

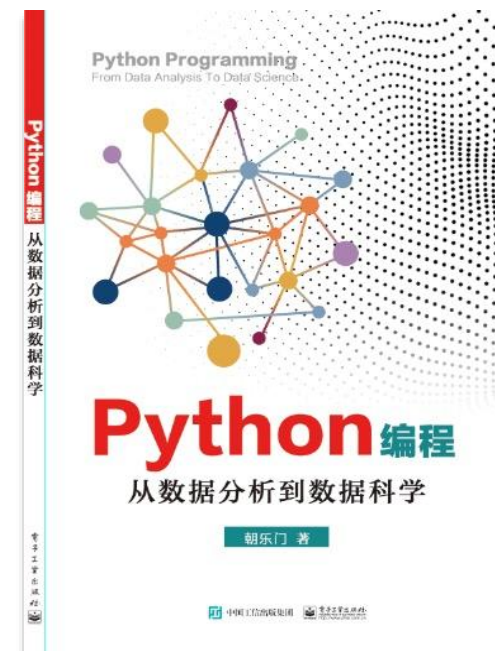
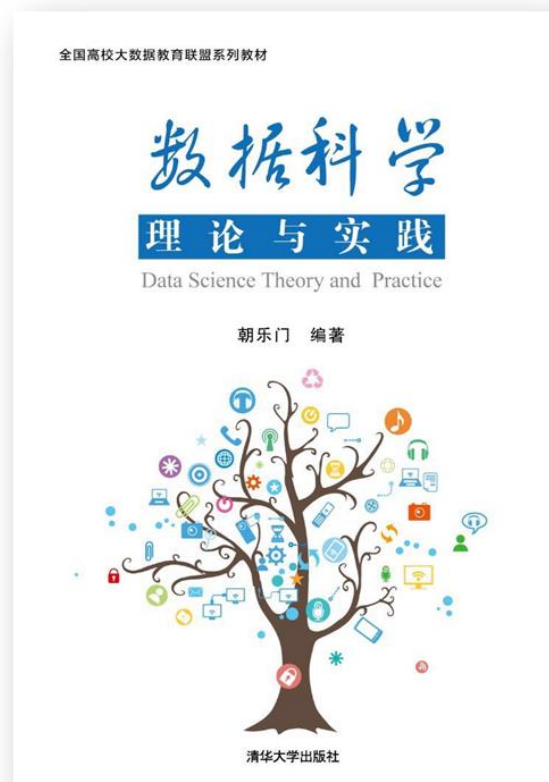
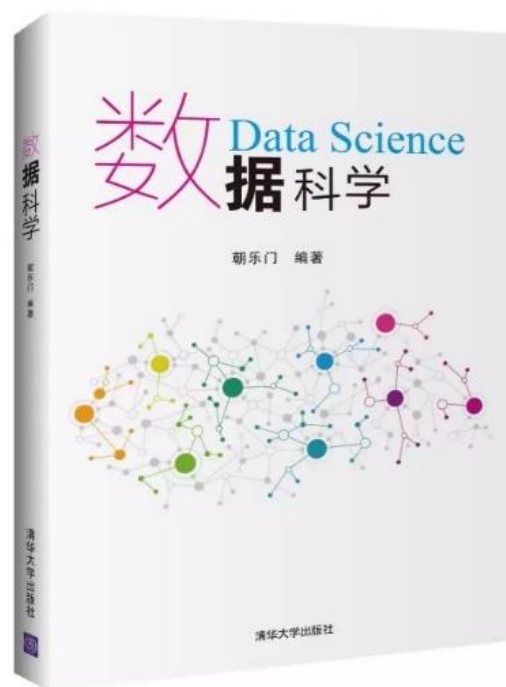
# 《数据科学导论》的建设与开源

朝乐门

2018年10月14日·大连



# 我的工作



# 报告提纲

## 《数据科学导论》建设与开源

1.跟踪调研

2.如何建设

3. 开源倡议

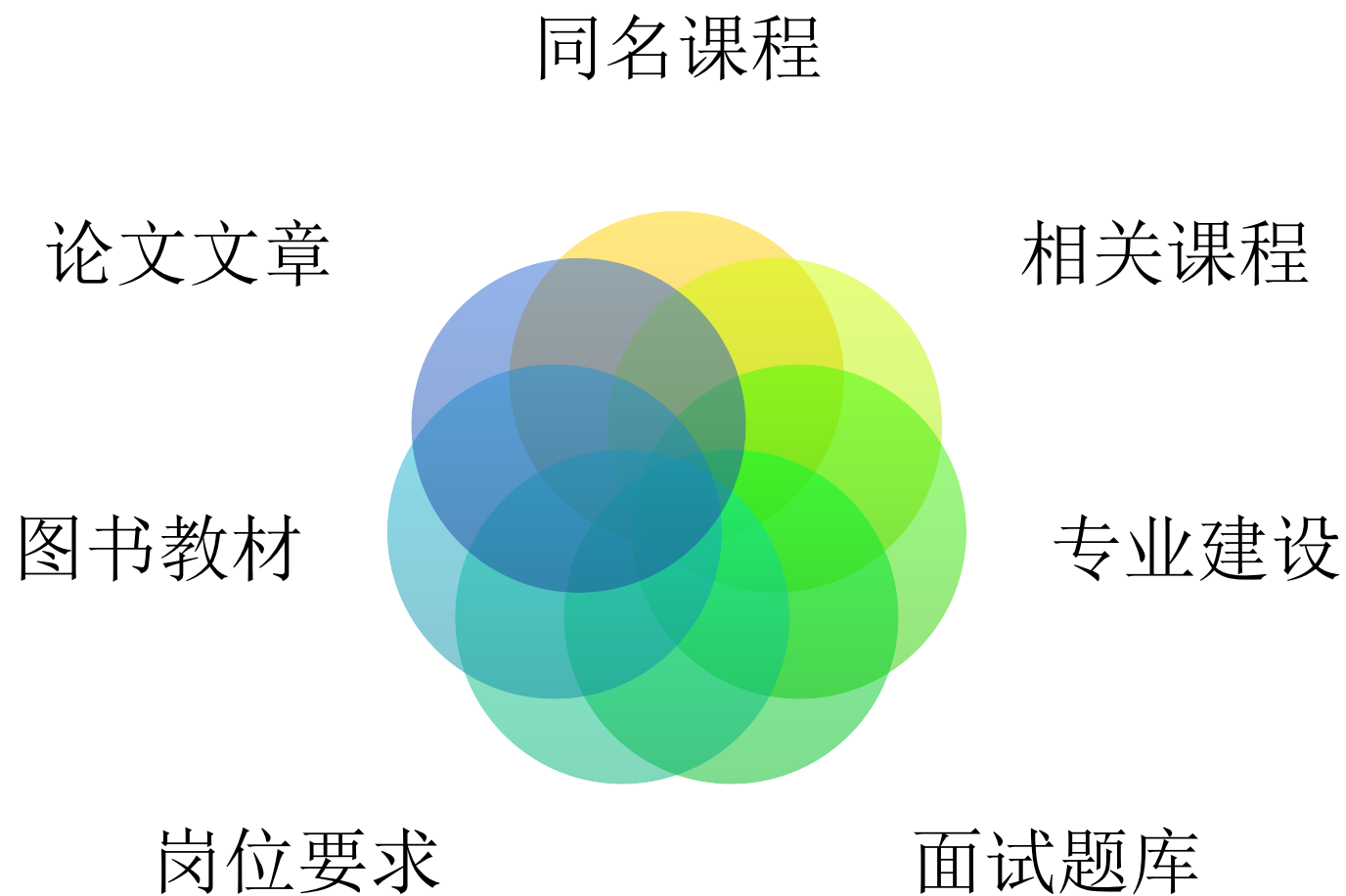
# 《数据科学导论》建设与开源

## —— 第一部分 跟踪调研

朝乐门



# 我们的跟踪调研



# 连续3年跟踪调研

表 1 数据科学的课程调研

课程名称	年份	形式	学校	开课教师	选课要求
Data Science: Large-scale Advanced Data Analysis	2011	面授	佛罗里达大学	Daisy Zhe Wang	硕士
Data Science and Analytics Thought Leaders Introduction to Data Science	2012	面授	加州大学伯克利分校	Ram Akella 等	不限
Introduction to Data Science	2012	面授	哥伦比亚大学	Rachel Schutt	不限
Introduction to Data Science	2013	面授	谢菲尔德大学	Paul Clough	数据相关/硕士
Data Science(Coursea)	2014	网授	约翰·霍普金斯大学	Roger D. Peng 等	不限
Executive Data Science(Coursea)	2014	网授	约翰·霍普金斯大学	Roger D. Peng 等	不限
Data Science at Scale (Coursea)	2014	网授	华盛顿大学	Bill Howe	不限
Data Science	2014	面授	哈佛大学	Rafael Irizarry 等	本科
Intro to Data Science	2014	面授	纽约大学	Brian D'Alessandro	不限
大数据科学与应用系列讲座(MOOC 学院)	2015	网授	清华大学	李军	不限
Foundations of Data Science	2015	面授	加州大学伯克利分校	John DeNero	不限
Data Sciences Basic	2015	面授	美国东北大学	Akira Suzuki	不限
Fundamentals of Data Science	2015	面授	慕尼黑大学	Goeran Kauermann	统计与科学相关
A Practical Approach to Data Science	2016	面/网授	哈佛大学	Ramon Mata-Toledo	不限
Introduction to Computational Thinking and Data Science (edx)	2016	网授	麻省理工学院(MIT)	Eric Grimson 等	不限
Process Mining: The Practice of Data Science (Coursea)	2016	网授	埃因霍芬理工大学	Wil van der Aalst	硕士
Data Science	2016	面授	法国圣艾蒂安大学	Marc Sebban	不限
Fundamentals of Data Science	2017	面授	牛津大学	Julian Gallop	不限
数据科学	2017	面授	中国人民大学	朝乐门	不限
Data Science	不详	面授	伦敦大学	Aysha Chaudhary	数据相关/硕士

朝乐门等. 全球数据科学课程建设现状的实证分析\*[J]. 数据分析与知识发现, 2017, 1(6): 12-21.

NYU



## Intro to Data Science (Fall 2014)

- Understanding Data and Data Science (DM process)
- Learning from Data (ML)
- Practicing the “Science” of Data Science (testing hypotheses ,technical design decisions)
- Special Topics in Applying Data Science (written and oral communication)
- URL:<https://cds.nyu.edu/ds-ga-1001-intro-data-science/>

## Introduction to Data Science(Spring 2018)

- Introduction
- Supervised learning
- Unsupervised learning
- Performance measures
- Text mining
- Feature extraction and selection, review for midterm
- Differentiable programming, neural networks
- Deep neural networks, non-linear PCA
- Bayesian and variational inference
- Gradient boosting
- Meta data and learning Lab
- Data visualization



Iddo Drori



# NYU

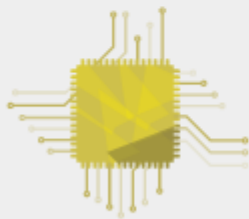


## Master of Science in Data Science

Learn More!

Data science creates meaning from vast amounts of complex data.

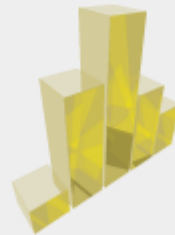
Using automated analytical methods, it reveals patterns humans alone might never see. Data science combines aspects of:



COMPUTER  
SCIENCE



APPLIED  
MATHEMATICS



STATISTICS



MACHINE  
LEARNING



VISUALIZATION

来源: <https://cds.nyu.edu/>



UW



## Introduction to Data Science (Coursera, by Bill Howe)

- Introduction
- **Data Manipulation at Scale**
- Analytics
- **Communicating Results (Visualization, Data Products, Privacy)**
- Special Topics (Graph Analytics + Guest Lectures)
- URL: <https://www.class-central.com/course/coursera-introduction-to-data-science-451#>



## Introduction to Data Science/LIS 572

- **experimental design**
- data collection and storage
- **basic analytics**
- machine learning
- **data visualization**
- URL: <https://www.washington.edu/students/crscat/lis.html>

# Oxford



## Data Science: An Introduction (21 Jan 2019 - , 10weeks )

- Applications, context, history, rationale. Summary of the course
- Further statistics, introduction to R and identifying some datasets.
- Data science techniques (and also in subsequent weeks); review of student chosen datasets; introduction to free data science package
- data preparation
- overview of workflow
- problems of big data
- Review of techniques and software not already covered in detail
- More on big data; web-based software, Wrap up, review
- URL:<https://www.conted.ox.ac.uk/courses/data-science-an-introduction>

## Introduction to Data Science (Short Courses, 2days)

- Range of possibilities of data science
- Role of statistics in data science and practical
- Outline of a free data science package and opportunity for an initial practical
- Selection of data science by techniques and practical
- Outline of further free data science package; further data science techniques
- Big data concepts



- URL:<https://www.conted.ox.ac.uk/courses/introduction-to-data-science>

# Columbia University in the City of New York



## Intro to DS

- What is Data Science
- Intro to statistical thinking
- Introduction to Bayesian modeling
- Exploratory Data Analysis and Visualization
- Algorithms
- Machine learning
- Privacy and data security
- Data Engineering
- Internet of Things
- Project presentations

## 相关课程

- Computer Systems for Data Science
- Machine Learning for Data Science
- Algorithms for Data Science
- Probability & Statistics for Data Science
- Exploratory Data Analysis & Visualization
- Data Science Capstone & Ethics

URL:<https://www.dropbox.com/s/oiv21udbnv8ieug/syllabus-g5705-fall2016.pdf?dl=0>

MIT



## Introduction to Data Science( 6.s07)

- Introduction
- **Learning of distributions and their parameters**
- Hypothesis testing
- **Regression and prediction**
- Classification
- **Dynamical models**
- Conclusion

### URL

- <http://web.mit.edu/6.s077/www/Intro-Doc-Stellar.pdf>

#### Dynamical models (3)

- 1. Linear models (autoregressive and state space models)
- 2. Nonlinear models: maximum likelihood and empirical risk minimization
- 3. Difficulties with nonlinear models (e.g., unobserved variables, **causality**)

# Stanford



## Data Science 101

- Data science: what is the buzz about?
- Visualization tools
- Numerical summaries of data
- Prediction (linear and logistic regression)
- Sampling variability and uncertainty of statistical estimates (Inference I)
- Testing statistical hypotheses (Inference II)
- Safeguarding reproducibility and multiple hypothesis testing
- **Causality and experimentation**
- Machine learning
- Ongoing Challenges to Data Science + Review
- URL: <https://web.stanford.edu/class/stats101/>

Summer Session

## SCI 01 A, Introduction to Data Science

- What is Data Science, Getting Started with R Malaria Detection
- Data Visualization, Basics in Programming with R Mice Project
- Predictive Modeling, Regression Analysis Data Science in Healthcare
- Classification Methods Time Series Prediction
- Feature Selection, Clustering Techniques Internet of Things
- Association Rule Mining, Web Scraping -
- Advanced Analysis of Models, Shiny: Interactive Web Apps in R Mining Health Insurance Data
- Advanced Topics, Project Presentation
- [URL:https://continuingstudies.stanford.edu/coursework/document/9356/?f=20181\\_SCI%2001%20A\\_Syllabus.pdf](https://continuingstudies.stanford.edu/coursework/document/9356/?f=20181_SCI%2001%20A_Syllabus.pdf)







# UC Berkeley

## Introduction to Data Science (2014)

- Introduction
- **Data Science Process**
- **Data Preparation**
- Tabular Data
- Data Cleaning
- Natural Language Processing
- Exploratory Data Analysis
- ML
- Visualization
- **Graph Processing**
- Project Posters
- <https://bcourses.berkeley.edu/courses/1267848/>

## Introduction to Data Science(2019)

- The data science **lifecycle**: from business understanding to customer acceptance
- Data science **roles**
- Basic R programming using Rstudio
- Exploratory data analysis
- Data visualizations
- Statistical concepts such as hypothesis testing (**A/B testing**) and inference
- Reproducible computation with Git, Github and R
- Overview of predictive analytics
- **Business use cases**
- <https://extension.berkeley.edu/search/publicCourseSearchDetails.do?method=load&courseId=39287923#outline>

Certificate Program  
in Data Science

# Harvard



## Data Science Professional Certificate Program

- R Basics
- Visualization
- Probability
- Inference and Modeling
- Productivity Tools
- Wrangling
- Linear Regression
- Machine Learning
- Capstone
- URL:<https://www.edx.org/professional-certificate/harvardx-data-science>

Professional Certificate Program

## Data Science: Wrangling

- Importing data into R from different file formats
- Web scraping
- How to tidy data using the tidyverse to better facilitate analysis
- String processing with regular expressions (regex)
- Wrangling data using dplyr
- How to work with dates and times as file formats
- Text mining
- URL:<https://www.edx.org/course/data-science-wrangling-harvardx-ph125-6x>



# CMU



## Data Science (67-364/95-885)

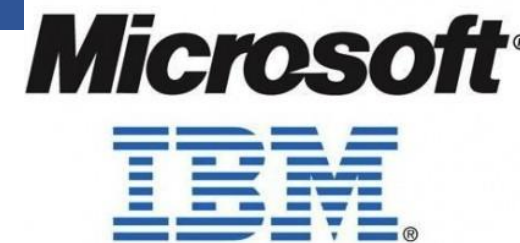
- **I. Data Science**
- 1. **The Data Science pipeline**
- 2. Exploratory Data Analysis (scraping and visualization)
- 3. Machine learning (supervised and unsupervised approaches)
- 4. **Analytics tasks**: classification, prediction, recommendation, clustering
- **II. Big Data**
- 5. The Hadoop platform (HDFS, map-reduce)
- 6. Cloud based platforms such as AWS
- 7. Data processing on Hadoop (MrJob and Pig)
- 8. Spark

URL:<https://www.cmu.edu/information-systems/images/syllabi/67364-95885-data-science-s17.pdf>

## Applied Data Science(95-852)

- Introduction to Applied Data Science. Definitions and application environment.
- **What challenges and needs could be addressed with Data Science? What can be achieved by applying Data Science in practice**
- and what are the dimensions of attainable impact in business scenarios? How to estimate and validate added value?
- **Reconciling prior knowledge against observations: Bayesian analytics.**
- Building efficient models from complex data. Regularization. Graphical models.
- **Predictive analytics with simultaneous use of multiple models.**
- Time series analysis, featurization, and forecasting. Statistical process control. Modeling sequential patterns.
- ....

[https://api.heinz.cmu.edu/media/95-852\\_Applied\\_Data\\_Sci\\_Syllabus\\_F17.pdf](https://api.heinz.cmu.edu/media/95-852_Applied_Data_Sci_Syllabus_F17.pdf)



## 其他

### IBM: introduction to data science

- **What** is Data Science?
- **Open Source tools** for Data Science
- Data Science **Methodology**
- **Databases and SQL** for Data Science
- <https://www.coursera.org/specializations/introduction-data-science>

### Microsoft: introduction to data science

- How the Microsoft **Data Science curriculum works**
- How to navigate the curriculum and **plan your course schedule**
- Basic **data exploration and visualization** techniques in Microsoft Excel
- **Foundational statistics** that can be used to analyze data
- <https://www.edx.org/course/introduction-to-data-science-3>

产品  
证书

# 分析启示



课程名称	A Practical Approach to Data Science	Intro to Data Science	Intro to Data Science	A Crash Course in Data Science	Intro to Data Science	Fundamentals of Data Science
开设学校	哈佛大学	华盛顿大学	Udacity 平台	约翰斯·霍普金斯大学	哥伦比亚大学	牛津大学
基础理论	√	√	√	√	√	√
数据加工	√	√	√	√	√	√
统计分析	√	√	×	√	√	√
机器学习	×	√	×	√	√	×
可视化与沟通	√	√	√	×	×	×
数据管理	√	√	×	×	√	√
数据计算	√	√	√	×	√	×
数据科学工具	√	√	√	√	√	√

# 到底什么是数据科学？



图片来源: <https://goo.gl/images/r8wSsm>

# 数据科学的定义

## ■ 数据科学 ≈ 大数据科学 (+ 小数据科学)

### 新兴科学

- 是一门将“现实世界”映射到“数据世界”之后，在“数据层次”上研究“现实世界”的问题，并根据“数据世界”的分析结果，对“现实世界”进行预测、洞见、解释或决策的**新兴科学**；

### 交叉性学科

- 是一门以“数据”，尤其是“大数据”为研究对象，并以数据统计、机器学习、数据可视化等为理论基础，主要研究数据预处理、数据管理、数据计算、数据产品开发等活动的**交叉性学科**；

### 独立学科

- 是一门以实现“从数据到信息”、“从数据到知识”和（或）“从数据到智慧”的转化为主要研究目的，以“数据驱动”、“数据业务化”、“数据洞见”、“数据产品研发”和（或）“数据生态系统建设”为主要研究的**独立学科**；

### 知识体系

- 是一门以“数据时代”，尤其是“大数据时代”面临的新挑战、新机会、新思维和新方法为核心内容的，包括新的理论、方法、模型、技术、工具、应用**最佳实践**在内的



朝乐门,邢春晓,张勇. 数据科学研究的现状与趋势[J]. 计算机科学, 2018, 45(1): 1-13

(来源: 朝乐门. 数据科学[M]. 北京: 清华大学出版社, 2016: 20-21)



# 数据科学的学科定位

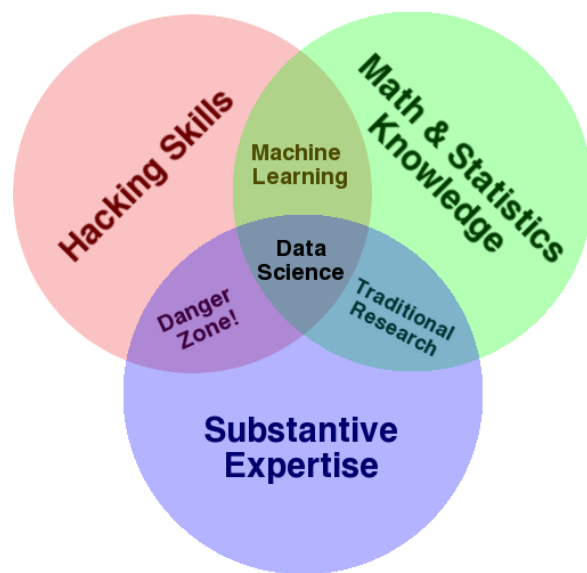


图1 Drew Conway的数据科学韦恩图（2010）

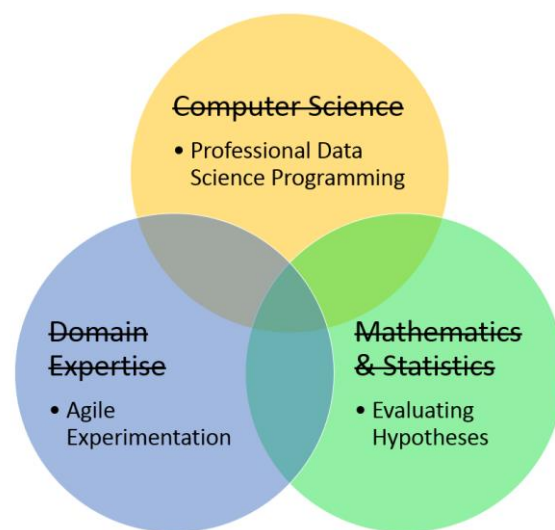


图2 Jerry Overton的数据科学韦恩图（2016）

来源:

图1: Schutt R, O'Neil C. Doing data science: Straight talk from the frontline[M]. O'Reilly Media, Inc., 2013:7.

图2: Jerry Overton.Going Pro in Data Science [M].O'Reilly Media, Inc,2016:12.

## 黑客精神

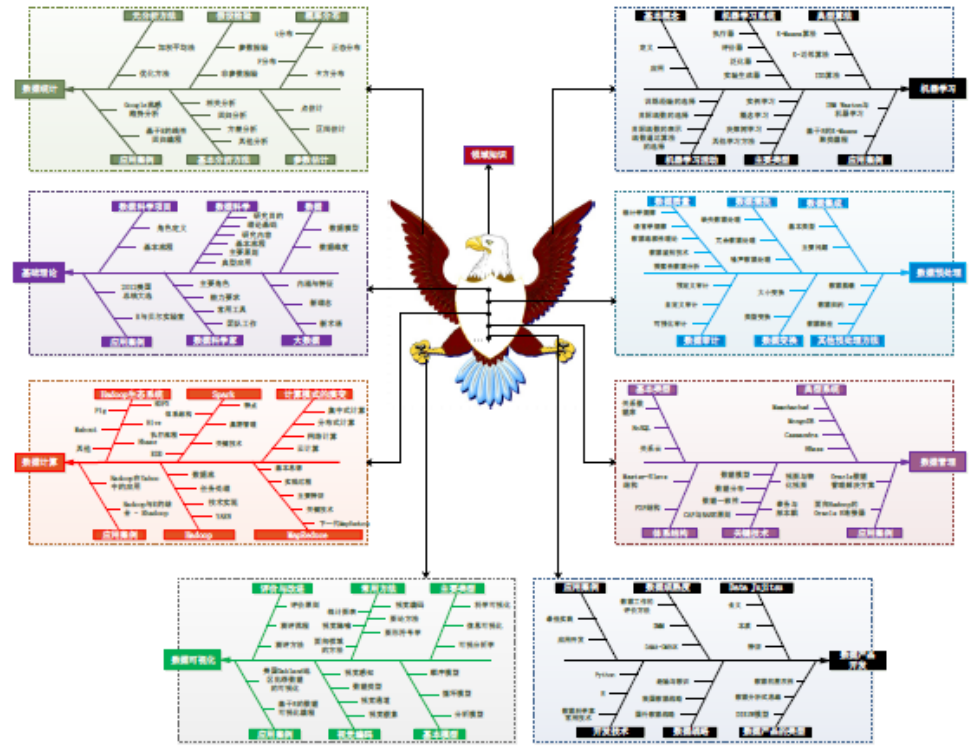
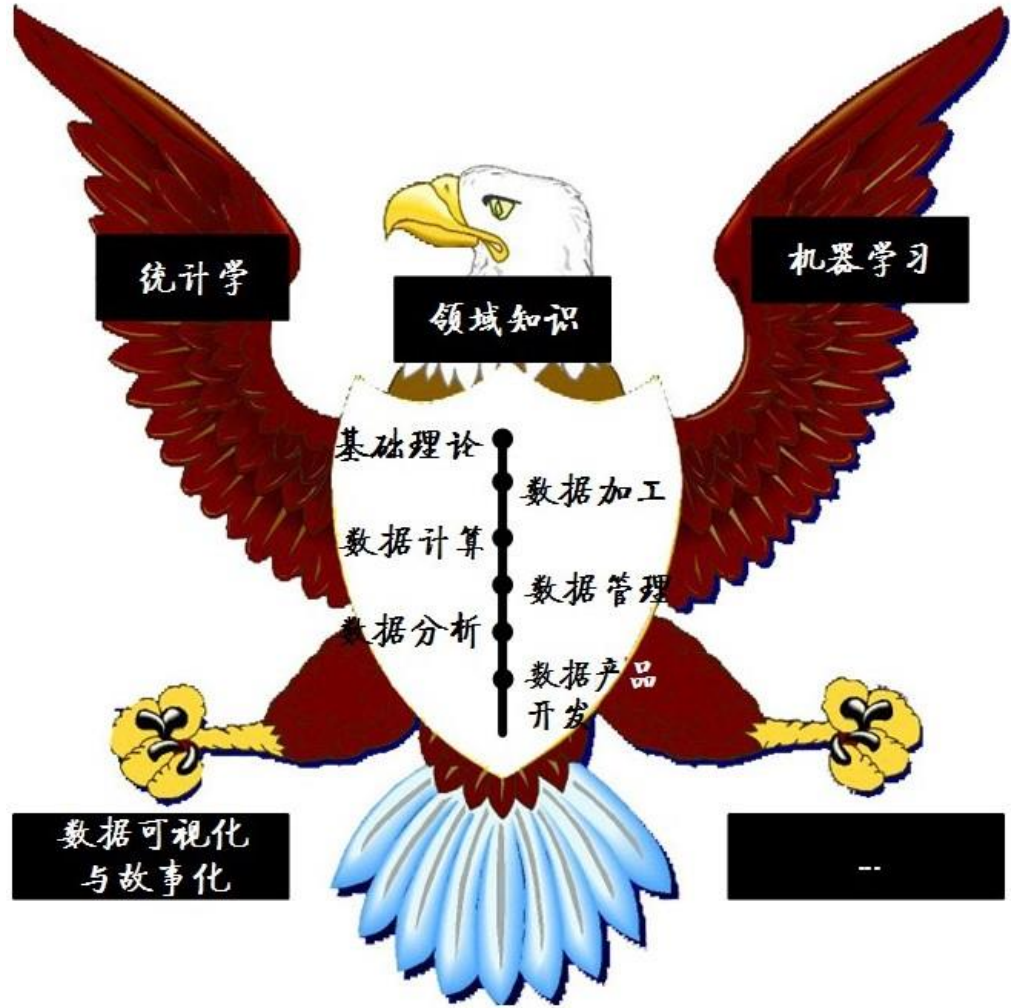
热衷挑战  
崇尚自由  
主张分享  
追求创新

## 黑客道德准则

The Hacker Ethic  
(Steven Levy, Hackers: Heroes of the Computer Revolution)

3个要素：理论+实战+素质

# 数据科学的知识体系



[来源] 朝乐门.数据科学[M].北京:清华大学出版社,2016.



# 数据科学的研究目的与任务

大数据及其运动规律的揭示

从数据到智慧的转化

数据洞察（Data Insights）

数据业务化

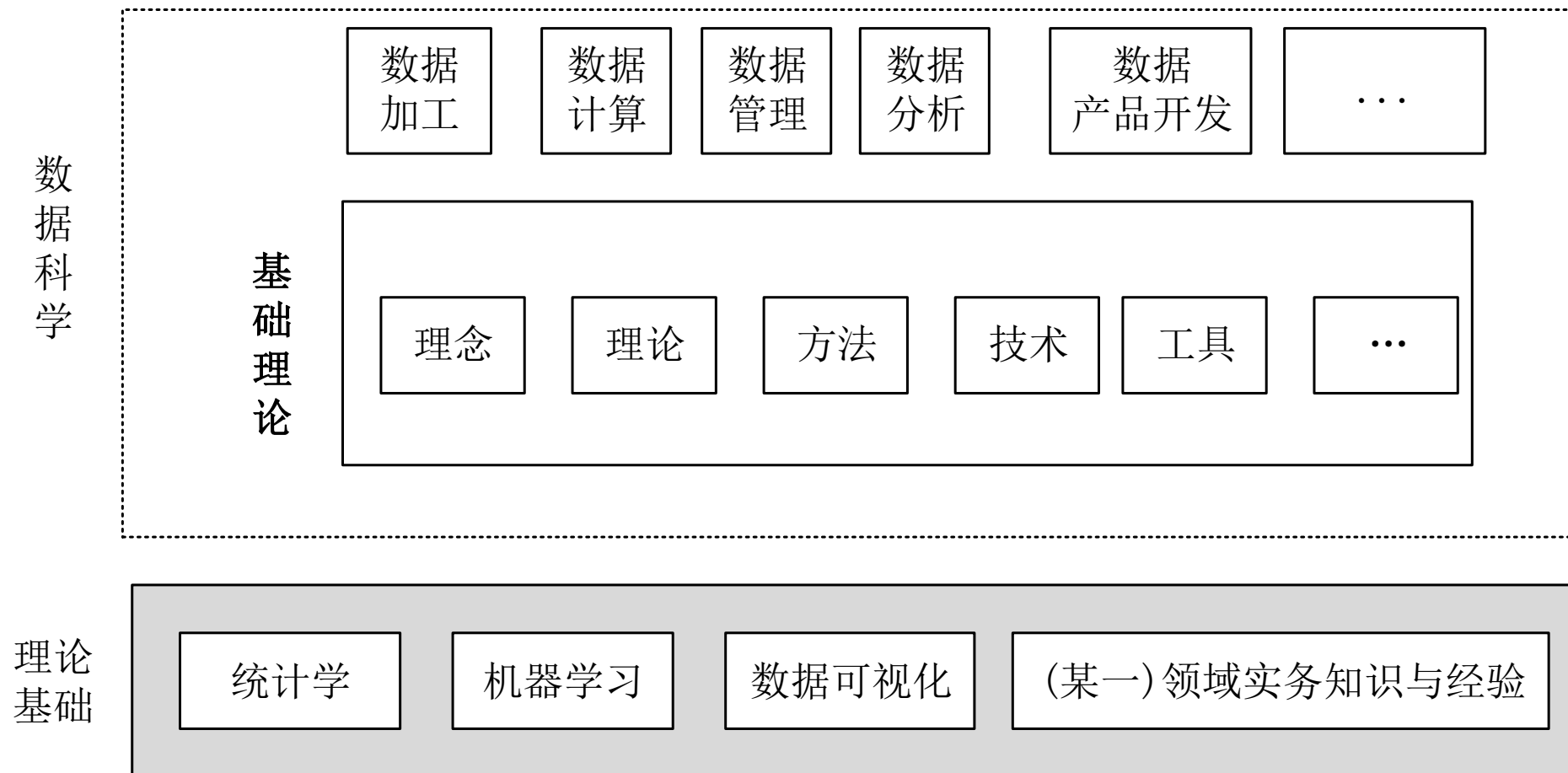
（大）数据分析

数据驱动型决策（支持）

数据产品的研发

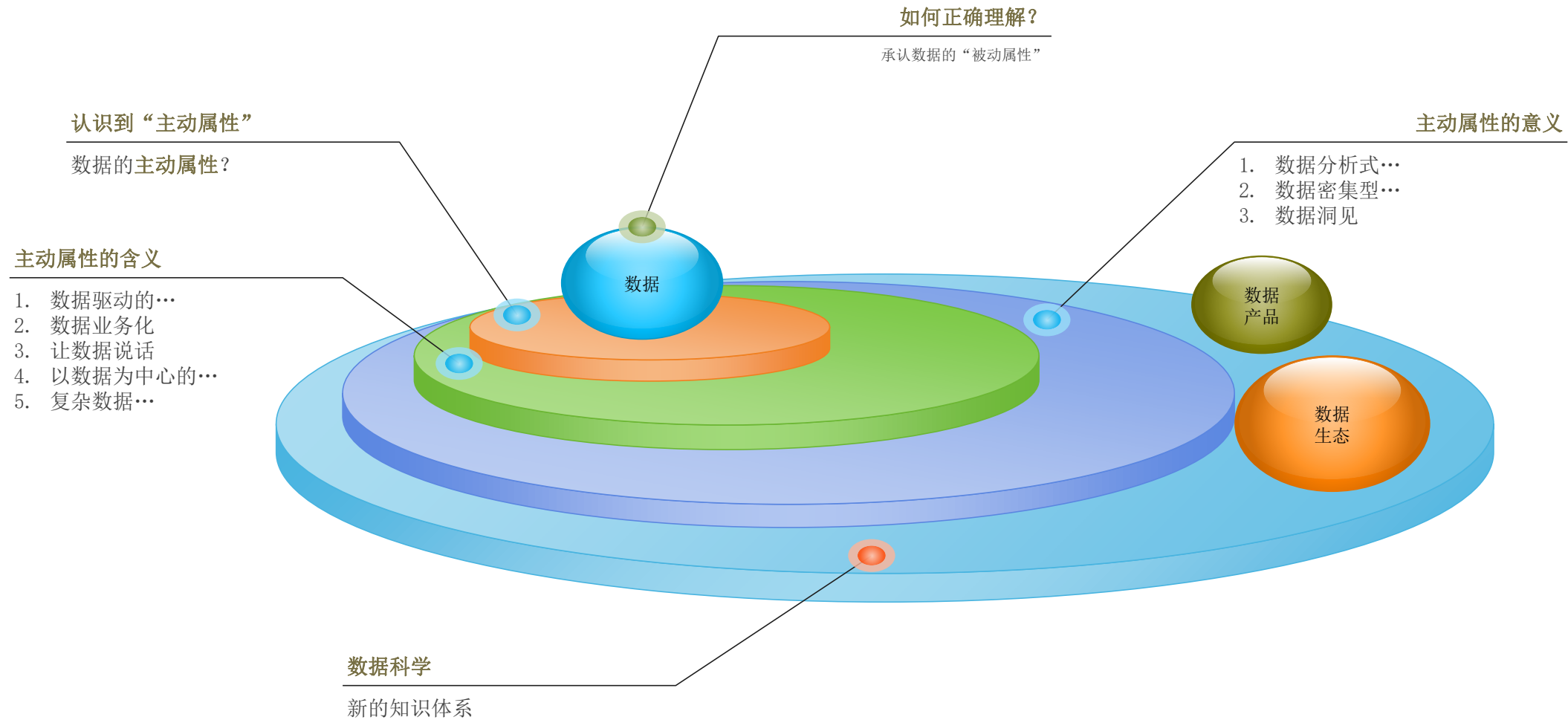
数据生态系统的建设

# 数据科学的研究内容



**【注意】基础理论与理论基础是两个不同的概念**

# 数据科学的新视角



# 数据科学的基本原则

三世界原则

三要素原则

数据驱动原则

数据复杂性  
原则

数据资产原则

DIKUW原则

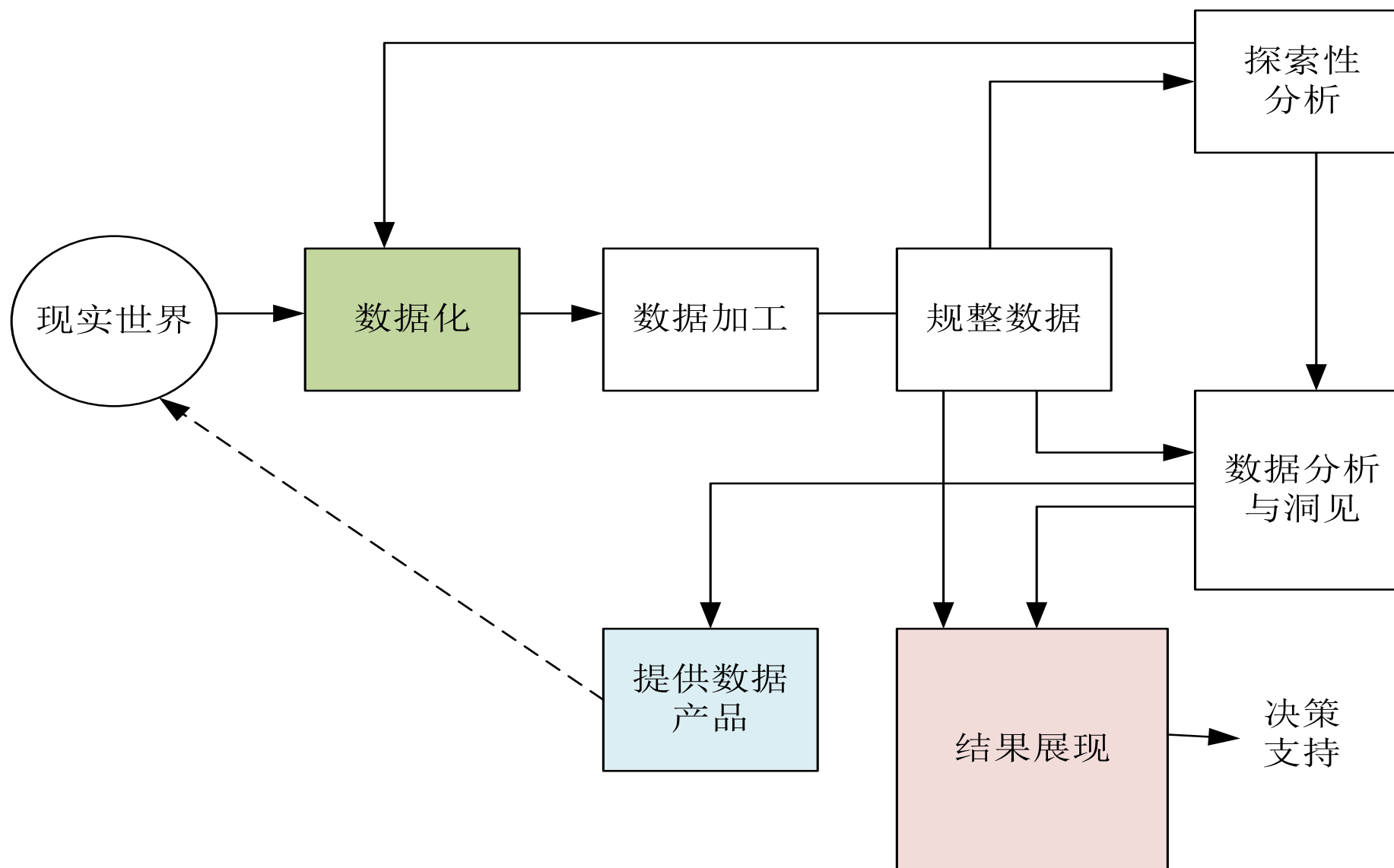
3C原则

协同原则

从简原则

数据范式原则

# 数据科学的基本流程



修改自：  
Schutt R,  
O'Neil C.  
Doing data  
science:  
Straight talk  
from the  
frontline[M].  
O'Reilly Media,  
Inc.2013:41.

# 《数据科学导论》建设与开源

## ——第二部分 如何建设

朝乐门



# 课程建设的重点与难点

1

- 坚持导论类课程的基本定位

2

- 凸显数据科学本身

3

- 培育兴趣与自学能力为主要任务

4

- 统筹《数据科学》课程链为基本前提

5

- 其他一些体会



## 2.1 坚持导论类课程的基本定位

### 数据科学专业

- 《数据科学与大数据技术导论》
- 《数据科学导论》
- 《数据科学基础》
- 《数据科学》

### 非数据科学专业

- 《数据科学方法与技术》
- 《大数据分析》
- 《大数据概论》



朝乐门等.数据科学与大数据技术专业特色课程研究[J].计算机科学,2018,45(3):1-8

## 为什么要开设这门课程？

“老知识”与“新数据”之间的矛盾

“数据”已经变了，“知识”还没有跟上

“系统性变革”与“局部性解读”的矛盾

大数据与AI对社会的影响是“整体性的”，但是很多就是“局部性”解读



# 坚持导论类课程的基本定位

经典理论  
×  
最佳实践

理论学习  
+  
动手操作

全集知识  
—  
领域差异性知识

最深奥理论  
÷  
最基本逻辑

学习“数据科学”的四则运算基本原则  
(来源: 朝乐门.数据科学[M].清华大学出版社,2016)

一图告诉你如何轻松学习《数据科学》



# 如何轻松学习数据科学

**1 掌握统计学、数学及机器学习**

- 数学: >可汗学院数学课程, >MIT OpenCourseware的线性代数
- 统计学: >Udacity的统计学, >OpenStax的统计学, >DataCamp的基于R的统计学概论
- 机器学习: >Stanford Online的机器学习, >Coursera的实用机器学习, >DataCamp的机器学习概论

**2 学会写代码**

- 学习计算机科学基础理论
- 学会端对端的开发
- 熟悉某一编程语言: >开源: R, Python等, >商业: SAS, SPSS等
- 互动中学习: >DataCamp, >Python-DataCamp

**3 理解数据库技术**

- 学会如何在数据库中存储和管理自己的数据
- 通过以下平台, 学习更多知识: MongoDB UNIVERSITY, Stanford ONLINE, DATASTAX, PostgreSQL, MySQL, ORACLE, cassandra

**4 探究数据科学流程**

采集 - 探索 - 加工 - 建模 - 验证 - 报告

**5 重视大数据**

- 大数据的3V特征: Variety, Volume, Velocity
- Hadoop框架: Hadoop, MapReduce
- Spark框架: Spark
- 了解大数据处理的特殊性
- 学会分布式存储与处理方法
- 理解内存集群计算框架的优点

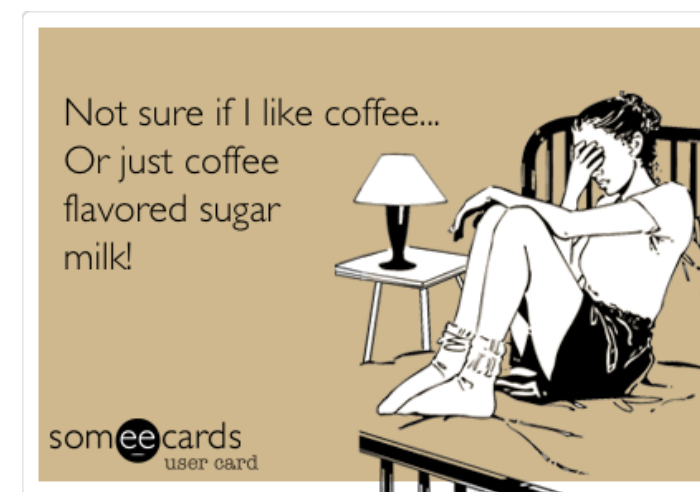
**6 成长、协作与学习**

- 积极参加相关竞赛: kaggle, DRIVEN DATA
- 要有自己的宠物项目: Meetup
- 与数据科学家合作: 挑战自己并拓宽自己的技能
- 与数据科学爱好者保持联系
- 培育数据科学家精神: 拥有自己的代表作, 提升讲故事的能力
- 理论联系实际

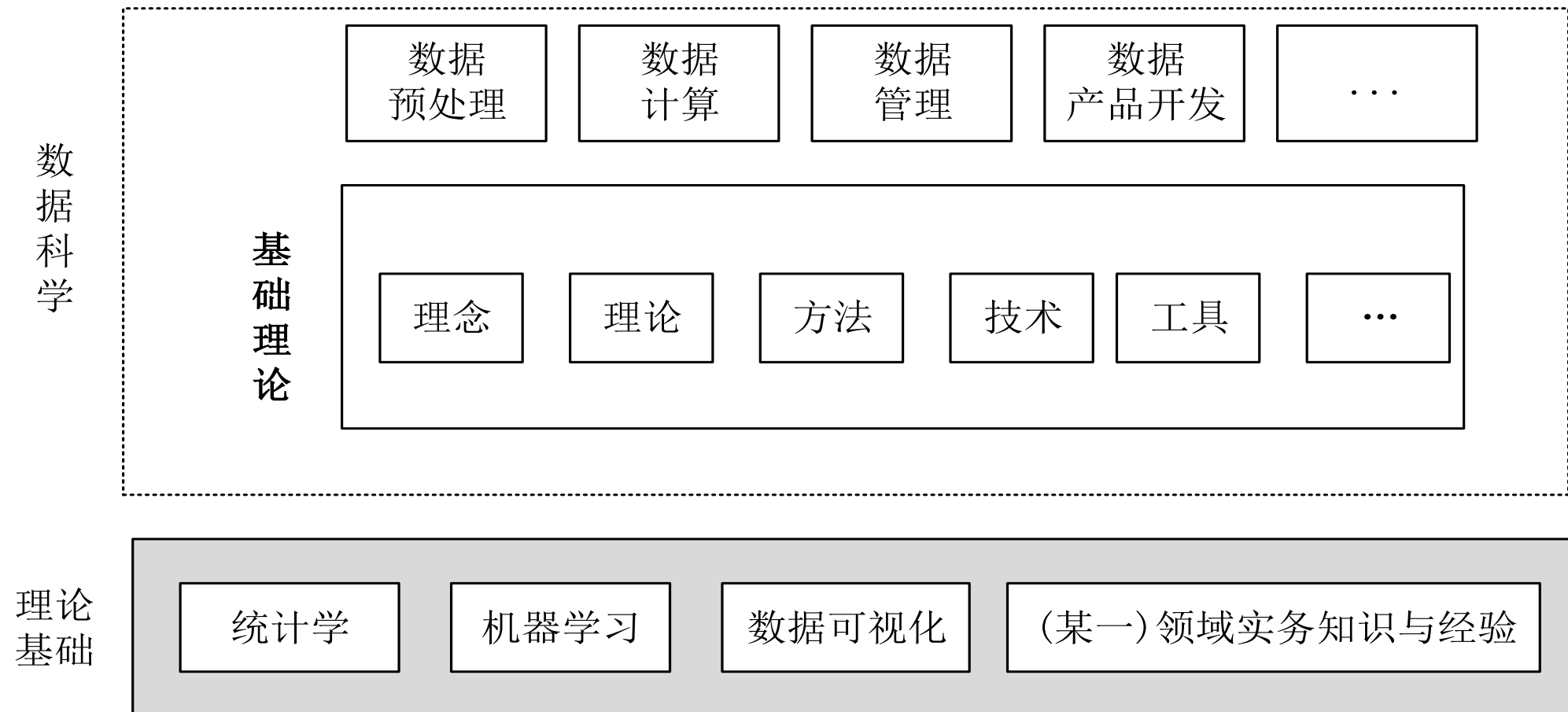
一图告诉你  
如何轻松学习《数据科学》



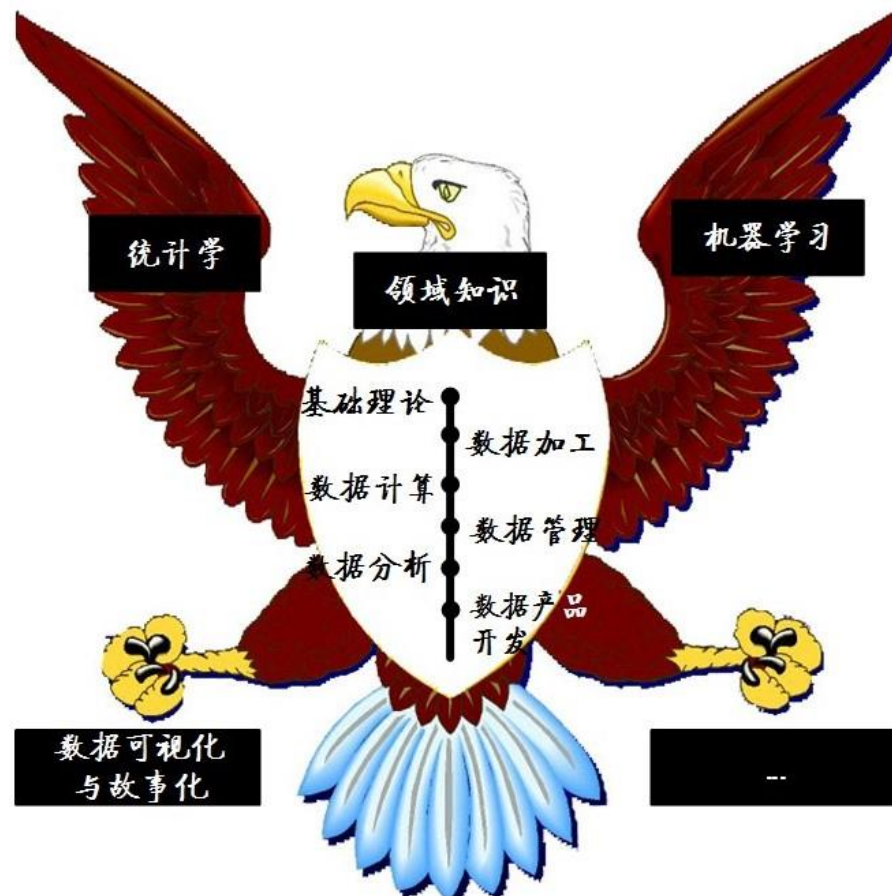
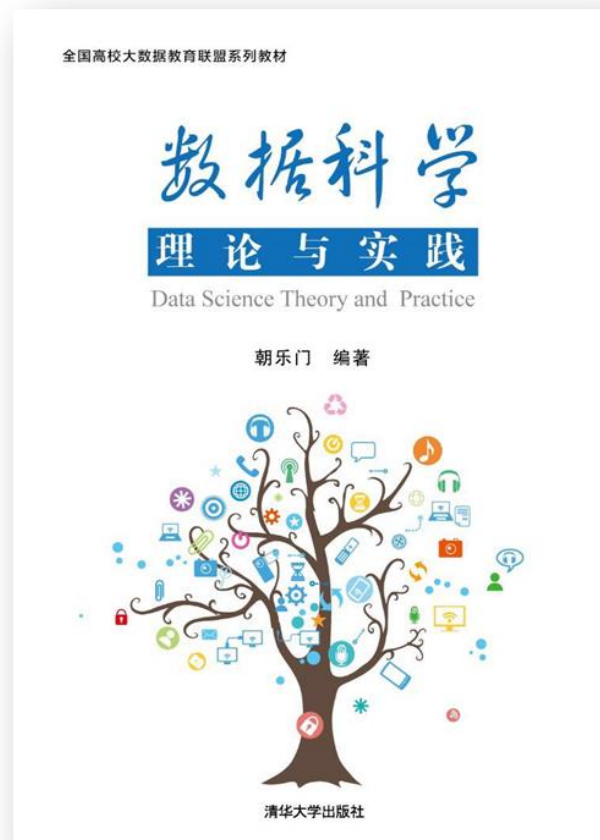
## 2.2 凸显数据科学本身



# 理论基础与基础理论的区别



# 凸显数据科学本身





# 数据科学导论课程中容易缺什么或错什么？

## 因果分析

- Causality (MIT, Stanford)

## 试验设计

- design of experiment (UW, Stanford)

## 数据产品开发

- Data products (UW, Harvard)

## 探索性分析与规范分析

- EDA(Columbia, UC Berkeley, Microsoft)

## 数据加工

- Wrangling, munging (Harvard)

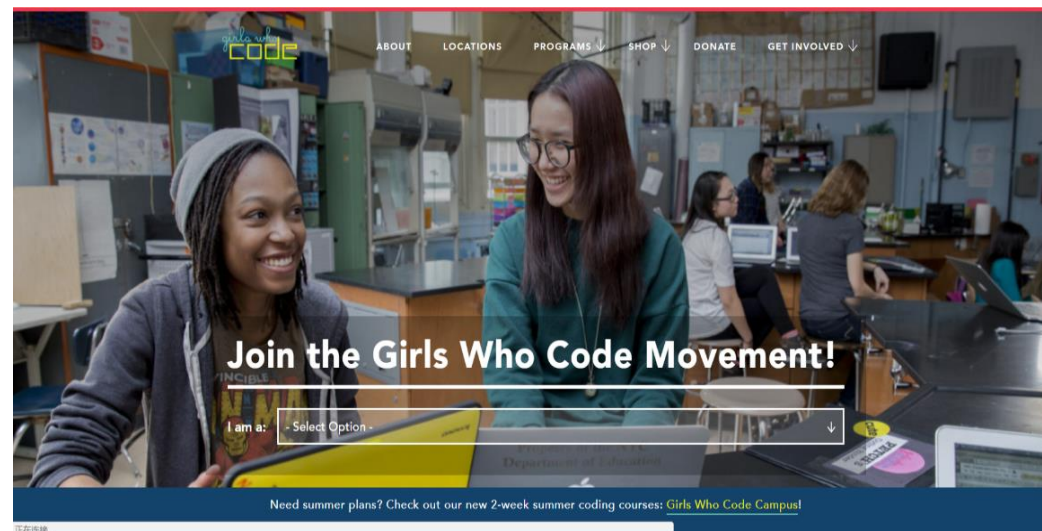
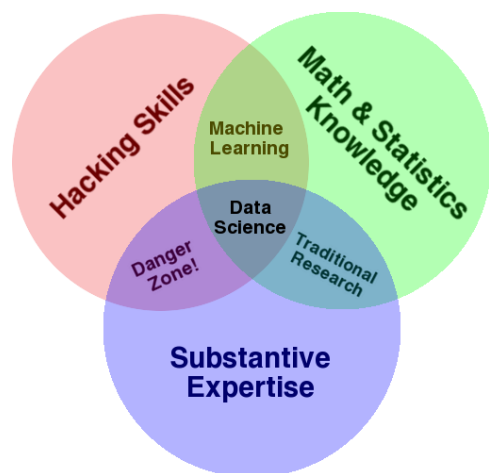
## 道德、伦理、法律等

- Ethics (Columbia)

基于数据的管理

## 2.3 培育兴趣与自学能力

- 数据科学的3个要素
- 数据科学的3C 原则



2018数据科学影响力报告  
|Top100个影响力人物、品牌  
与出版物 by Analytica

# Python数据科学实践系列

13.用Python讲解Spark基本原理

12.Twitter情感分析（共11篇文章，待发布）

11.2006-2017央视春晚主持人用词特征分析

10.数据可视化领域的6个著名实践及其源代码

9.基于Markov Chain Monte Carlo的智能手表睡眠数据分析

8.按主题抓取多个网页及建立自己的数据集

7.Jupyter Notebook/Lab中添加R Kernel的详细步骤

6.Web信息爬取 | 详解 + Reddit等2个案例实践

5.基于MovieLens的影评趋势分析|详解

4.Windows和PC机上搭建Spark+Python开发环境的详细步骤

3.Jupyter Notebook/Lab中添加R Kernel的详细步骤

2.盘点数据科学领域常用的Python库

1.如何用Python学习数据科学

Python数据科学实战系列



## 2.4 主要瓶颈及应对思路

### 缺少平台

- IBM Workbench
- 全国大数据教育联盟
- SPSS Analytic Server
- IBM Bluemix
- ...

### 缺少数据

- 开放数据
  - data.gov.uk
  - data.nasa.gov
  - fars.nhtsa.dot.gov  
(美国交通事故数据)
- 挑战赛数据
  - Kaggle
  - github.com
- Capstone 类项目
- 与企业合作
- ...

### 缺少经验

- 科研项目
- 挑战赛
- 开源项目
- 同行交流
- 跨行合作
- ...

### 缺少参考材料

- 国外教材
- 本土教材
- 协同与共享



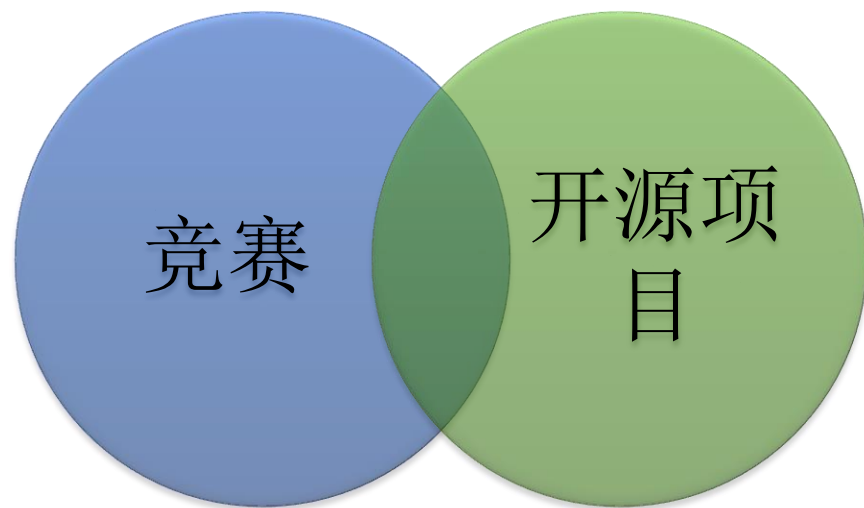
# 开源课程

Open-source course

## 2.5 重视专业差异性



## 2.6 注意两条捷径



国内外数据科学领域的重要竞赛





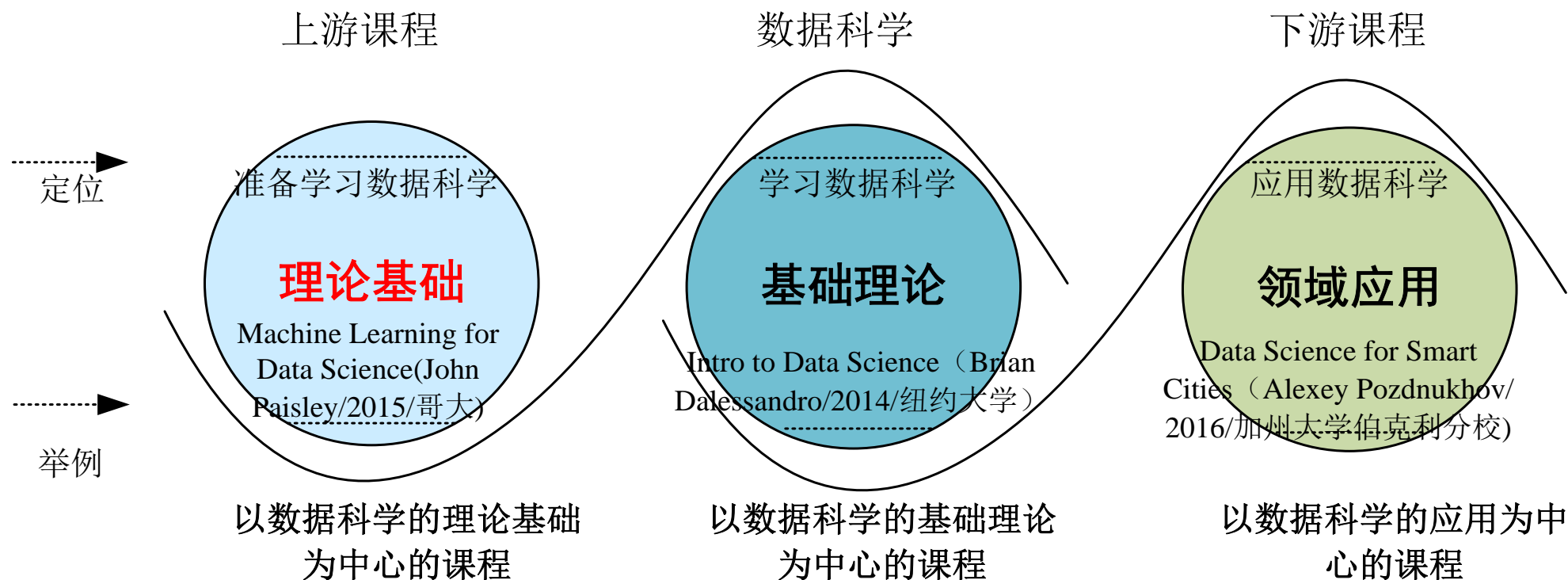
# 重视学生的就业能力：行业动态及面试题的深度分析



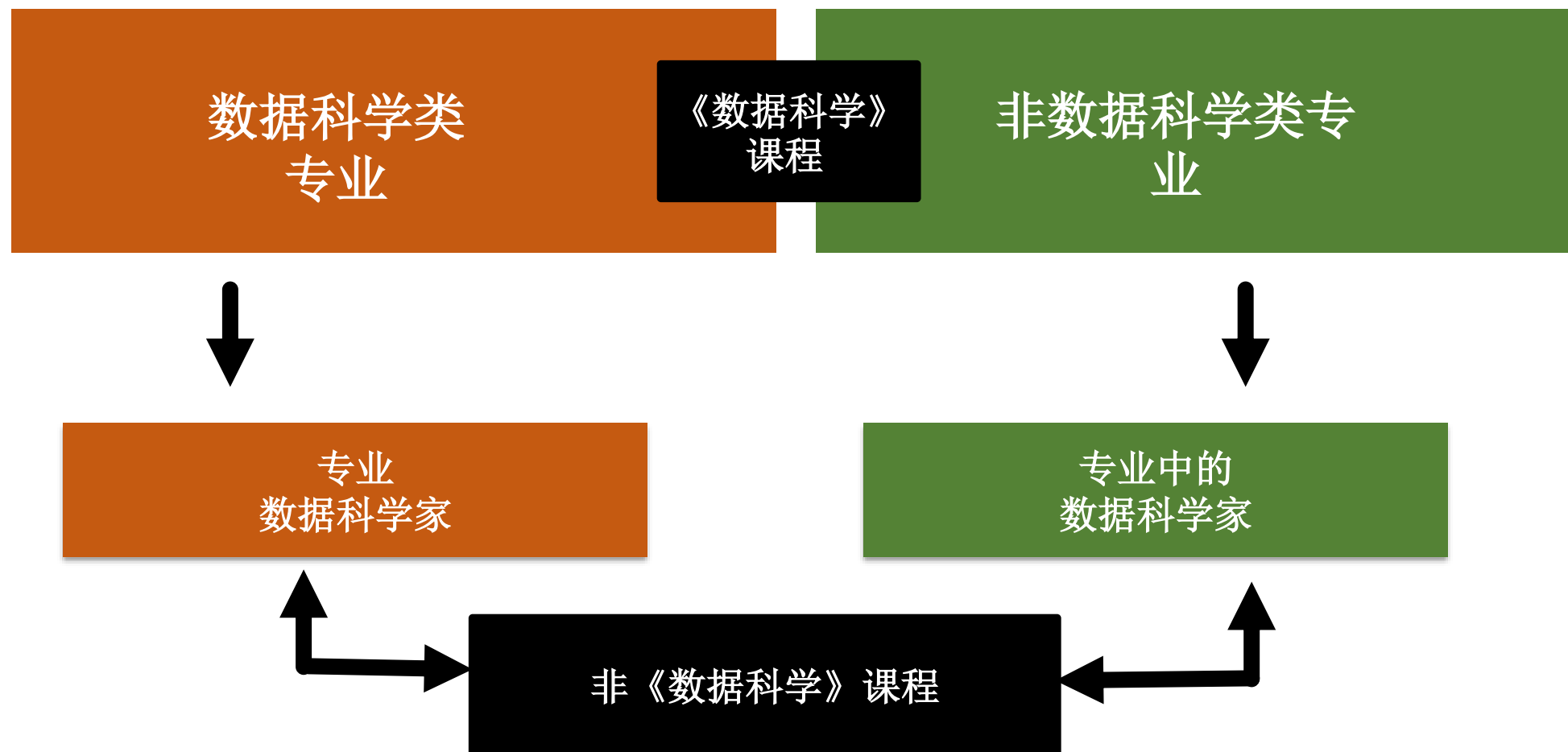
数据科学相关岗位面试题库  
500题



## 2.7 统筹课程链



## 2.8 课程建设与专业建设之间的协同



# 小结

2.1 坚持导论类课程的性质

2.2 凸显数据科学本身

2.3 培育兴趣与自学能力

2.4 主要瓶颈及应对思路

2.5 重视专业差异性

2.6 注意两条捷径

2.7 统筹课程链

2.8 课程建设与专业建设之间的协同

# 《数据科学导论》建设与开源

## —— 第三部分 开源倡议

朝乐门

## 3.1 为什么要提出【开源课程】的倡议

### 教师的困境

- 一些低级重复性劳动耗尽了老师们的绝大部分**备课时间**
- 技术和环境变化太快，**每年需要更新数据**
- **有疑问**，不知道请教谁；就算知道应该找“他”，但也请不动“他”
- **有好的经验与做法**，做了很好的课程建设，也没有机会展示给同行
- **有没有Bugs**，自己也不知道
- 花费几年的时间，备好一门课，结果**上一两次就被叫停**

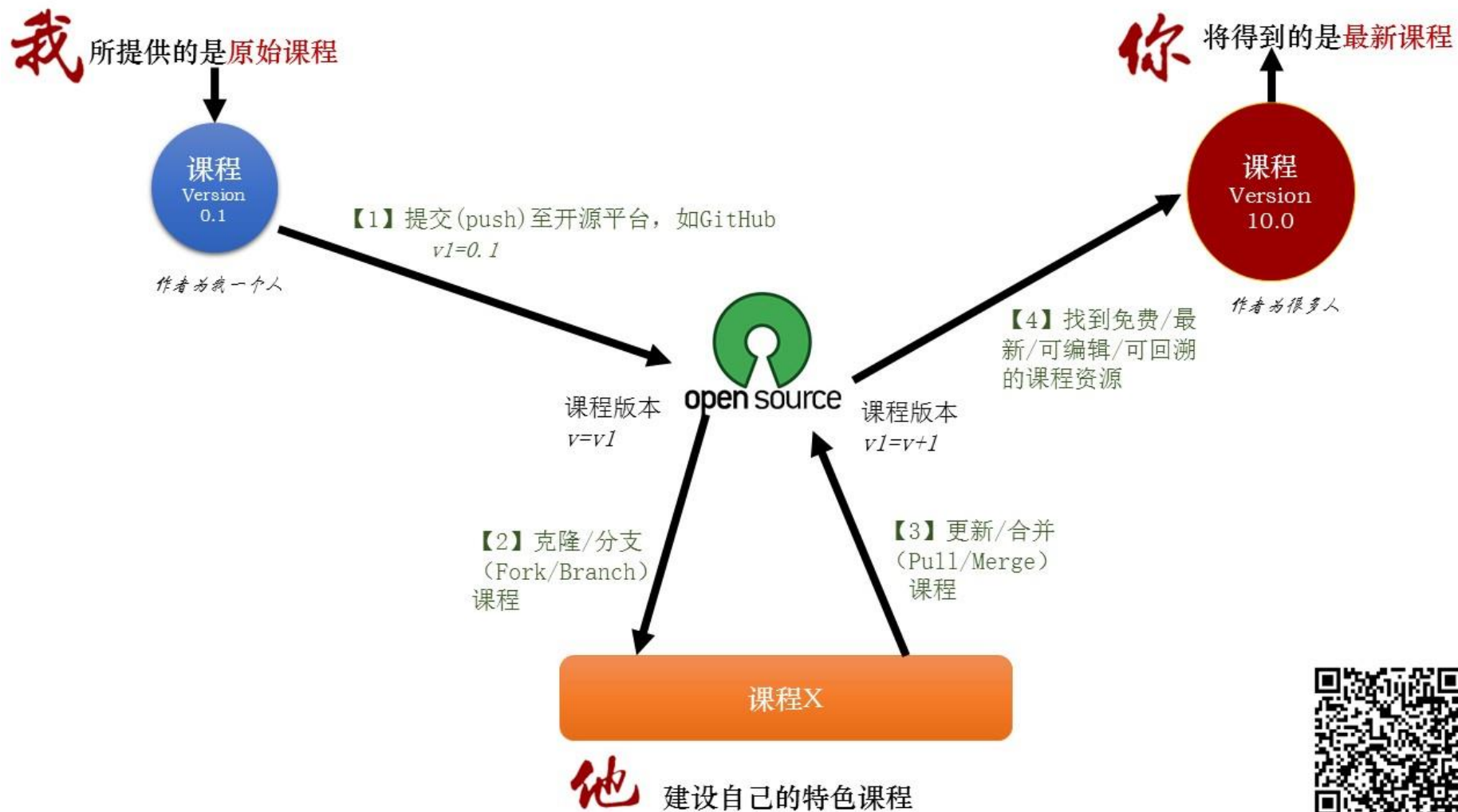
### 教育部门的痛点

- **所谓的优质资源和精品课程都是针对学生的**，对于学生来说确实是精品课程
- 其他老师们只能看到这些示范课程的**“可执行程序级的最终结果”**，根本拿不到真正有价值的东西——**“源代码级的可编辑备课材料”**
- 目前，教改立项惠及的**仅仅是少数（优秀）老师**
- 缺乏一种**确保绝大多数老师的备课质量的保障机制**，老师们都在“作坊式备课”，备课质量参差不齐。



What

# 3.2 什么是【开源课程】



Mailing lists  
Issues  
wiki



# 什么是【开源课程】

## 是什么

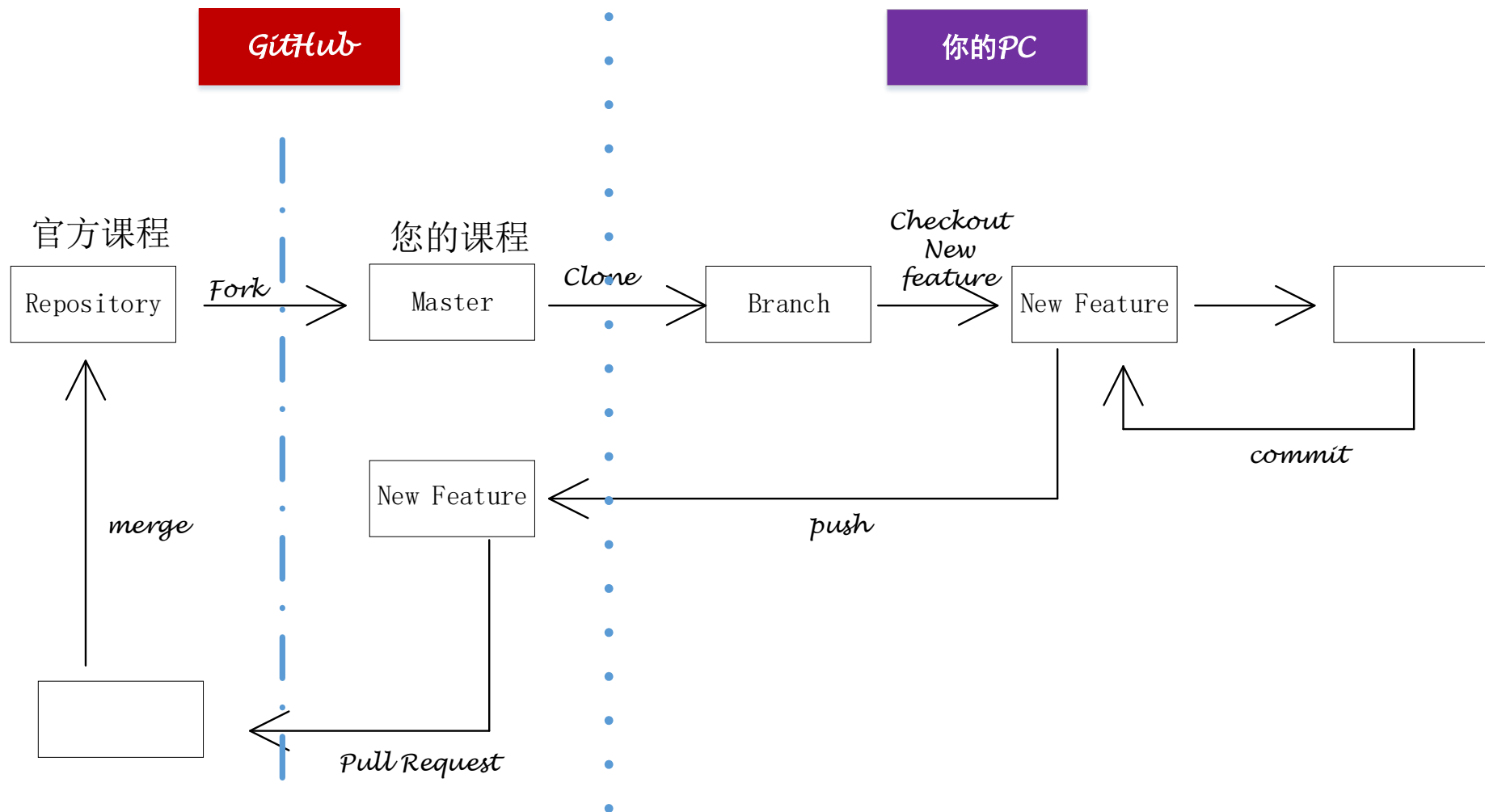
- 是一种文化
- 是一种社区
- 是一种生态环境

## 不是什么

Open-source  $\neq$  source-available  
不仅是资源共享，更是合作共建模式  
不仅是资源库，更是社区  
不仅是PPT，而是课程全套资源  
不是PDF，而是源代码



### 3.3 如何建设【开源课程】



## 3.4 谁可以参与【开源课程】？

### 无歧视

- 每位老师都是平等的
- 不再是少数人的“特权”

### 有受益

- 共同维护
- 超越限制
- 展示自己
- 结识同行
- 帮助他人的同时，也在帮助自己

*Where*

## 3.5 【开源课程】在哪里？

### 开源社区

- GitHub、Git ,...
- 【例】<https://github.com/LemenChao/Introduction-to-Data-Science>

### 邮件列表

- Google groups,...
- 【例】DataScienceChina

### 行业联盟

- 全国高校大数据教育联盟

# 《数据科学导论》的开源

## 对外公布

- 2017年11月18日
- 微信公众号“数据科学DataScience”

## 正式启动

- 2017年12月23日
- 全国高校大数据教育联盟“数据科学与大数据技术专业核心课程建设系列研讨会”

## 第一个示范课程

- 《数据科学导论》：朝乐门，chaolemen@ruc.edu.cn
- 订阅订方法：发送邮件至 [datasciencechina+subscribe@googlegroups.com](mailto:datasciencechina+subscribe@googlegroups.com)
- 发帖方法：发送邮件至 [datasciencechina@googlegroups.com](mailto:datasciencechina@googlegroups.com)
- 退订方法：发送邮件至 [datasciencechina+unsubscribe@googlegroups.com](mailto:datasciencechina+unsubscribe@googlegroups.com)





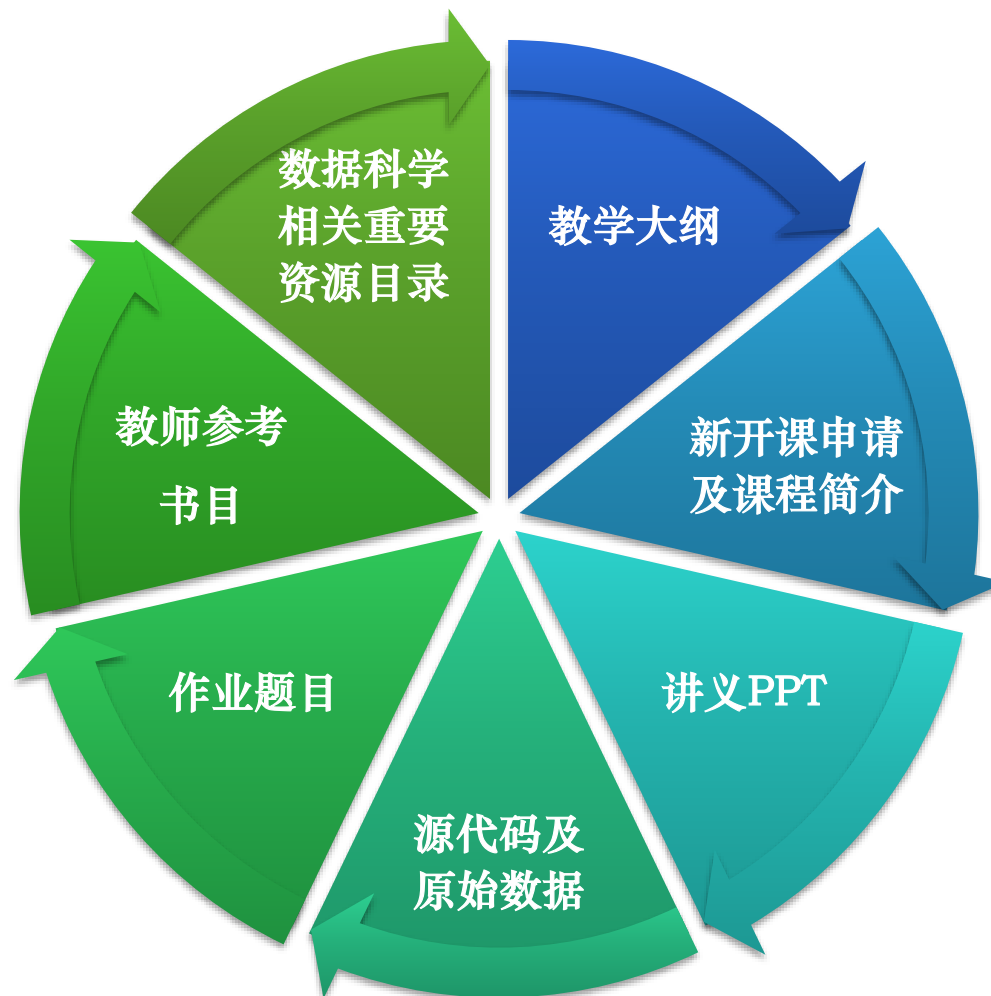
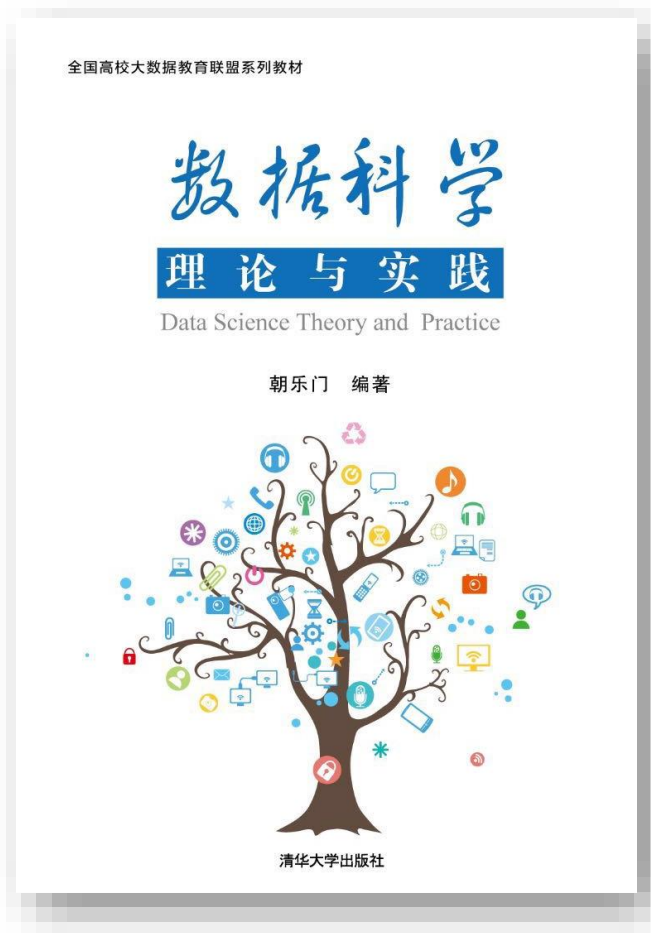
# 《数据科学导论》的开源



Open Sourcing 高校课程资源的倡议及《数据科学导论》的开源



# 《数据科学导论》的全套教学资料



## 更多参考资料（1）

2015-2018大数据与数据科学领域68篇推荐文章



未来数据科学家必备的【核心算法】与【常用模型】



Python  
数据科学实战系列



## 更多参考资源 (2)

数据科学相关岗位面试题库  
500题



数据科学领域常用的Python  
库

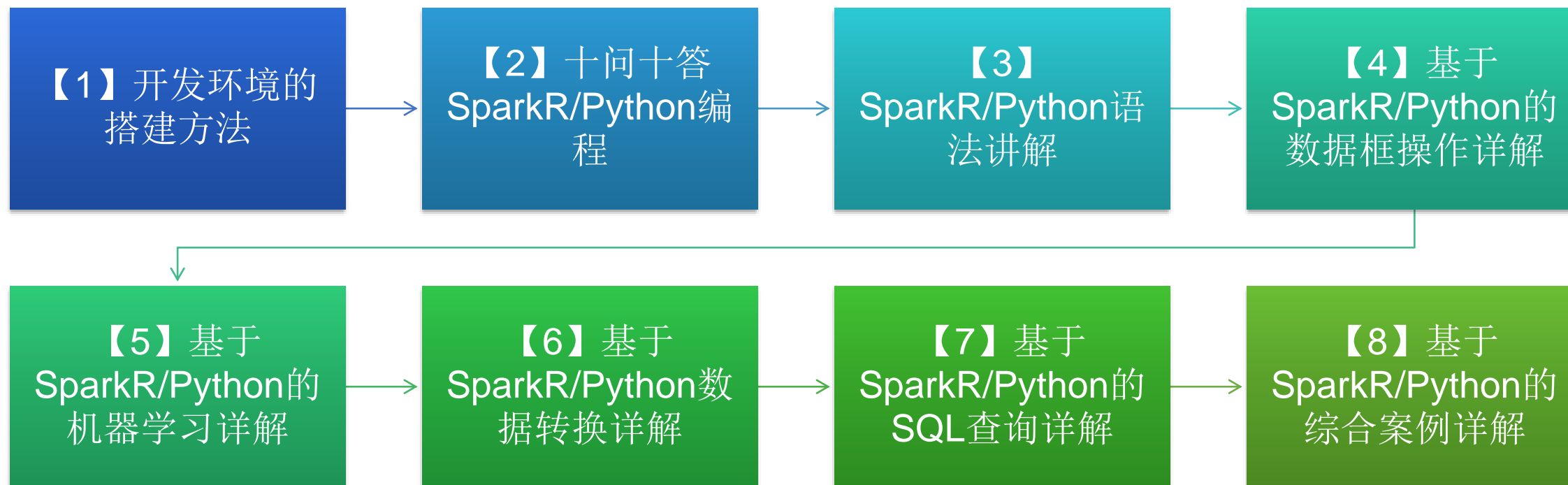


TOP20 | Python人工智能  
与机器学习开源项目





# Spark+R/Python编程系列教程





参考书目



教学支撑平台



主讲人联系方式



主讲人微信