



# 《数据采集与预处理》

教材官网：<http://dblab.xmu.edu.cn/post/data-collection/>

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

## 《数据采集与预处理》课程介绍

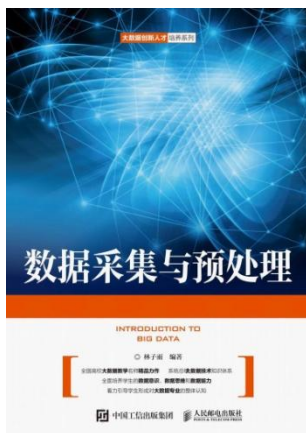
(PPT版本号：2022年1月版本)

林子雨 副教授

厦门大学计算机科学与技术系

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

主页：<http://dblab.xmu.edu.cn/linziyu>





# 提纲

- 1.主讲教师
- 2.先修课程
- 4.教材介绍
- 5.内容提要
- 6.教学大纲
- 7.配套资源

本PPT是以下教材的配套讲义  
林子雨编著《数据采集与预处理》  
人民邮电出版社

教材官网：  
<http://dbllab.xmu.edu.cn/post/data-collection>



## 数据采集与预处理

INTRODUCTION TO  
BIG DATA

◎ 林子雨 编著

全国高校大数据教学名师精品力作 系统总结大数据技术知识体系  
全面培养学生的数据意识、数据思维和数据处理能力  
着力引导学生形成对大数据专业的整体认知



中国工信出版集团



人民邮电出版社  
POSTS & TELECOM PRESS



# 主讲教师



主讲教师：林子雨

中国高校首个“数字教师”提出者和建设者

2009年7月从事教师职业以来

累计**免费**网络发布超过**1500万**字高价值教学和科研资料

网络浏览量超过**1500万**次



数字教师LOGO



# 先修课程

在学习本课程之前，建议学生已经学习过Python编程  
(不是必须)



# 教材介绍

## 《数据采集与预处理》

厦门大学 林子雨 编著

名师精品，多年大数据教学实践的厚积薄发  
深入浅出，清晰呈现数据采集与预处理技术体系  
实例丰富，快速掌握各种技术的基本使用方法  
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

教材官网：<http://dbllab.xmu.edu.cn/post/data-collection/>







# 内容提要

- 本课程详细阐述了大数据领域数据采集与预处理的相关理论和技术
- 内容包括概述、大数据实验环境搭建、网络数据采集、分布式消息系统Kafka、日志采集系统Flume、数据仓库中的数据集成、ETL工具Kettle、使用pandas进行数据清洗



# 教学大纲

章（或节）	主要内容	学时安排
第1章 概述	数据的概念、类型、组织形式，数据分析过程以及数据采集与预处理的任務，数据采集、数据清洗、数据转换和数据脱敏	2
第2章 大数据实验环境搭建	Python的安装和使用方法，JDK的安装以及MySQL的安装和使用方法，Hadoop的安装和使用方法	4
第3章 网络数据采集	网络爬虫的基本概念，网络爬虫、网络爬虫的类型以及反爬机制，网页基础知识，如何使用Python实现HTTP请求，如何定制requests以及如何解析网页，3个网络爬虫的具体实例	4
第4章 分布式消息系统 Kafka	Kafka在大数据生态系统中的作用以及Kafka与Flume的区别与联系，Kafka的相关概念、Kafka的安装和使用以及如何使用Python操作Kafka，Kafka与MySQL的组合使用	2
第5章 日志采集系统 Flume	Flume的安装和使用方法以及Kafka和Flume的组合使用方法，采集日志文件到HDFS以及采集MySQL数据到HDFS的方法	2
第6章 数据仓库中的数据集成	数据仓库的概念，包括传统的数据仓库和实时主动数据仓库，数据仓库中的数据集成，包括数据集成方式、数据分发方式和数据集成技术，两种具有代表性的数据集成技术，即ETL和CDC	2
第7章 ETL工具Kettle	Kettle的基本概念、基本功能和安装方法，如何使用Kettle进行数据抽取、转换和加载	4
第8章 使用pandas进行数据清洗	如何使用pandas进行数据清洗，NumPy的基本使用方法pandas的数据结构和一些基本功能，如何使用pandas进行汇总和描述统计、处理缺失数据等	4
合计		24



# 课程配套教学资源

所有资料全部免费共享  
支持电脑和手机浏览



## 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精华资源<http://dblab.xmu.edu.cn/post/8197/>



扫一扫访问平台主页

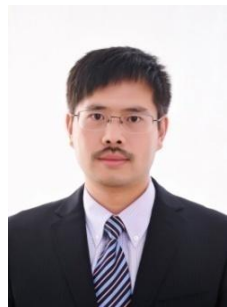


扫一扫观看3分钟FLASH动画宣传片





# 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与技术系

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），全国高校知名大数据教师，现为厦门大学计算机科学系副教授，厦门大学信息学院实验教学中心主任，曾任厦门大学信息科学与技术学院院长助理、晋江市发展和改革局副局长。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、2017年度和2020年度厦门大学教学类奖教金获得者，荣获2019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1000万字高价值的研究和教学资料，累计网络访问量超过1000万次。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400万次，累计访问量超过1500万次。



# 附录B：大数据学习路线图



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



# 附录C：林子雨大数据系列教材



林子雨大数据系列教材

用于导论课、专业课、实训课、公共课

了解全部教材信息：<http://dbllab.xmu.edu.cn/post/bigdatabook/>





# 附录D：《大数据导论（通识课版）》教材

## 开设全校公共选修课的优质教材



本课程旨在实现以下几个培养目标：

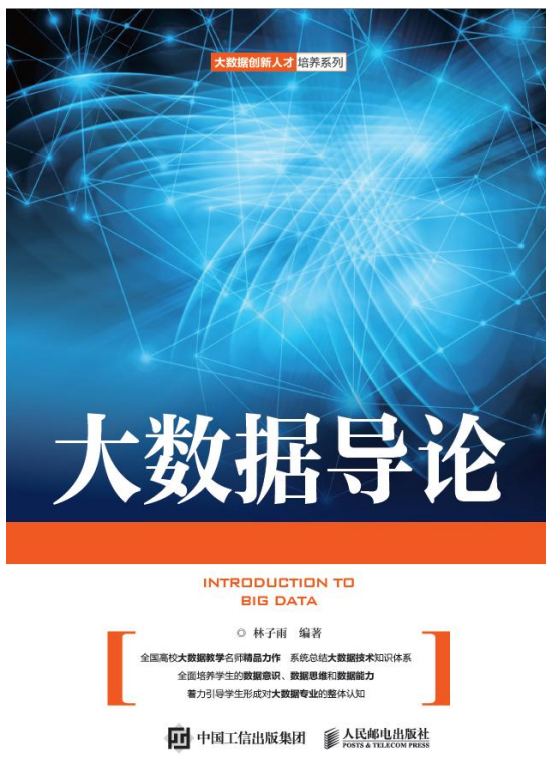
- 引导学生步入大数据时代，积极投身大数据的变革浪潮之中
- 了解大数据概念，培养大数据思维，养成数据安全意识
- 认识大数据伦理，努力使自己的行为符合大数据伦理规范要求
- 熟悉大数据应用，探寻大数据与自己专业的应用结合点
- 激发学生基于大数据的创新创业热情

高等教育出版社 ISBN:978-7-04-053577-8 定价：32元 版次：2020年2月第1版  
教材官网：<http://dbl原因.xmu.edu.cn/post/bigdataintroduction/>



# 附录E：《大数据导论》教材

- 林子雨 编著 《大数据导论》
  - 人民邮电出版社，2020年9月第1版
  - ISBN:978-7-115-54446-9 定价：49.80元
- 教材官网：<http://dbl原因.xmu.edu.cn/post/bigdata-introduction/>



开设大数据专业导论课的优质教材



扫一扫访问教材官网



# 附录F：《大数据技术原理与应用（第3版）》教材

《大数据技术原理与应用——概念、存储、处理、分析与应用（第3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:978-7-115-54405-6 定价：59.80元

全书共有17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、Flink、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、MapReduce、Spark和Flink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。

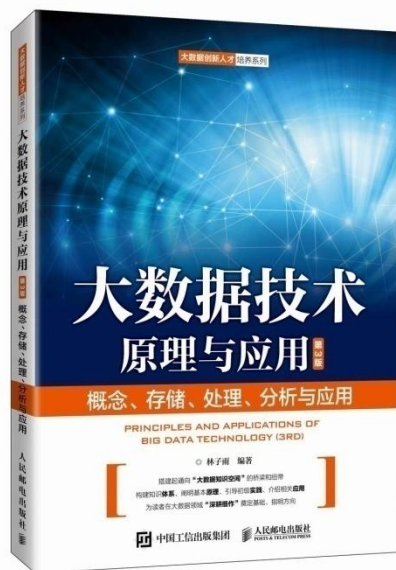
本书可以作为高等院校计算机专业、信息管理等相关专业的大数据课程教材，也可供相关技术人员参考、学习、培训之用。

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://dbllab.xmu.edu.cn/post/bigdata3>



扫一扫访问教材官网







# 附录G：《大数据基础编程、实验和案例教程（第2版）》

本书是与《大数据技术原理与应用（第3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，八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

林子雨编著《大数据基础编程、实验和案例教程（第2版）》

清华大学出版社 ISBN:978-7-302-55977-1 定价：69元 2020年10月第2版



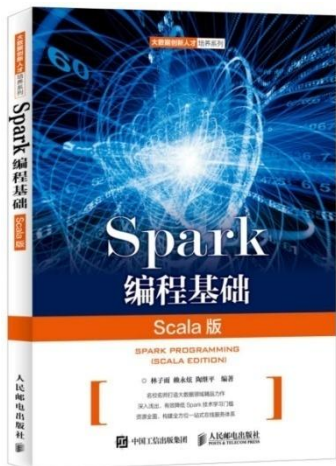
# 附录H：《Spark编程基础（Scala版）》

## 《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/>

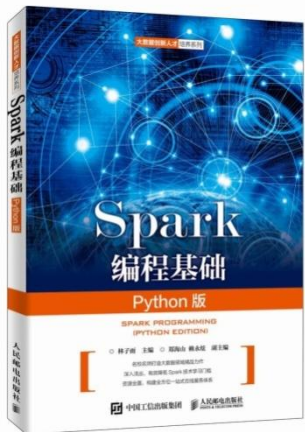


本书以Scala作为开发Spark应用程序的编程语言，系统介绍了Spark编程的基础知识。全书共8章，内容包括大数据技术概述、Scala语言基础、Spark的设计与运行原理、Spark环境搭建和使用方法、RDD编程、Spark SQL、Spark Streaming、Spark MLlib等。本书每个章节都安排了入门级的编程实践操作，以便读者更好地学习和掌握Spark编程方法。本书官网免费提供了全套的在线教学资源，包括讲义PPT、习题、源代码、软件、数据集、授课视频、上机实验指南等。



# 附录I: 《Spark编程基础 (Python版)》

## 《Spark编程基础 (Python版)》



厦门大学 林子雨, 郑海山, 赖永炫 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:978-7-115-52439-3

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



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、上机实验指南等。



# 附录J：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



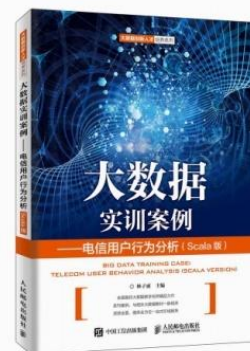


# 附录K：高校大数据实训课程系列案例教材

为了更好地满足高校开设大数据实训课程的教材需求，厦门大学数据库实验室林子雨老师团队联合企业共同开发了《高校大数据实训课程系列案例》，目前已经完成开发的系列案例包括：

- 《电影推荐系统》（已经于2019年5月出版）
- 《电信用户行为分析》（已经于2019年5月出版）
- 《实时日志流处理分析》
- 《微博用户情感分析》
- 《互联网广告预测分析》
- 《网站日志处理分析》

系列案例教材将于2019年陆续出版发行，教材相关信息，敬请关注网页后续更新！  
<http://dbllab.xmu.edu.cn/post/shixunkecheng/>



扫一扫访问大数据实训课程系列案例教材主页

The background of the slide features a blue gradient with several white 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, resting their head on their hand. In the bottom left corner, two more people are shown in profile, one appearing to be speaking or gesturing towards the other.

**Thank You!**

**Department of Computer Science, Xiamen University, 2022**