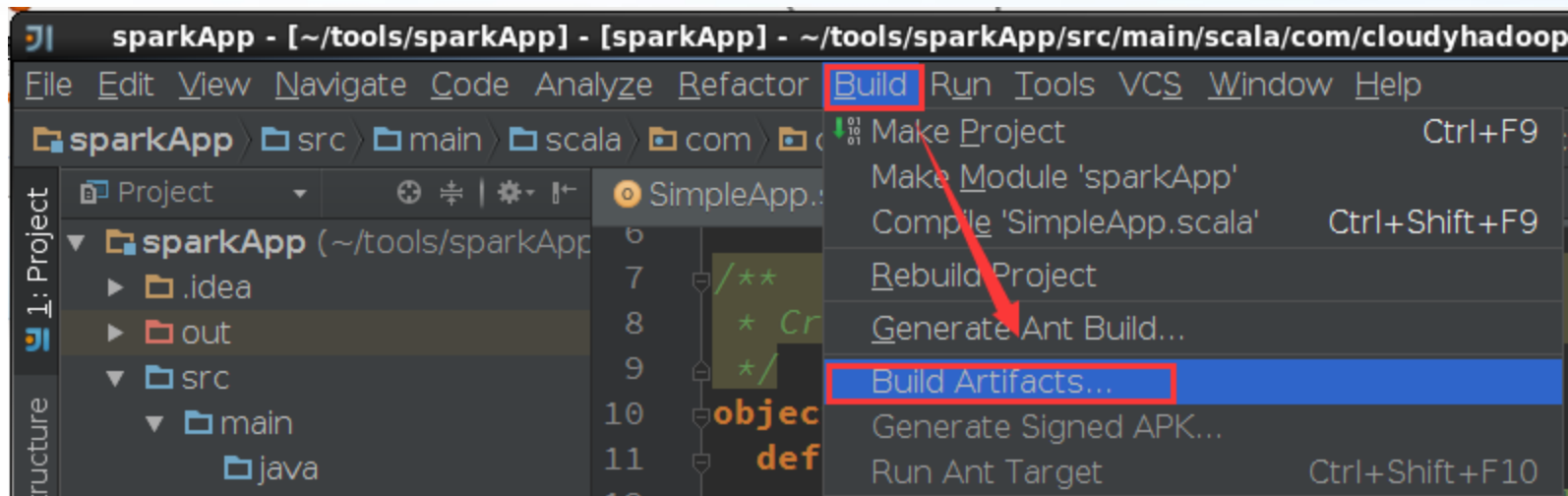
◆ 云帆大数据是国内首家坚持实时在线授课、提供高端开发课程网络培训机构。采用新兴的互联网教育模式，坚持实时在线授课模式，既继承传统教育的学习交流特点，又发挥互联网的无处不在的时空特性，将天南地北有志向的人才组织在一起学习交流，使原先孤立的个体学习，组合成有组织的学习探讨，并且把原先的学费用降低到十分之一左右，使更多的人能学习到最新的高端课程技术。云帆大数据同时是一个平台，如果你是一个学员，可以尽情的学习和交流；如果你是一个有梦想有才华的人，可以联系我们，给你提供一片驰骋的原野。

◆ 云帆大数据已推出国内首家《企业级Hadoop 2.x 应用开发课程》、《企业级Hadoop 2.x 项目实战课程》和《企业级Spark 1.x 应用开发课程》，更多其他详情，请登录我们的培训网站http://www.cloudyhadoop.com。

F.A.Q. 常见问题

实时在线授课，专业课程辅导