



欢迎参加

第 1 届全国高校大数据教学研讨会

2017年5月12日-13日 厦门大学



第1届全国高校大数据教学研讨会 (BDTS2017)
大会特邀报告



2017年5月12日至13日，第1届全国高校大数据教学研讨会（BDTS2017）在厦门大学科艺中心音乐厅隆重举行。本届研讨会由教育部高等学校计算机类专业教育指导委员会主办，厦门大学、厦门理工学院、贵州师范大学、人民邮电出版社联合承办，旨在搭建专业的大数据教学交流平台，汇聚全国高校大数据教学精英力量，共同探讨大数据专业和课程体系建设，为加快推进全国高校大数据教学发展贡献力量。来自全国300多所院校的400余名教师参加了本次研讨会。

厦门大学谭绍滨校长助理、人民邮电出版社教育中心营销部肖稳副主任，北京大学、中国科学院、厦门大学、华东师范大学、同济大学等重点院校的6位大数据教学知名专家，以及来自国内知名大数据企业的3名业界专家出席会议并做特邀大会报告。厦门大学林子雨助理教授主持会议。

更多内容请访问大会官网：<http://dbl原因lab.xmu.edu.cn/post/bigdata2017/>



同济大学 王伟 副教授 在做大会特邀报告

如何建设大数据与数据科学 通识实践类课程？

王伟

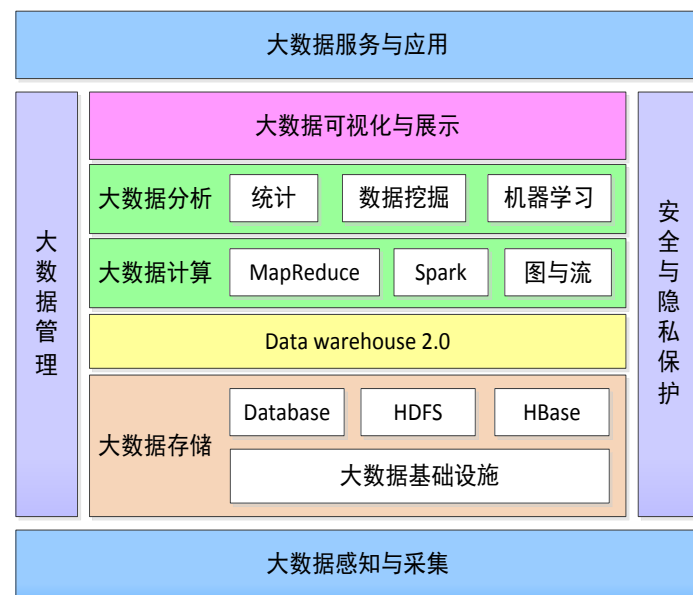
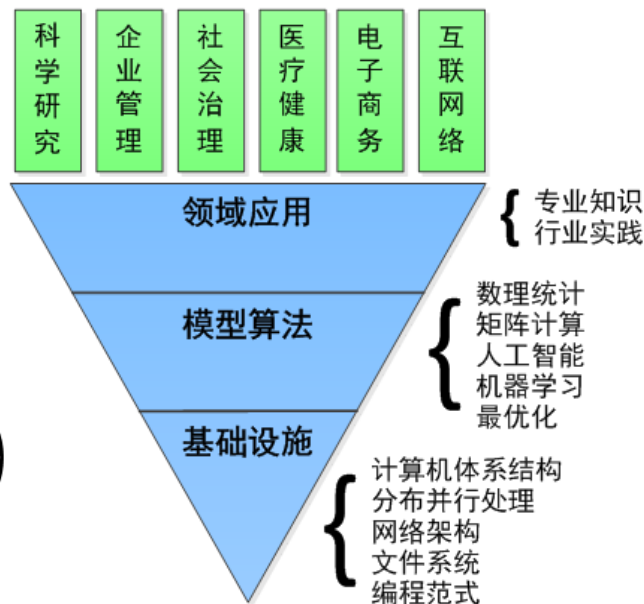
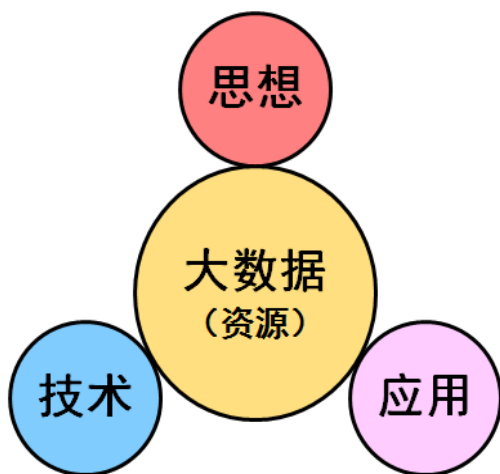
同济大学，计算机科学与技术

2017年5月于厦门

Outline

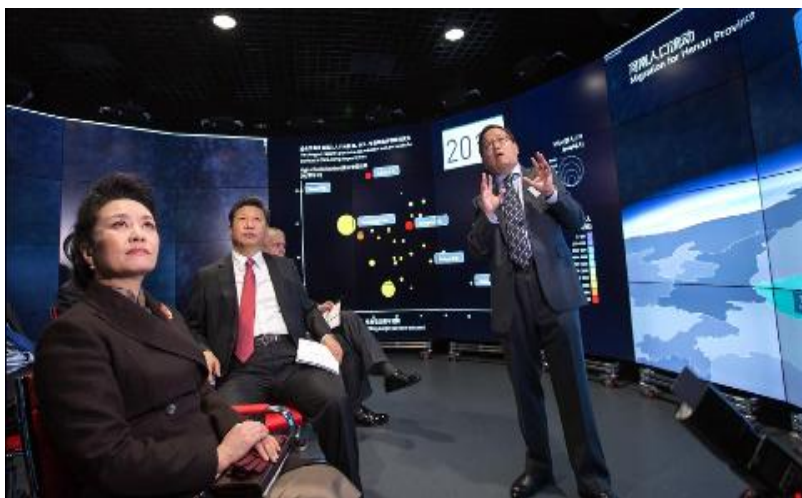
- 对大数据与数据科学的思考
- 导论类课程建设与教学实践
- 数据科学与大数据实践平台
- 总结：未来的挑战

大数据的冲击



时代的呼唤

国家



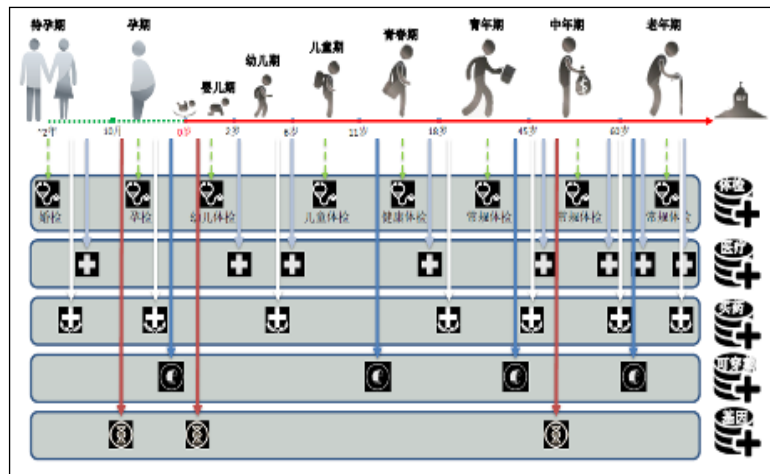
企业



机器



人类



教育的变革



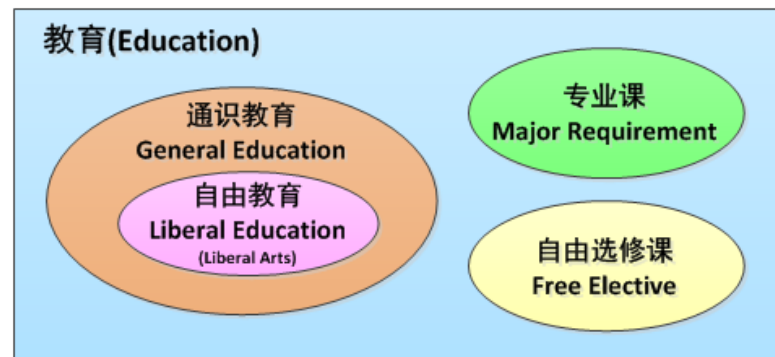
大数据教育改革

- 2016年大数据与数据科学的教育改革
 - 《数据科学与大数据技术》本科专业（专业代码：080910T）
 - 《大数据技术与应用》高职专业（专业代码：610215）
- 2017年3月，教育部公布第二批“大数据专业”获批高校，两批共35所
 - 第一批：北京大学、对外经济贸易大学、中南大学3所
 - 第二批：中国人民大学、复旦大学等32所
- 全国高校大数据教育联盟、数据中国“百校工程”等等

通识教育的属性

- 通识教育，既是为常识服务的教育。
- 通识教育的精神支柱和思想来源： Liberal arts（自由博雅）
 - 旨在培养一流的头脑、一流的心灵；
 - 一种间接的思维训练，传授的是“软实力”；
 - 关注审辨式思维、交流能力和解决问题的能力。

中国通识教育的意义：为了解决中国社会的实际问题，发展具有中国特色的通识教育。



数据科学作为一门通识课

- 数据科学有利于培养信息时代一个健全的人；
- 数据科学有利于培养跨学科视野；
- 数据科学有利于培养表达自我所必备的技能；
- 数据科学有利于培养个人的科学思维方式；
- 数据科学有利于围绕数据开展实践。

A liberal arts background gives data science its soul.



Liberal Arts理念向数据科学注入“博雅”之心

- 问正确问题的能力
- 科学方法观
- 团队协作的精神
- 沟通交流的能力
- 三观正确的决策



Ultimately lead to a better understanding of the natural world.

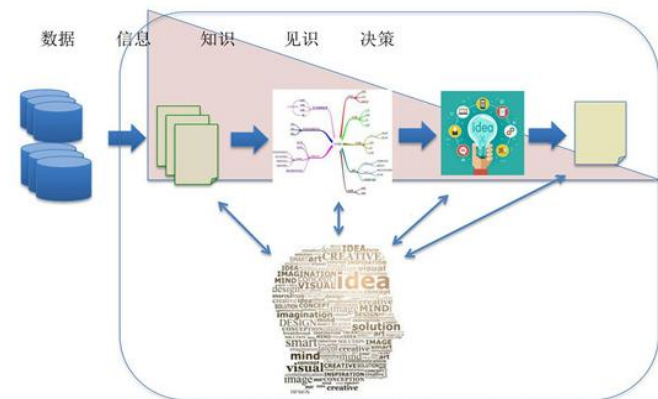
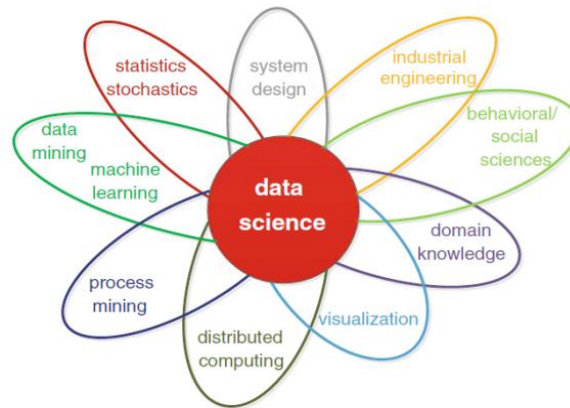
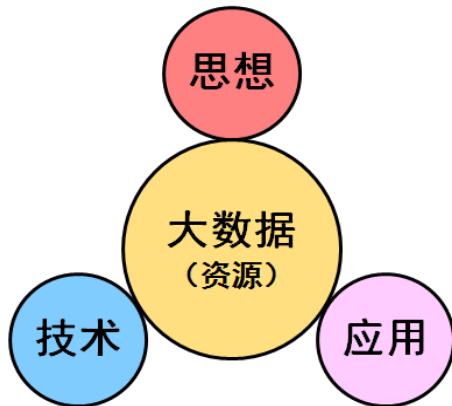
思考2： 数据科学的知识体系



数据科学：从大数据到行动

数据科学 ≈ 思维 + 计算机科学 + 统计 + 应用

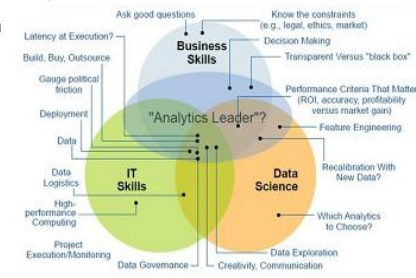
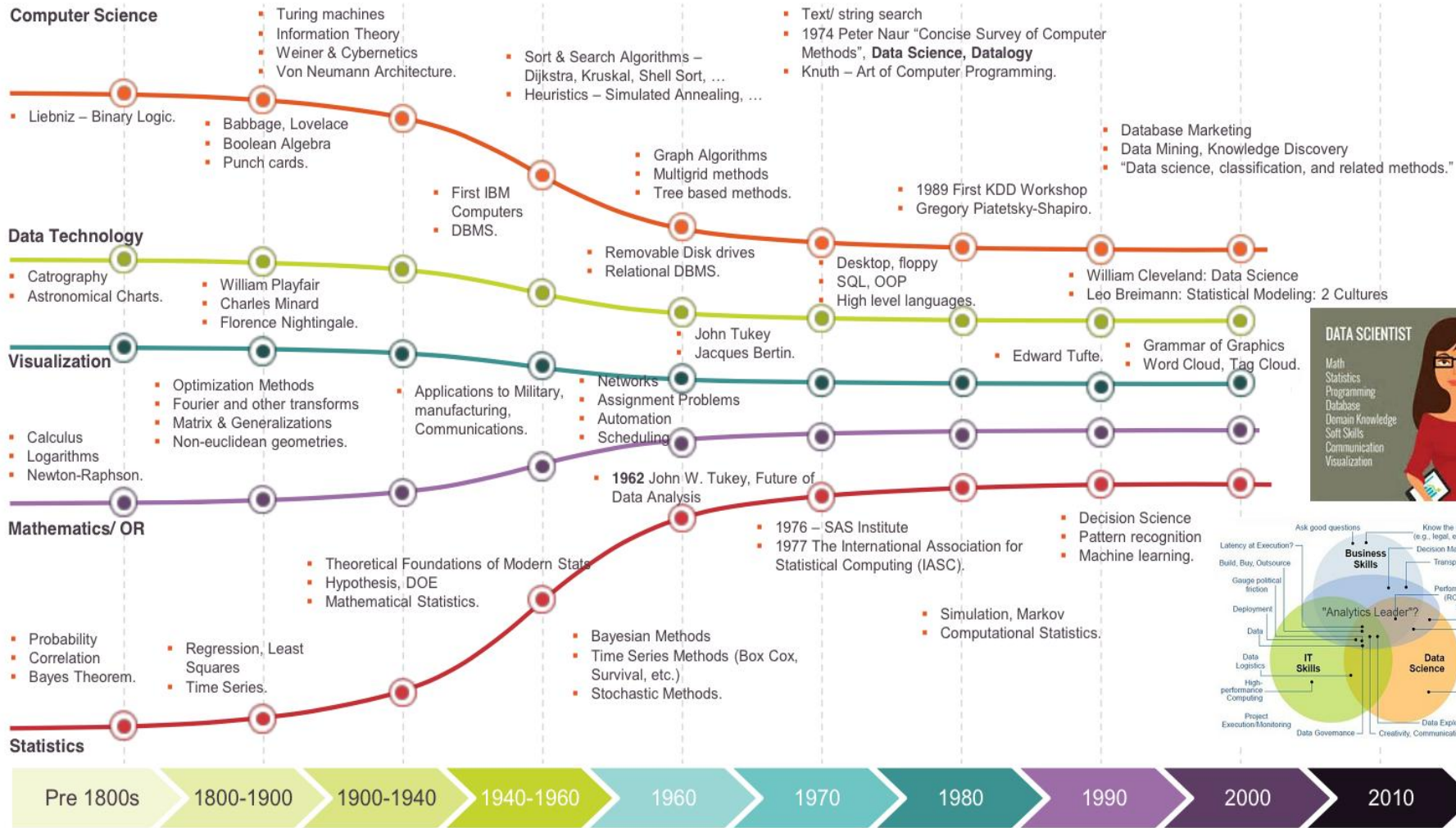
- 首先，建立数据思维方式，学习怎样利用数据；
- 其次，应该了解数据清理、集成、探索等相关技术；
- 最后，洞见和商业意识也至关重要。



数据科学的三大学科支柱

- 数据科学天生就是一个**交叉学科**，和数据科学最为密切的一些学科包括：计算机科学与技术、数学、统计学、信息管理、情报学等。数据科学的三大主要支柱为：
 - **Datalogy (数据学)**: 对应数据管理 (Data management)
 - **Analytics (分析学)**: 对应统计方法 (Statistical method)
 - **Algorithmics (算法学)**: 对应算法方法 (Algorithmic method)

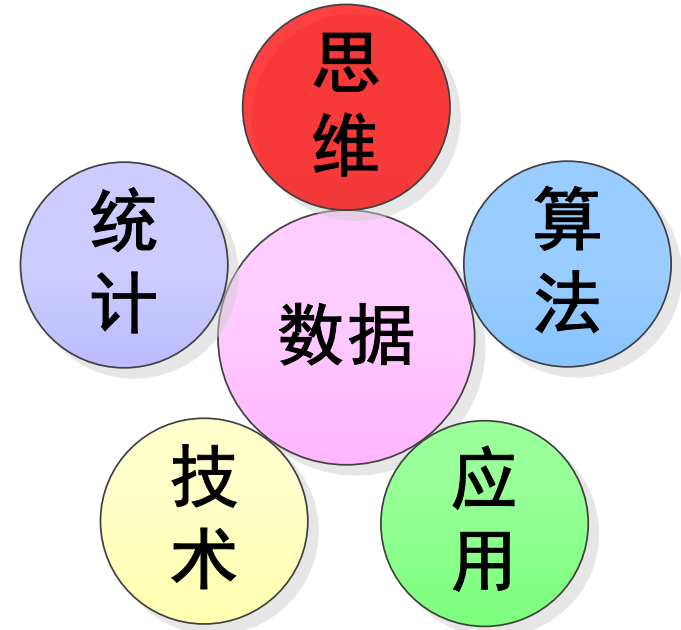
数据科学发展历程



数据科学的五大要素

- A-SATA模型

- 分析思维 (Analytical Thinking)
- 统计模型 (Statistical Model)
- 算法计算 (Algorithmic Computing)
- 数据技术 (Data Technology)
- 综合应用 (Application)



数据科学的核心知识点

- **分析思维** (Analytical thinking): 包括**计算思维** (Computational thinking)和**统计思维** (Statistical thinking);
- **数学基础**: 微积分、线性代数、概率统计、离散数学等;
- **数据建模与评估**: 统计模型、实验设计、模型评估等;
- **算法实现**: 问题求解能力和算法设计;
- **数据管理**: 设计数据的整个生命周期, 包括感知、存储、计算、分析、可视化等;
- **知识转化**: 沟通交流, 道德规范等。

思考3：数据科学的实践方法



信息化、大数据与教育

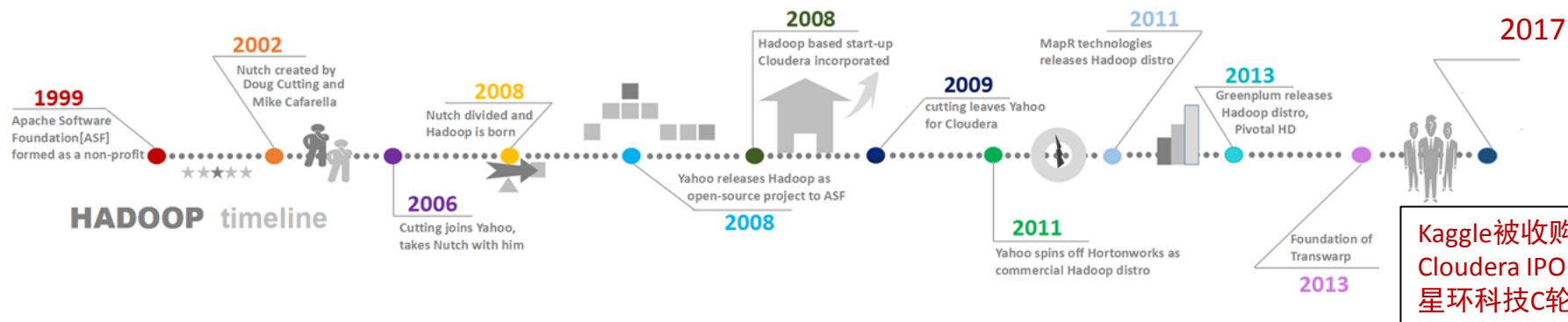


习近平：

“当今世界，科技进步日新月异，**互联网、云计算、大数据**等现代信息技术深刻改变着人类的思维、生产、生活、学习方式，深刻展示了世界发展的前景。”

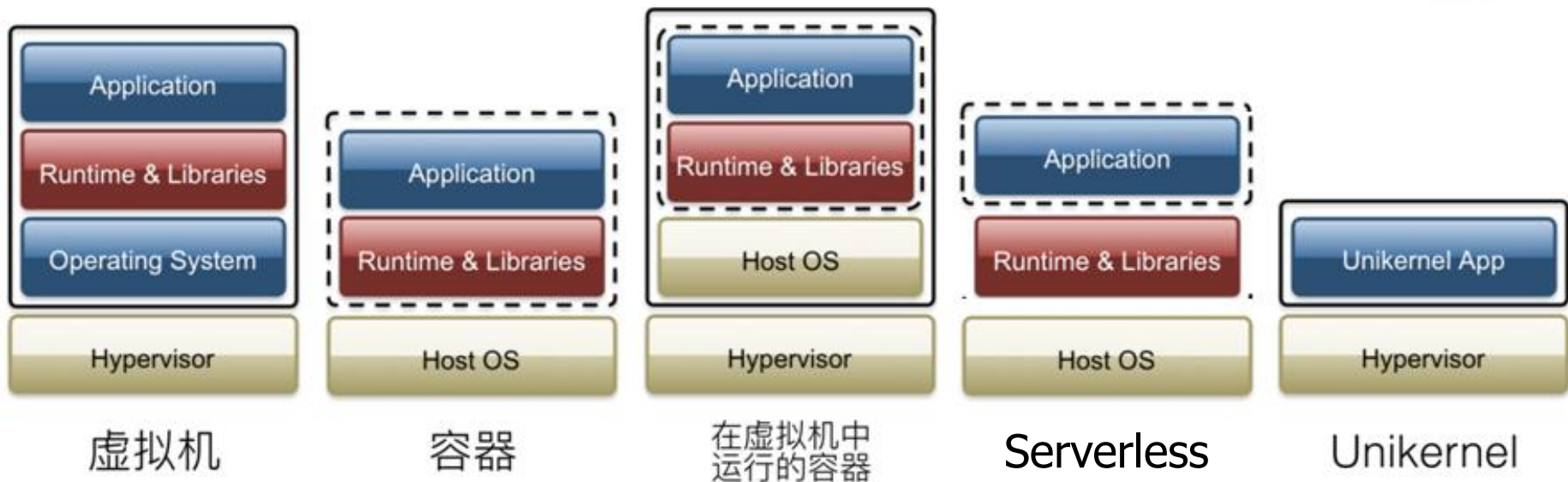
“因应信息技术的发展，推动教育变革和创新，构建**网络化、数字化、个性化、终身化**的教育体系，建设“**人人皆学、处处能学、时时可学**”的学习型社会，培养大批创新人才，是人类共同面临的重大课题。”

大数据技术的演变



ETL 数据装载工具	Workflow 工作流开发工具	数据质量 管理工具	可视化 报表工具	机器学习 建模工具	统计挖掘 开发工具	资源 管理工具	分析管理工具
SQL批处理 Batch Processing	交互式分析 OLAP Analysis	实时数据库 OLTP Transactional Processing	数据挖掘 机器学习 算法库 / 框架 Machine Learning	深度学习 Deep Learning	图分析引擎 Graph Analysis	流处理引擎 Streaming Processing	
批处理框架 Map/Reduce2, Tez		高性能处理框架 Spark		向量处理框架 TensorFlow			通用计算引擎
短时任务资源管理框架 YARN		长时任务资源管理框架 Mesos		资源隔离 / 调度 / 管理框架 Kubernetes			资源管理框架
分布式文件系统 HDFS	分布式大表 HBase	搜索引擎 Elastic Search	分布式缓存 Redis	消息队列 Kafka	分布式协作服务 Zookeeper		分布式存储引擎

云计算技术的演变



数据科学与大数据实践平台

- 围绕**大数据和数据科学**提供**全在线**课程教学、习题练习、动手实践、视化案例库、数据竞赛等服务。
- 通过**校企合作**，对接**广泛**的大数据基础软件；
- 通过自主研发的**云件技术**，在线学习使用**任意**数据分析工具。
- 形成一个**学、做、测、竞、评、奖**的大数据在线学习**闭环**。

导论类课程建设与教学实践

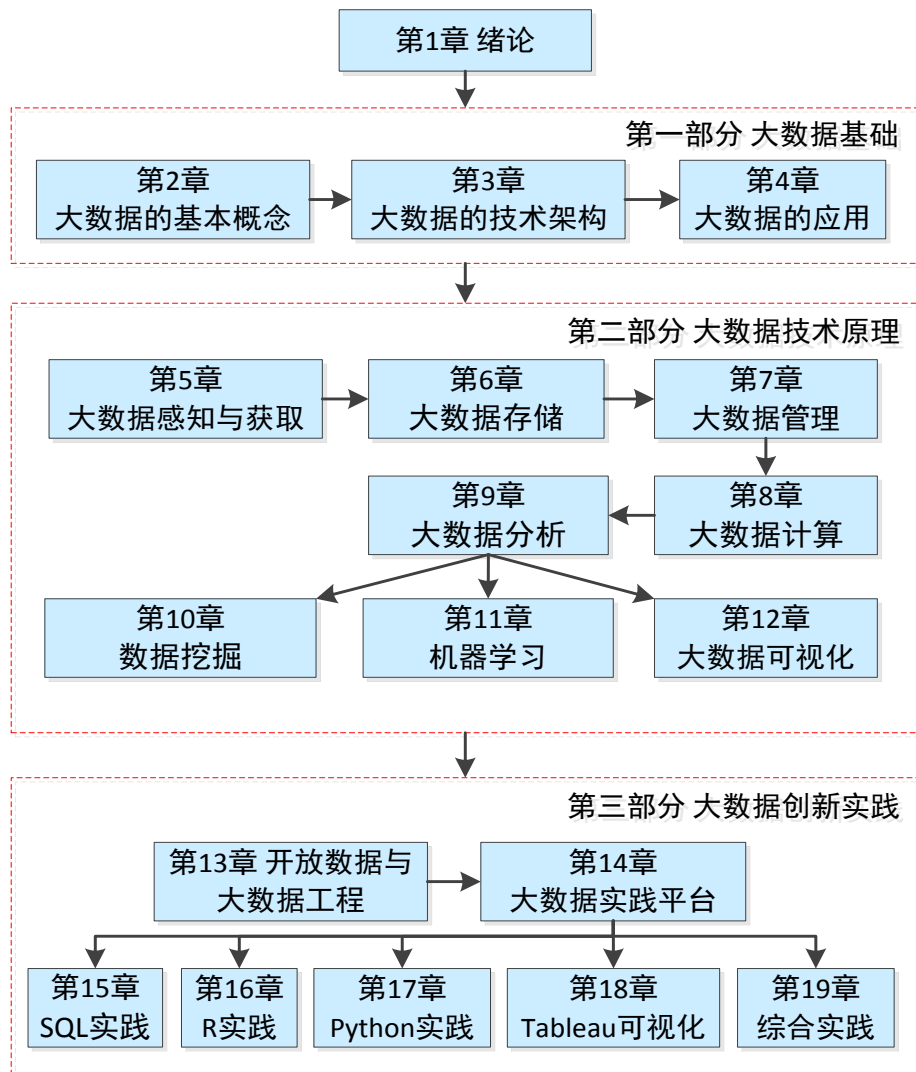
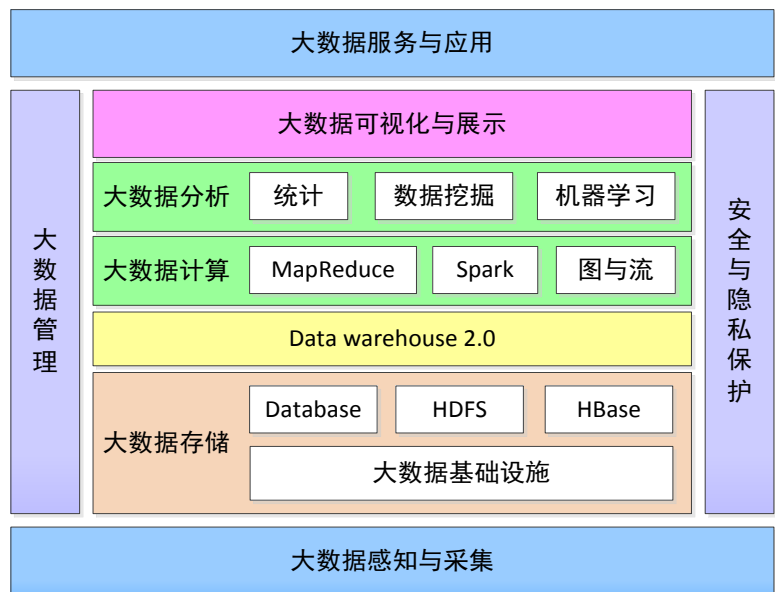


2016年

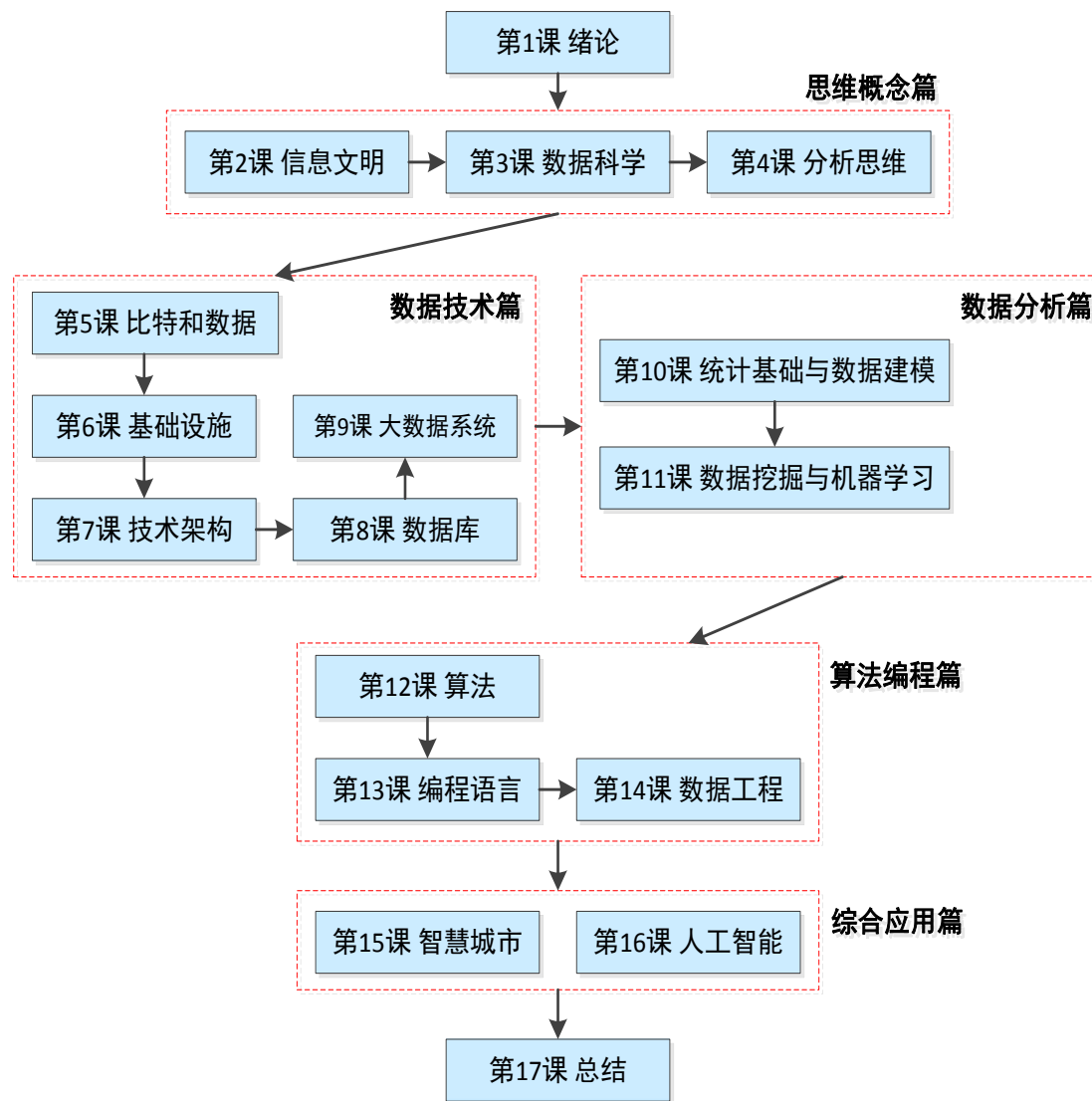
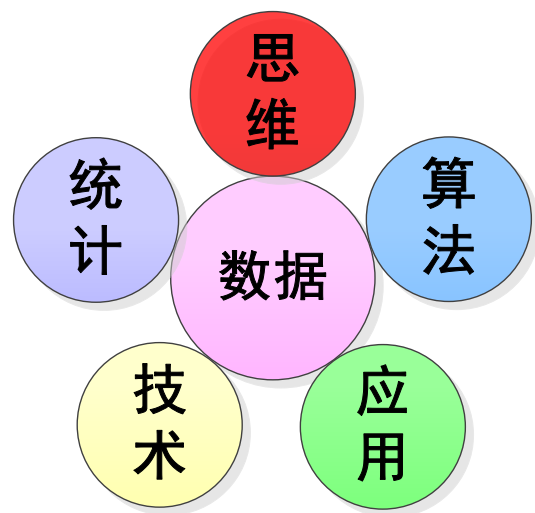


2017年

《大数据原理与实践》课程总览

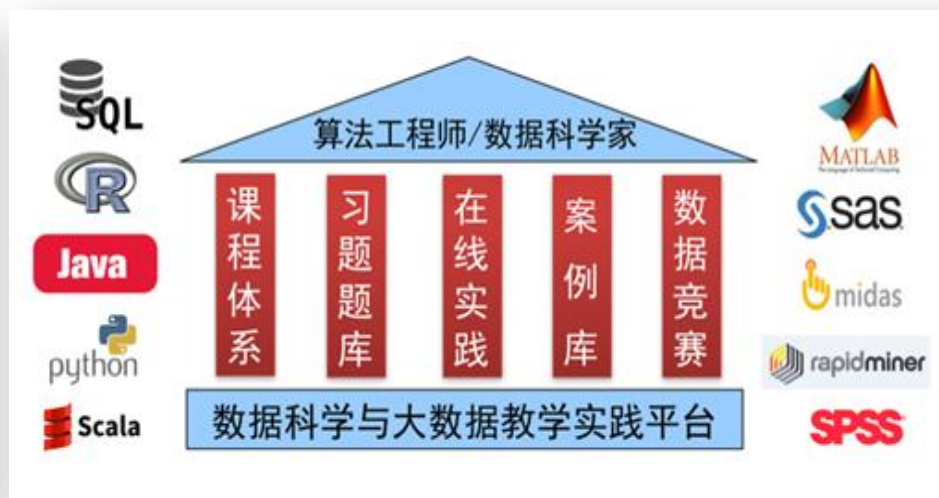


《数据科学通识导论》课程总览



“自由博雅”的16个字实践

- 建立对话；激发思辨；
- 协作交流；动手实践。



1. 建立对话、激发思辨



课件



文章



互动

嘉数汇



微信公众号每周更新模式

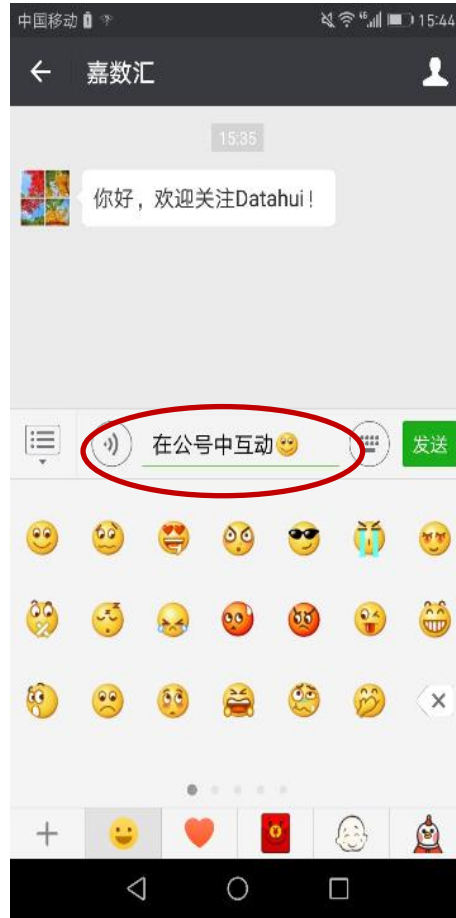
- **周二**：发布本周课件初稿，授课；
- **周三、周四**：互动、点评与问答；
- **周五**：发布最终版课件以及相应文本注释；
- **周末**：课外阅读文章。



如何回答与互动



在公号中输入：**姓名+学号**



在公号中互动



在文章中留言

互动交流

《数据科学通识导论》互动交流 - 201701

期



2017年第01期互动交流。



《数据科学通识导论》互动交流 - 20170

2期

《数据科学通识导论》互动交流 - 201703

期



每周我们会选取和汇总几条比较有代表性的学员问答进行反馈，希望通过这种方式达到交流的目的，启发大家进一步的思考。



《数据科学通识导论》互动交流 - 201704期



《数据科学通识导论》互动交流 - 201705



编者按：每周我们会选取和汇总几条比较有代表性的学员问答进行反馈，希望通过这种方式达到交流的目的，启发大家进一步的思考。

《数据科学通识导论》互动交流 - 201705-

06期



编者按：每周我们会选取和汇总几条比较有代表性的学员问答进行反馈，希望通过这种方式达到交流的目的，启发大家进一步的思考。

《数据科学通识导论》互动交流 - 201708

期



编者按：每周我们会选取和汇总几条比较有代表性的学员问答进行反馈，希望通过这种方式达到交流的目的，启发大家进一步的思考。

《数据科学通识导论》互动交流 - 201709

期



编者按：每周我们会选取和汇总几条比较有代表性的学员问答进行反馈，希望通过这种方式达到交流的目的，启发大家进一步的思考。

更新于 星期一 - 12:03



钟形曲线的陷阱：我们能够冲破阶层固化的束缚吗？



建立对话、激发思辨

“ZZZZB”:

“数据科学”的
分析，就其本身
直观地展现出
边可见的例子
到积累，等积累足够多的



张翔 1452229 (Ice)

昨天 23:16

第一周讨论内容: 1. 通识课 我认为就是以科普为目的的课程 用于拓宽知识

“suyi”:

我认为通识教育是培养一个人如何理解世界、分析世界和表达自我的能力，培养一个人在满足温饱之后想要达到更高通识教育是为了达到不断学习



张云笛+

1, 通识
业, 或方
一起上的
课我觉得
五六节课
成一些自
参与度,
修课, 时
些概念而
时候讲一
个课程。
老师就可

会很好。然后“数据科学”嘛，还是希望老师能够讲一些关于数据的算法之类的。 ps.老师，我想问一问，我现在有一点r和python的编程基础，但是直接接触kaggle又会有些困难，这个时候，我应该怎么学习啊？



朱涵林+1352356 (朱涵林)

星期四 10:10

数据科学:以数据为中心的科学。其涵盖面也非常广,存储,索引,科学计算,可视化等方面都是数据科学的范畴。学习数据科学,首先应该知道数据的分类,结构化数据,半结构化数据,非结构化数据。再根据数据的特点,找到适合它的存储方式,是传统关系数据库还是分布式存储系统。针对不同的系统,制定一定的策略进行索引和信息检索。为了从数据中得到我们所需的信息,需要对数据进行清洗,使用机器学习的方法对其进行分类,特征提取,或是建立相关模型进行预测,为了有更高的效率也许需要分布式计算。当有了数据及其分析结果,为了将其展示出来,向别人更好的表达自己的成果,数据可视化也必不可少。。。。数据科学四个字,其实是计算机领域和统计学领域众多学科相互协作的集合体。想要好好学习,既要有概念性的认识和总体的认识,比如数据的分类,什么样的数据适合什么样的策略,针对现有数据和算法我们能干什么,不能干什么,不同架构不同模块有什么特点和优劣等。有了对数据科学总体的认识后,可以挑选自己感兴趣的方向进行进一步的学习和研究,在算法效率,准确度等方面进行严谨的思索和探索。在“通”和“专”上不断的往复学习,才能学好数据科学。

2. 协作交流、动手实践

数据科学与大数据教学实践平台

在这里拥抱大数据



< >



教学

告别传统的“白板”教学方法，创造全新模式，在这里会发布与课堂内容同步的课后教学任务，学生们可要及时查看哦！



实践

我们提供各种数据，利用课堂上学到的数据分析的方法，在平台的在线编辑器中亲自实践，免去各种软件工具的安装。



竞赛

有竞争才有进步，我们会不定期的举行各种数据集下的大数据竞赛，让大家通过比赛发现隐藏在数据之后的真相。

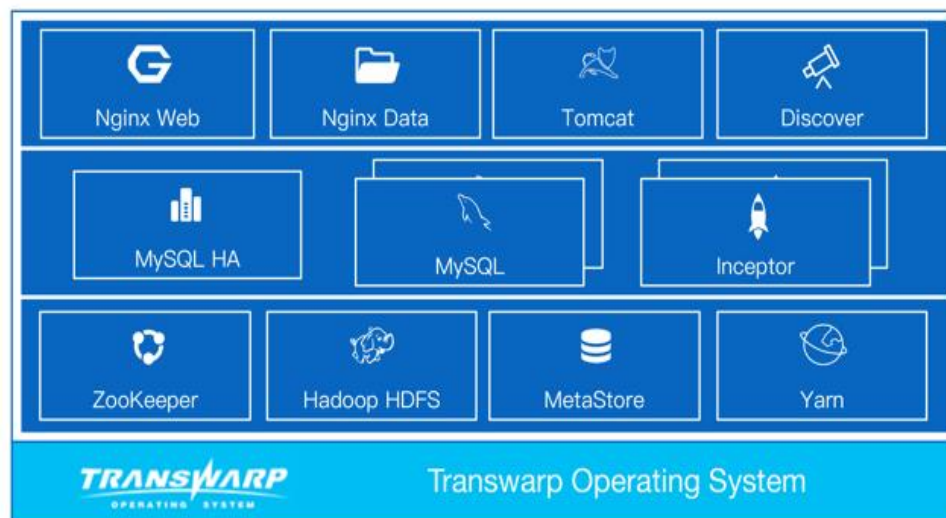
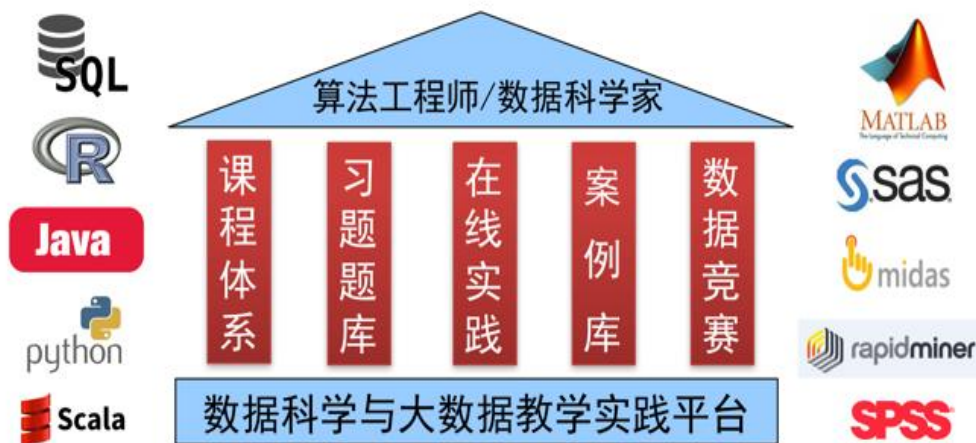


案例

所谓“教学相长”，我们会定期更新有关大数据的各种案例，看看别人的思路对你是否有所启发，通过学习案例实现进步。

数据科学与大数据实践平台

- 围绕**大数据和数据科学**提供课程教学、习题练习、动手实践、视化案例库、数据竞赛等服务。
- 通过**校企合作**，对接广泛的大数据基础软件；
- 通过自主研发的**云件技术**，在线学习使用任意数据分析工具。
- 形成一个**学、做、测、竞、评、奖**的大数据在线学习闭环。

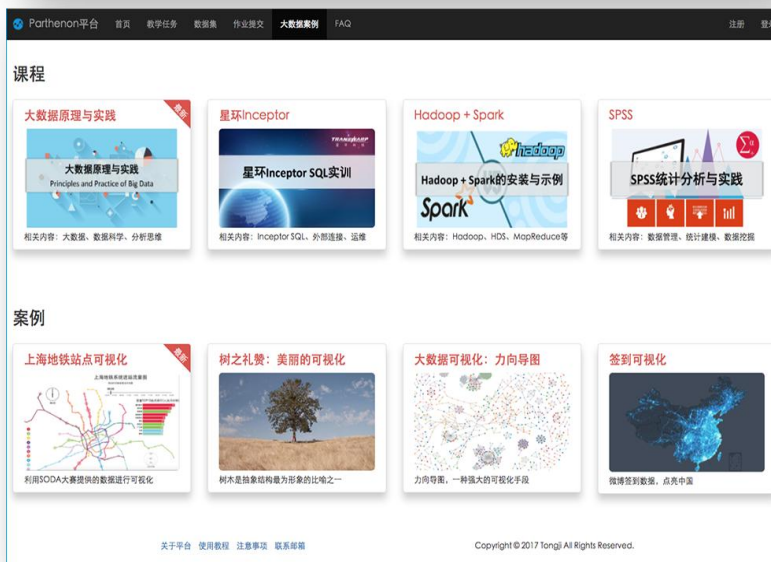


课程

- 通识基础
- SQL实践
- 平台搭建
- 数据分析
- 工具使用

案例

- 经典案例
- 学生作品
- 创新创意



在线实践

- 课程联动
- 方便直观
- 多语言支持
- 学练测闭环
- 在线竞赛
- 数据马拉松

The screenshot displays the Parthenon platform interface. At the top, there is a navigation bar with links for 'Parthenon平台', '首页', '教学任务', '数据集', '作业提交', '大数据案例', and 'FAQ'. On the right side of the navigation bar, there are links for '注册' and '退出'. Below the navigation bar, the page title is '数据集: 口碑客流量'. There are three buttons: 'SQL', '下一步', and '全部提交'. A dropdown menu is open, showing 'SQL', 'R', and 'Python'. The main content area is a code editor with a dark background. It contains two SQL scripts. The first script is in Chinese and provides instructions for using the SQL editor. The second script is in English and performs a complex SQL query. Below the code editor, there is a '导出成CSV' button and a search bar. The search bar contains the text '搜索:'. Below the search bar is a table with 10 columns: 'shop_id', 'city_name', 'location_id', 'per_pay', 'score', 'comment_cnt', 'shop_level', 'cate_1_name', 'cate_2_name', and 'cate_3_name'. The table contains 6 rows of data.

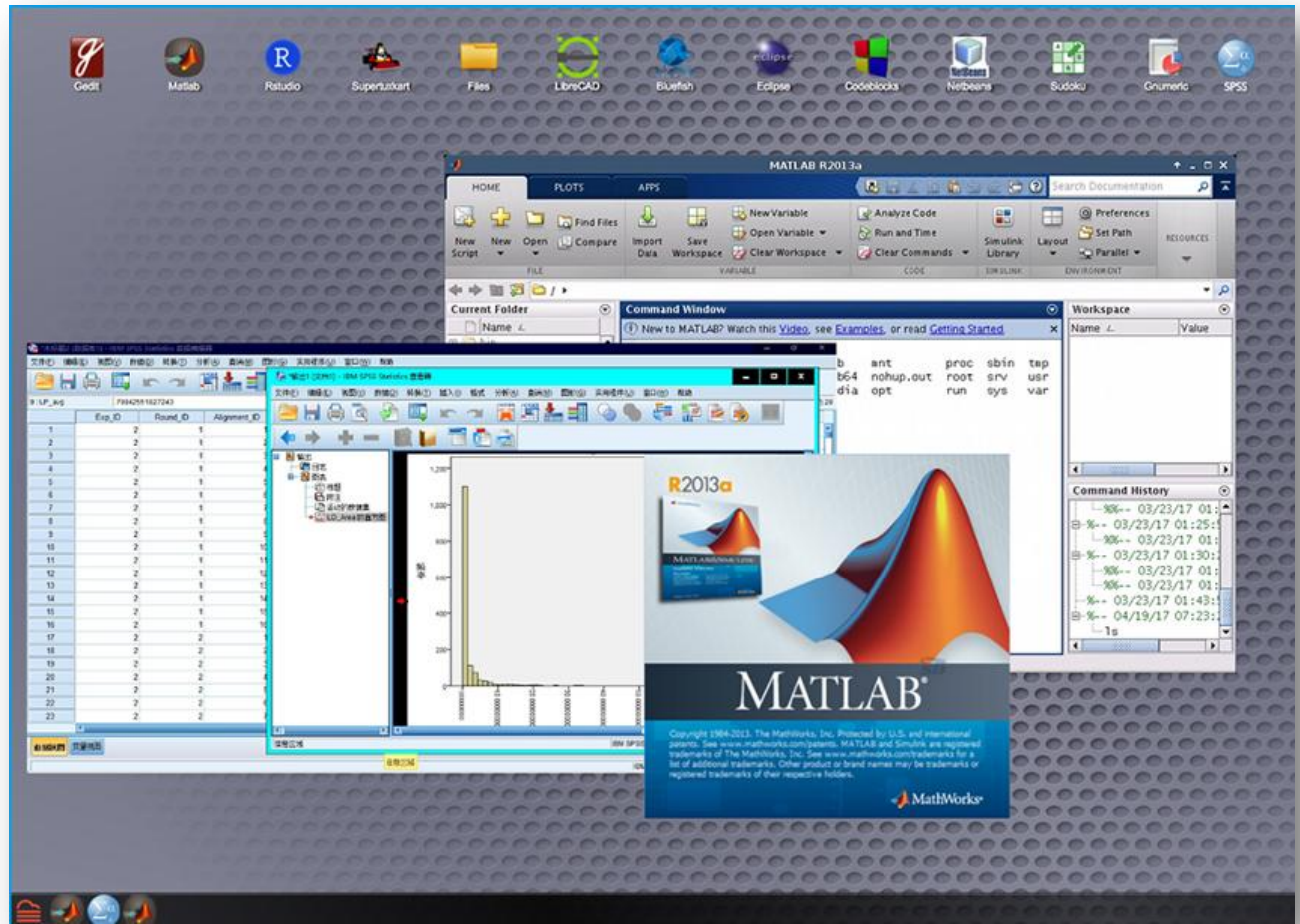
```
1 -- SQL编辑器 Beta 1
2 -- 1.所有数据表为只读，且不能创建新表
3 -- 2.一次提交多条语句时，分步每步执行一条
4 -- 3.修改未执行的语句后，需要重新全部提交才能生效
5
6 -- 以下操作查看该数据集的所有表
7 show tables;
8 select * from shop_info limit 20;
9 select * from user_pay_sum_day limit 20;
10
11 select c.shop_id,c.city_name,c.time_stamp,c.sum_day,d.high_tmp,d.low_tmp,d.weat
    .time_stamp,a.sum_day from user_pay_sum_day as a left join shop_info as b o
    city_weather as d on c.city_name=d.city_name and c.time_stamp=d.time_stamp
    .shop_id,c.time_stamp;
12
13 select shop_info.shop_id,shop_info.city_name,city_weather.time_stamp,city_weath
    .weather from shop_info right join city_weather on (shop_info.city_name=ci
    .time_stamp>='2016-11-01' and city_weather.time_stamp<='2016-11-14') where
    city_weather.time_stamp<='2016-11-14' order by shop_info.shop_id,city_weath
14
15
```

```
1 -- R 编辑器 Beta 1
2 # 1.一次提交多条语句时，分步每步执行一条
3 # 2.修改未执
4 for (i in 1:length(shangjia)) #每家商铺逐个加入训练集，整体训练
5 {
6   cur_shop<-shangjia[i]
7   xunlian_data<-ori_data[which(ori_data$shop_id==cur_shop),]
8   c<-xunlian_data$sum_day[1:21] #三个月流量
9   jingji=ecno[cur_shop] #经济指数
10  mean=round(mean(c),3) #三个月均值
11  sd=round(sd(c),3) #三个月标准差
12  cmax=max(c) #三个月最大值
13  cmin=min(c) #三个月最小值
14  cmedian=round(median(c),3) #三个月中位数
15  #win周es=NULL
16  #for(i in 1:2) #滑动窗口特征
17  #[
18  #win_week=xunlian_data$sum_day[(22-7*):21] #滑动窗口
19  #mean=round(mean(win_week),3) #滑动均值
20  #median=round(median(win_week),3) #滑动中位数
21  #max=max(win_week) #滑动最大值
```

shop_id	city_name	location_id	per_pay	score	comment_cnt	shop_level	cate_1_name	cate_2_name	cate_3_name
1	湖州	885	8	4	12	2	美食	休闲茶饮	饮品/甜点
2	哈尔滨	64	19			1	超市便利店	超市	
3	南昌	774	5	3	2	0	美食	休闲茶饮	奶茶
4	天津	380	18			1	超市便利店	超市	
5	杭州	263	2	2	2	0	美食	休闲食品	生鲜水果
6	大连	1139	13	3	1	0	美食	烘焙糕点	面包

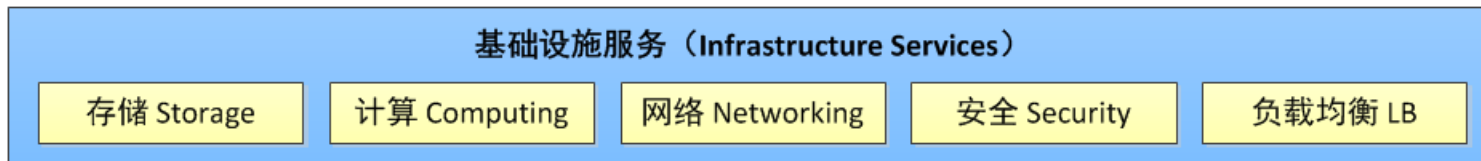
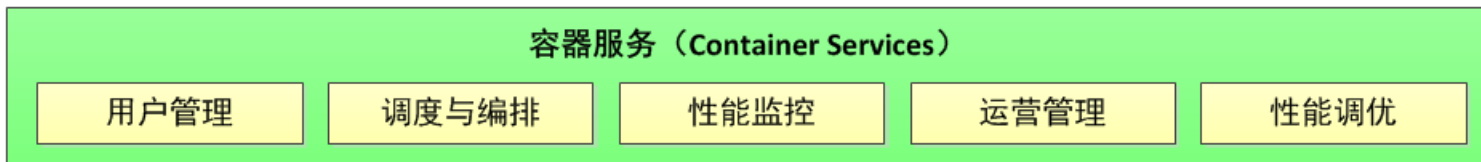
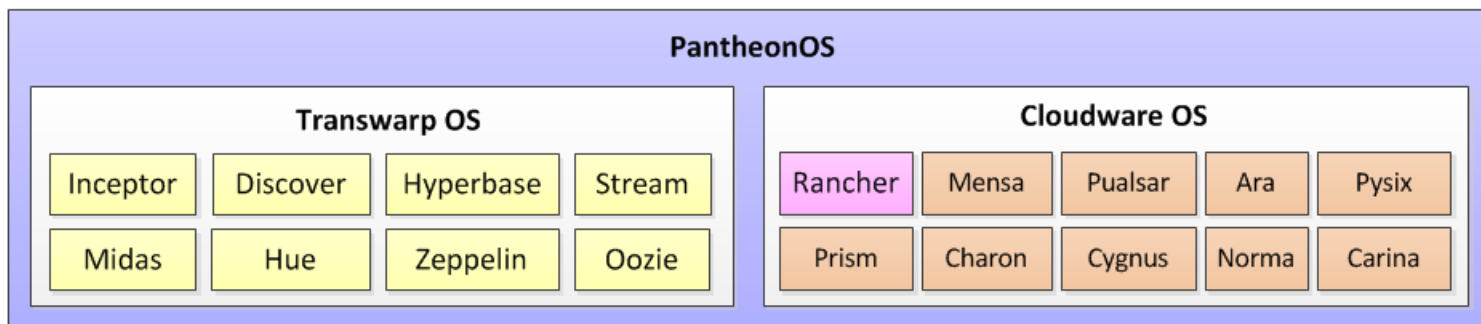
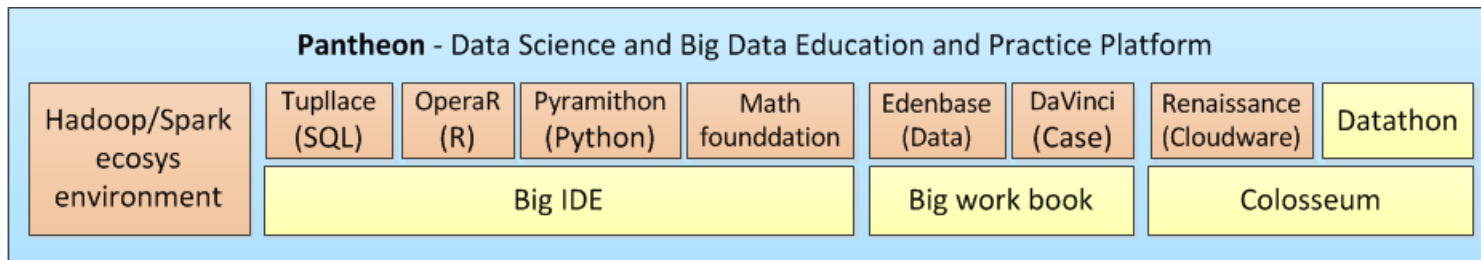
云件工具

- 简洁方便
- 桌面工具
- 数据工具
- 存储服务
- 在线指导
- 团队协作



Dong Guo, Wei Wang*, et al., Cloudware: An Emerging Software Paradigm for Cloud Computing, In Proceedings of the Internetwork 2016, Beijing, China, September 18, 2016.

技术架构



作品展示：上海地铁系统进站流量图

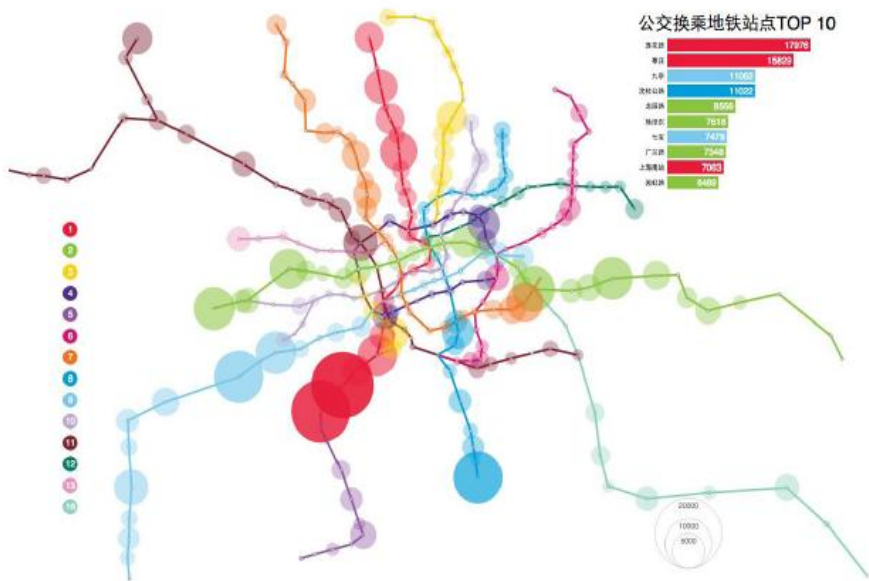
工作日

公交换乘地铁站点统计图

工作日 休息日

公交换乘地铁站点TOP 10

南京路	17976
曹路	16409
大丰	11009
沈杜公路	11002
漕河泾	8954
梅川路	7818
七宝	7473
广兰路	7243
上海曹路	7043
颛桥路	6403



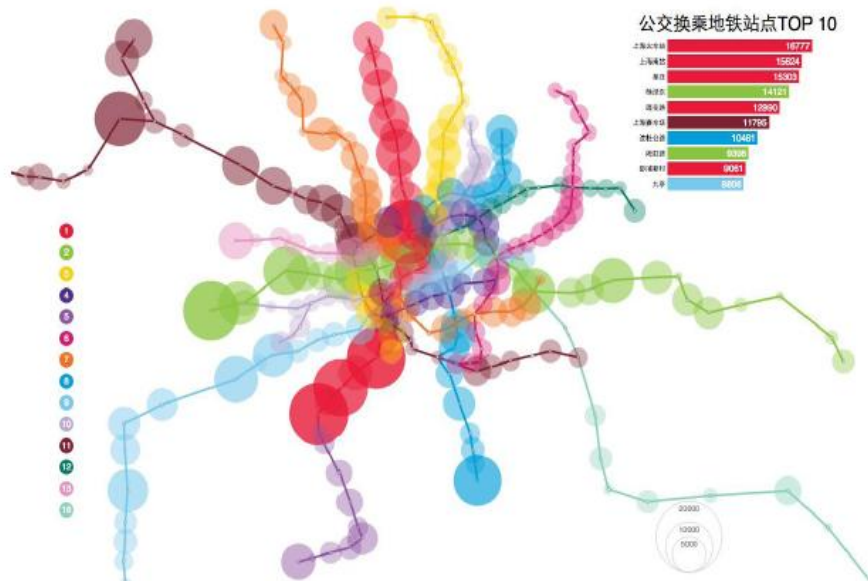
休息日

公交换乘地铁站点统计图

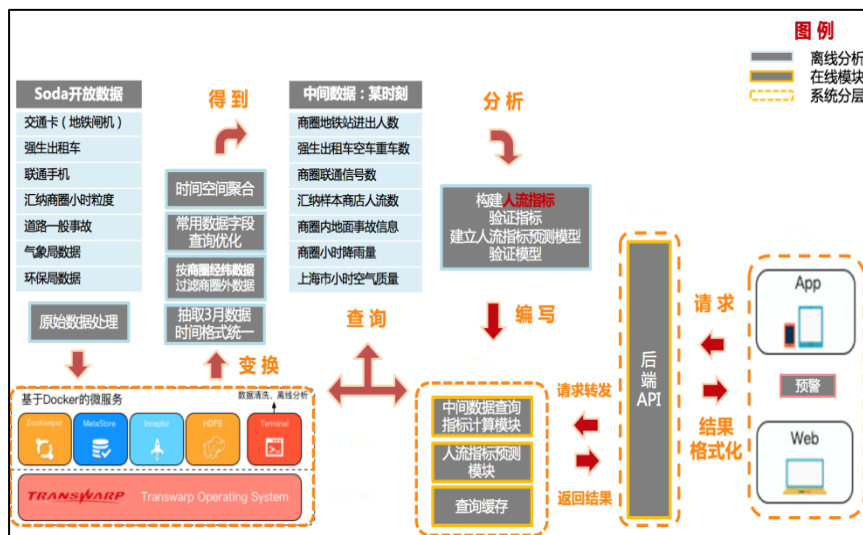
工作日 休息日

公交换乘地铁站点TOP 10

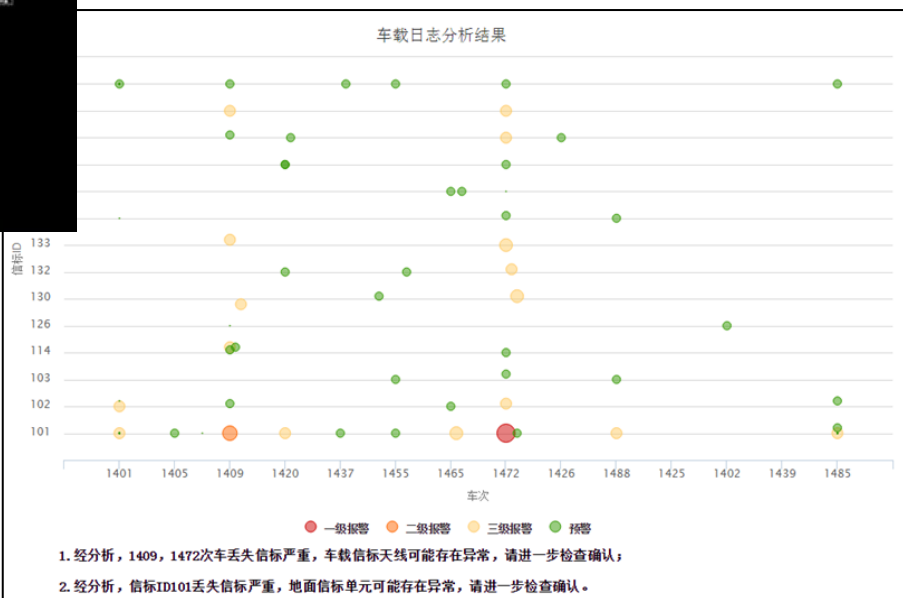
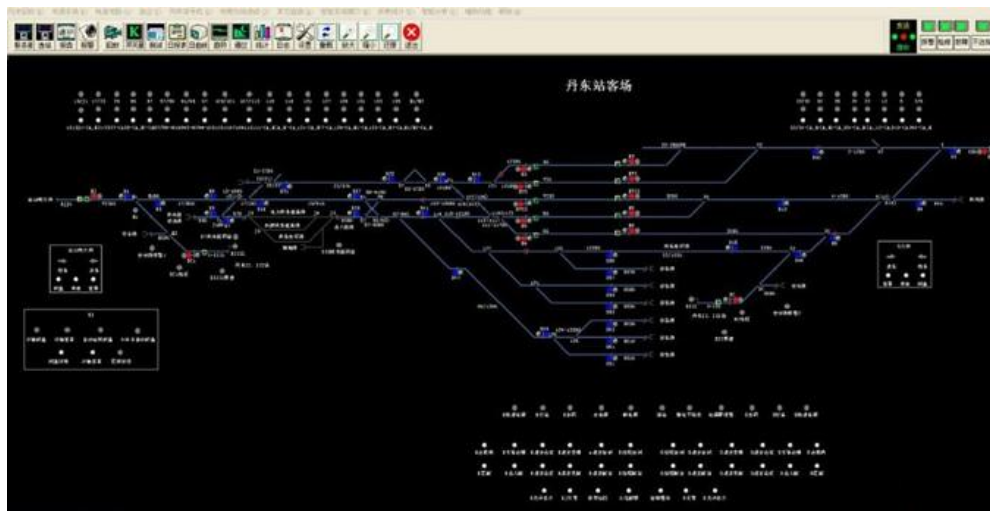
上海火车站	16777
上海城隍庙	15204
曹路	14933
梅川路	14113
漕河泾	12890
上海曹路	11785
静安寺	10461
梅川路	8338
静安寺	8081
大丰	5883



作品展示：基于人流指数预测的商圈公共安全预警系统



作品展示：轨道交通运维大数据分析



总结：未来的挑战

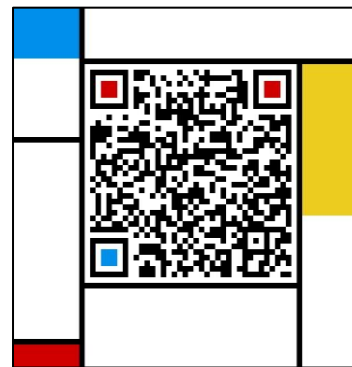
- 大数据与数据科学知识体系的沉淀
- 通识类课程的改革（取代计算机基础？）
- 新工科背景下大数据与数据科学实践平台
- 真正做到以学生为本

Thanks!

课程公众号



个人微信



一个好的教育，是一个灵魂对另一个灵魂的呼唤；
一门好的课程，是一个生命与另一个生命的碰撞。