

# 大数据学科建设的**关键因素**

---

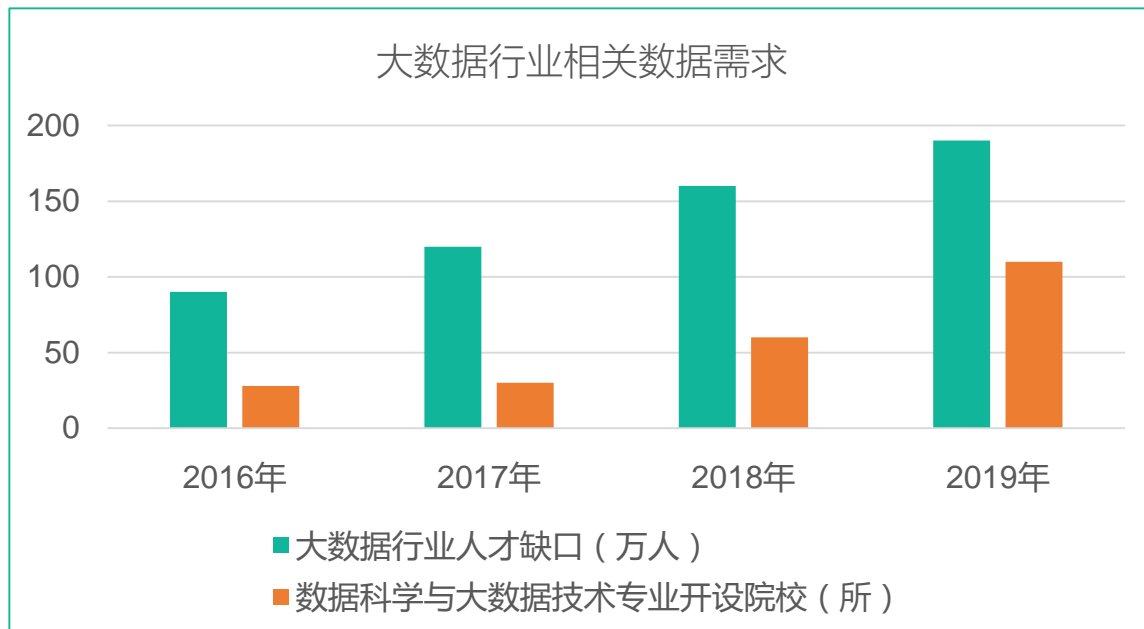
王涛

北京大学·北京大数据研究院博雅大数据学院 执行院长

大数据教育联盟 副秘书长

# 大数据学科是什么？

据全球最顶尖管理咨询公司麦肯锡 (McKinsey) 出具的一份详细分析报告显示, 预计到2018年, 大数据或者数据工作者的岗位需求将激增, 其中大数据科学家的缺口在14万到19万之间, 对于懂得如何利用大数据做决策的分析师和经理的岗位缺口则将达到**150万**。



北京大学在鄂维南院士的带领下, 积极推动数据科学的学科建设, 并于2016年在北京大学率先开设“**数据科学与大数据技术**”专业 (专业代码080910T), 数据科学本科专业, 第一届学生2017年毕业。北京大学、对外经济贸易大学、中南大学成为首批获教育部批准开设本专业的高校。

2017年3月17日, 教育部公布第二批**32所获批高校**。包括人民大学、复旦大学、华东师范大学、北京信息科技大学等。预计2018-2019年, 国内将有上百所高校获批“数据科学与大数据技术”专业。2016年9月, 教育部公布新增“**大数据技术与应用**”**专科专业** (专业代码610215)。

数据科学作为一门新兴的大数据学科，所依赖的两个因素是：**一是数据的广泛性和多样性；二是数据研究的共性。**现代社会的各行各业都充满了数据，数据类型多种多样，包括传统的结构化数据，也包括网页、文本、图像、视频、语音等非结构化数据。

**数据科学主要包括两个方面：用数据的方法来研究科学和用科学的方法来研究数据。**前者包括生物信息学、天体信息学、数字地球等领域；后者包括统计学、机器学习、数据挖掘、数据库等领域。

1



全新的专业，大数据学科如何建设？

2



多学科知识交叉，对授课教师提出了新的挑战，大数据师资如何培养？

3



迫切需要全套课程产品，教材、讲义、课件、授课视频等，大数据教学产品如何研发？

4



大数据专业注重实践和应用，传统课堂教学方式如何变革，集成数据、分析环境、实际问题和讲义案例的大数据教学实训平台如何搭建？



北京大数据研究院于2015年8月27日，在北京市委市政府的支持与指导下，由北京大学、中关村管委会、海淀区政府、北京工业大学四方共同筹建。北京大数据研究院是国内首个整合了政府、大学和市场三方面资源的大数据研究机构，目标是吸引国际一流大数据研究人员来京发展，用五到十年的时间，建成国际一流的大数据教育、科研创新和产业化平台，成为中国乃至世界大数据产业发展的一面旗帜。



博雅大数据学院是由北京大数据研究院成立，专注于大数据教育产品研发、服务的机构。作为全国最优秀的大数据教育机构之一，博雅大数据学院在中国科学院院士、北京大数据研究院院长、北京大学元培学院院长鄂维南院士的带领下已拥有超过百名专业研发人员以及众多提供大数据优质授课服务的专职教师。





## 院长：鄂维南院士

千人计划，北京大学教授、北京大学元培学院院长，美国普林斯顿大学教授。973项目“非结构化数据分析”首席科学家。



## 高文院士

大数据研究院学术委员会主任，国家自然科学基金委员会副主任，ACM/IEEE Fellow。



## 张平文院士

大数据研究院学术委员会主任、北京大学学科建设办公室主任。

## 国际人才引进

邵 骋	普林斯顿大学	汤林鹏	普林斯顿大学
章思鑫	纽约大学	张 立	爱荷华大学
朱占星	爱丁堡大学	李千骁	普林斯顿大学
周亚俊	哈佛大学		

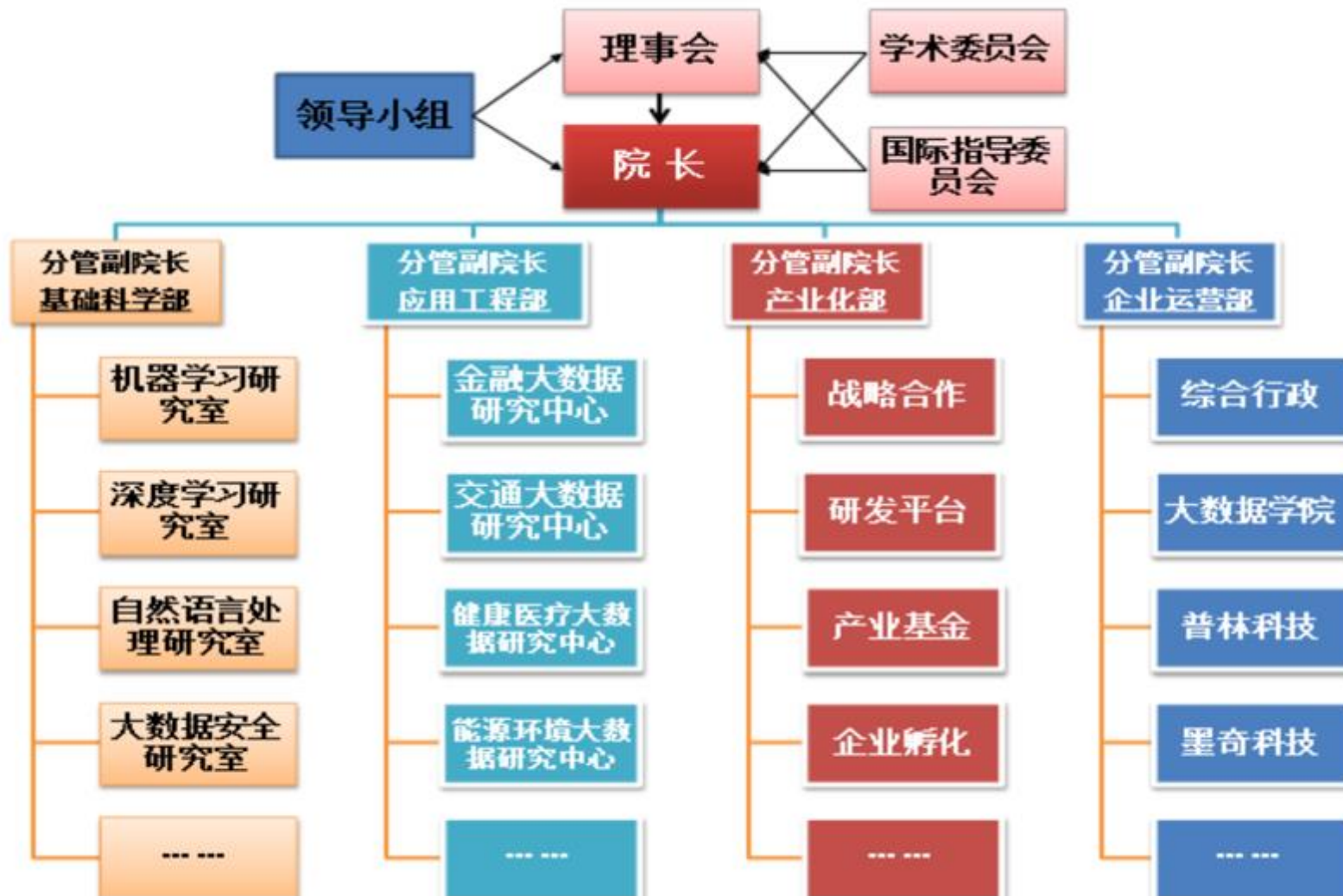
...

## 国内人才引进

张志华	上海交通大学	刘云淮	公安部三所
王亦伦	电子科技大学	严 睿	百度

...

# 北京大数据研究院组织架构





# 大数据学科建设的关键因素一

## “数据科学与大数据技术”（080910T）专业建设

“数据科学与大数据技术”专业强调培养具有多学科交叉能力的大数据人才。该专业重点培养具有以下三方面素质的人才：一是理论性的，主要是对数据科学中模型的理解和运用；二是实践性的，主要是处理实际数据的能力。三是应用性的，主要是利用大数据的方法解决具体行业应用问题的能力。

### 培养目标：

“数据科学与大数据技术”专业，培养德、智、体、美全面发展，掌握数据科学的基础知识、理论、及技术，包括面向大数据应用的**数学、统计，计算机**等学科基础知识，数据建模、高效分析与处理，统计学推断的基本理论、基本方法和基本技能。对自然科学和社会科学等**应用领域中大数据**的了解，具有较强的专业能力和良好外语运用能力，能胜任数据分析与挖掘算法研究和大数据系统开发的**研究型和技术型人才**。

# 大数据学科建设的关键因素一

## “数据科学与大数据技术”专业课程体系

传统课程



整合课程

相关模块核心知识

模块知识深入学习

《高等数学》  
《矩阵运算》  
《线性代数》  
...

数学

《大数据数学基础》

《概率论》  
《数理统计》  
...

统计学

《大数据统计基础》

《C程序设计》  
《Java程序设计》  
《操作系统原理》  
《数据结构》  
《数据库概论》  
...

计算机科学

《大数据分析的Python基础》

《数据存储》

数据分析模块

《数据科学导引》

⋮

《机器学习》

计算模块

《最优化算法》

计算机技术模块

《数据采集》

⋮

《分布式概论》

大数据应用模块

《大数据应用导论》

数据分析模块

《文本数据分析》

《图与网络数据分析》

⋮

《深度学习》

《人工智能导论》

计算模块

《大数据分析的算法》

计算机技术模块

《Web技术概论》

《云计算与大数据平台》

⋮

《数字图像处理》

大数据应用模块

《医疗健康大数据应用》

⋮

《金融大数据应用》

专业基础课

专业核心课

专业选修课

# 大数据学科建设的关键因素一

## “大数据技术与应用”（610215）专业建设

“大数据技术与应用”专业强调培养具有大数据实践能力的大数据人才。该专业重点培养具有以下两方面素质的人才：一是工具的掌握，掌握数据采集和数据分析的基本工具；二是数据分析能力，掌握实用数据分析和初步数据建模能力。

培养目标：

“大数据技术与应用”专业，培养掌握数据科学的**基础知识**及大数据相关技术，掌握大数据清洗和分析常用**工具**的使用，具有卓越的实践能力，能胜任**数据清洗、数据存储、数据分析与挖掘、大数据系统开发与构建**等工作的**专业应用型人才**。

# 大数据学科建设的关键因素一

## “大数据技术与应用”专业课程体系

### 语言和专业基础

《大数据的语言基础Python》  
《Linux系统基础》  
《大数据的统计基础》  
...

专业基础课

### 数据采集、存储与处理

《数据存储(MySQL)》  
《数据采集与网络爬虫》  
《数据清洗》  
《大数据处理工具  
(Excel/Weka/Pandas)》  
...

专业核心课

### 大数据分析、开发与应用

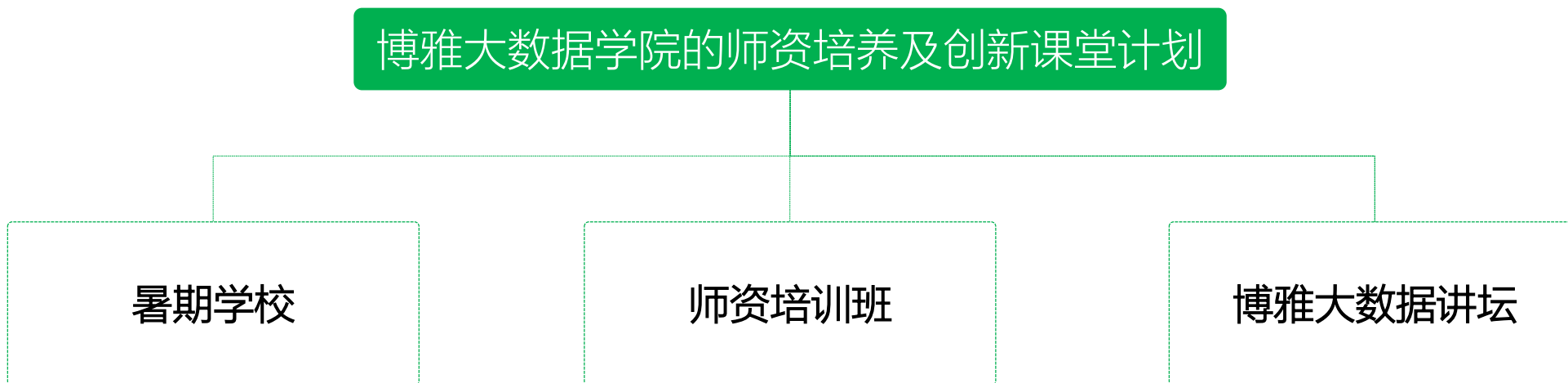
《数据分析导论》  
《大数据行业应用导论》  
《数据可视化》  
《Hadoop大数据平台基础》  
...

专业选修课

# 大数据专业建设的关键因素二

## 大数据师资队伍培养

作为新兴的交叉型学科，相关的大数据课程涉及数学、统计和计算机等各种知识的学习与实践，对教师有着较高的要求。众多高校针对本专业的建设尚处于摸索阶段，能否有效地帮助高校培养专业的师资团队并顺利开展教学工作，成为了现阶段的首要任务。



# 大数据专业建设的关键因素二

## 大数据师资队伍培养

### 暑期学校

- 2016年7月-8月 在北京大学举办了第一届大数据暑期学校“大数据的模型和算法”
- 500多名高校青年教师报名500余名
- 每年举办1次
- 2017年暑期班将在7月开始，鄂维南教授将亲自讲授《机器学习的数学导引》课程



### 师资班

- 北京大数据研究院博雅大数据学院专业师资团队授课
- 结课：授予北京大数据研究院认证证书（大数据分析师、大数据讲师）
- 课程包括《数据科学导论》《大数据分析的原理与技术》《机器学习的数学导引》《大数据行业应用解析（交通、医疗、金融、工业等）》《大数据分析的Python基础》《数据清洗技术与工具》《数据采集与网络爬虫》《数据可视化》。



# 大数据专业建设的关键因素二

## 大数据师资队伍培养

### 博雅大数据讲坛

- 定期邀请国际国内知名大数据专家和行业专家来北大进行大数据讲座，同时在大数据教育实训平台数据嗨客上进行网络直播
- 讲座作为院校大数据课程的一部分，融入课程计划

部分往期讲座回顾：

讲座题目	演讲嘉宾	嘉宾介绍
大数据对未来商业的改变	鄂维南院士	北京大数据研究院院长
大数据在金融行业的应用	李铭	中国人民银行征信中心专家
大数据在交通行业的应用	陈艳艳	北京工业大学城市交通学院院长
大数据在环保行业的应用	张平文	中科院院士
大数据在医疗健康行业的应用	李全政	哈佛医学院教授
人工智能	汪军	伦敦大学学院计算机系教授

# 大数据专业建设的关键因素三

## 大数据教学产品

### 博雅大数据学院大数据系列课程产品

- 《数据科学导论》《大数据分析的原理与技术》《机器学习的数学导引》  
《大数据行业应用解析（交通、医疗、金融、工业等）》  
《大数据分析的Python基础》《数据清洗技术与工具》  
《数据采集与网络爬虫》《数据可视化》等
- 2017年6月，由鄂维南院士主编的《数据科学导论》即将出版，高等教育出版社

### 其他教学产品

- 讲义、实战案例、线上实训题库、课程教学视频（与智慧树网合作）



# 大数据专业建设的关键因素四

## 大数据教学实训平台 - 数据嗨客

数据嗨客是北京大数据研究院和博雅大数据学院经过两年多研发的大数据教育实训平台。为高校大数据教学及企业数据人才培养提供线上实训环境及教学资料。让学生通过线上自主学习及实战演练，理解大数据科学的原理，掌握数据科学的知识体系，真实体验大数据建模分析的实际操作与演练过程。

2016年9月，数据嗨客率先在北京大学和南方科技大学投入实训教学。

2017年，数据嗨客在武汉大学、北京信息科技大学、华中科技大学、贵州大学等高校将投入教学。



<http://www.hackdata.cn/>

- 权威的知识体系
- 在线大数据实训
- 随时随地线上练习
- 支持主流语言（Python和R）
- 自动效果评估
- 便捷的教学管理和在线作业管理

# 大数据专业建设的关键因素四

## 大数据教学实训平台 - 数据嗨客

### 数据嗨客的5大特点

全面的大数据知识讲义和丰富的案例资源。将大数据知识点进行有机整合，配套丰富行业案例，并提供大数据建模全流程指导

组织多种奖金丰厚的大数据竞赛，让大数据人才将所学知识和实战技能真正用于解决企业面临的挑战性问题

练习

学习

教室

竞赛

工作

基于海量题库在线进行行业数据建模演练。涵盖多种体系，支持讨论和数据探索功能，实现自动模型训练、测试和评估，支持主流数据分析语言

专业的大数据教学管理和作业管理功能。除基本的教辅功能，支持针对大数据教学的在线作业发布、提交、评测和互动答疑等功能

利用用户画像技术构建大数据人才知识和技能图谱，实现人才和企业的有效匹配

# 大数据专业建设的关键因素四

## 大数据教学实训平台 - 数据嗨客

深入浅出的系统讲义

课程讲义：数据科学导引

The screenshot shows the Data Hacker website interface. At the top, there is a navigation bar with the logo and menu items: 学习 (Learn), 练习 (Practice), 教室 (Classroom), 排名 (Ranking), and 竞赛 (Competition). Below the navigation bar, there is a grid of course cards. Each card features a cover image, a title, and a button to start learning. The visible cards are:

- 数据科学导引 (Data Science Introduction) - 开始学习 (Start Learning)
- 大数据分析的Python基础 (Python Basics for Big Data Analysis) - 开始学习 (Start Learning)
- 数据清洗 (Data Cleaning) - 敬请期待 (Coming Soon)
- 数据采集 (Data Collection) - 开始学习 (Start Learning)
- 大数据分析的R基础 (R Basics for Big Data Analysis) - 开始学习 (Start Learning)

The screenshot shows the 'K近邻算法' (K-Nearest Neighbor) page. On the left is a sidebar menu with various machine learning topics, with 'K近邻算法' highlighted. The main content area includes:

- 标题:** K近邻算法
- 收藏:** 收藏
- 正文:** K近邻 (K-Nearest Neighbor, KNN) 是一种最经典和最简单的有监督学习方法之一。当对数据的分布只有很少或者没有任何先验知识时，K近邻算法是一个不错的选择。K近邻算法既能够用来解决分类问题，也能够用来解决回归问题。该方法有着非常简单的原理：当对测试样本进行分类时，首先通过扫描训练样本集，找到与该测试样本最相似的那个训练样本，根据这个样本的类别进行投票确定测试样本的类别。也可以通过这个样本与测试样本的相似程度进行加权投票。如果需要以测试样本对应每类的概率的形式输出，可以通过这个样本中不同类别的样本数量分布来进行估计。K近邻算法的预测过程示意图如图1所示：
- 图1 K近邻算法预测示意图:** A diagram showing a test sample (x) being classified based on its K nearest neighbors (k=5). The neighbors are represented by different colored shapes (blue squares, green circles, orange triangles) and their majority class is used for prediction.
- 正文:** K近邻算法有着坚实的理论基础，研究者已经证明，当  $K = 1$  时，K近邻算法的泛化错误率上界为贝叶斯最优分类器错误率的两倍。具体证明如下：  
假设样本独立同分布，且对任意测试样本  $x$  和任意小正数  $\delta$ ，在  $x$  附近  $\delta$  距离范围内总能找到一个训练样本  $z$ ，令  
$$c^* = \arg \max_{c \in Y} P(c|x)$$
  
表示贝叶斯最优分类器的结果，则最近邻分类器出错的概率是  $x$  与  $z$  类别标记不同的概率，即  
$$\begin{aligned} \epsilon &= 1 - \sum_{c \in Y} P(c|x)P(c|z) \\ &\leq 1 - P^2(c|x) \\ &\leq 1 - P^2(c^*|x) \\ &\leq 2 \times (1 - P(c^*|x)) \end{aligned}$$
  
在对测试样本进行预测时，因为只用到训练样本中与其最接近的  $K$  个样本，K近邻算法的偏差 (Bias) 往往很低，而方差 (Variance) 则很高。当训练集较小的时候，K近邻算法容易出现过拟合。  
算法的详细流程：

# 大数据专业建设的关键因素四

## 大数据教学实训平台 - 数据嗨客

课程讲义：大数据分析Python基础

MagicFrame：理论和实操融为一体

数据嗨客 HackData

学习 练习 教室 排名

初识Python

基本概念

科比投篮数据集介绍

**变量和常量**

注释：帮助理解代码

print函数

数据类型

算术运算符

类型转换

数据的容器

数据中的字符编码

控制结构

数据文件的读写操作

编写函数处理数据

Python的工具箱：模块

Python语言的类

解除和控制警报

### 变量和常量

收藏

变量这个概念在中学时期我们就已经非常熟悉了，但是中学时期的变量往往只是数字。在计算机科学领域，变量不仅可以是数字，还可以是其他数据类型。变量是所有编程语言重要的组成部分。

我们常常使用变量来存储一些之后可能会变化的值。在创建变量时，我们需要给变量取一个名称，然后给变量赋值。在这之后，我们就可以通过变量名称来代替所赋予的值。

如果我们需要对科比投篮ID为 1 的一次投篮进行分析，那么我们就可以创建一个名称为 `shot_id` 的变量，并且将 1 值储存在变量 `shot_id` 中。之后如果我们在程序中需要使用投篮ID，则我们直接可以使用 `shot_id`。使用变量的一个好处在于：如果之后我们想要分析科比的另外一次投篮，比如投篮ID为 2 的投篮，我们只需要修改变量 `shot_id` 的赋值，将 `shot_id` 赋值为 2 即可，而不需要在每一个使用投篮ID的位置进行修改。

#### 变量命名规则

变量在程序中通过一个变量名表示，Python语言的变量命名规则为：变量名必须是大小写英文字母、数字或下划线 `_` 的组合，不能用数字开头，并且对大小写敏感。变量名应简单易懂，尽量保证在数周甚至数月之后看到变量名仍然能够想起来它所代表的含义。

需要注意的是，Python语言中存在一些关键字，这些关键字不能用于命名变量：`and`、`as`、`assert`、`break`、`class`、`continue`、`def`、`del`等。我们可以通过 `import keyword` 导入 `keyword` 包，并使用方法 `keyword.kwlist` 查看写关键字。

开始答题

① 说明

- 我们已经将 1 值赋予给了变量 `shot_id`，请将变量 `shot_id` 的赋值改为 2

显示范例

```
1 shot_id = 2
```

清空 提交

👍 答对了，真棒!

结果输出

# 大数据专业建设的关键因素四

## 大数据教学实训平台 - 数据嗨客

### 案例学习：真实数据的行业案例

#### 使用K近邻算法诊断乳腺癌 (Python)

乳腺癌已成为当前社会的重大公共卫生问题。全球乳腺癌发病率自20世纪70年代末开始一直呈上升趋势。美国8名妇女一生中就会有1人患乳腺癌。中国不是乳腺癌的高发国家，但不宜乐观，近年我国乳腺癌发病率的增长速度却高出高发国家1~2个百分点。据国家癌症中心和卫生部疾病预防控制局2012年公布的2009年乳腺癌发病数据显示：全国肿瘤登记地区乳腺癌发病率位居女性恶性肿瘤的第1位，女性乳腺癌发病率（粗率）全国合计为42.55/10万，城市为51.91/10万，农村为23.12/10万。

早期的乳腺癌检测主要检查乳腺组织的异常肿块。如果发现一个肿块，那么就需要进行细针抽吸活检，然后在显微镜下检查细胞，从而确定肿块是良性还是恶性。如果能够通过机器学习建模，通过乳腺肿块的检测数据自动进行诊断，将会给医疗系统带来很大的益处。一方面，自动诊断能够大大提高检测效率。另一方面，结合大量不同历史案例的自动诊断能够使辅助医生进行决策，降低误判的风险。

本案例中，基于公开的乳腺癌诊断数据，我们将使用机器学习中一种简单的算法--K近邻算法来构建乳腺癌自动诊断模型。

#### 1 数据源

我们将使用来自UCI的**乳腺癌诊断数据集**。该乳腺癌数据包括569例乳腺细胞活检样本，每个样本包含32个变量。其中id变量是样本识别ID，diagnosis变量是目标变量（M代表恶性，B代表良性）。其他30个变量都是由10个数字化细胞核的10个不同特征的均值、标准差和最大值构成。这10个基本特征为：

- radius（半径）

#### 2 数据探索和预处理

首先，使用pandas的read\_csv()函数将CSV格式的乳腺癌数据集导入数据框中。

```
import pandas as pd
breast_cancer = pd.read_csv("datasets/wisc_bc_data.csv")
print breast_cancer.shape
breast_cancer.head(10)
```

(569, 32)

	id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smooth
0	842302	M	17.99	10.38	122.80	1001.0	0.11840
1	842517	M	20.57	17.77	132.90	1326.0	0.08474
2	84300903	M	19.69	21.25	130.00	1203.0	0.10960
3	84348301	M	11.42	20.38	77.58	386.1	0.14250
4	84358402	M	20.29	14.34	135.10	1297.0	0.10030
5	843786	M	12.45	15.70	82.57	477.1	0.12780
6	844359	M	18.25	19.98	119.60	1040.0	0.09463
7	84458202	M	13.71	20.83	90.20	577.9	0.11890
8	844981	M	13.00	21.82	87.50	519.8	0.12730
9	84501001	M	12.46	24.04	83.97	475.9	0.11860

10 rows x 32 columns

第一个变量为ID变量，无法为实际的模型构建提供有用的信息，所以需要将其删除。

# 大数据专业建设的关键因素四

## 大数据教学实训平台 - 数据嗨客

海量题库：提供大量实训练习题，进行大数据实战演练

共46题

全部

#	题目列表	提交次数	通过率
420	钓鱼网站的识别	45	77.78%
414	南非西开普省冠心病分类	223	78.48%
401	酵母菌的蛋白质位置预测	85	97.65%
399	电离层反射的雷达波质量分类	54	100.00%
398	根据背景对人群的分类	55	98.18%
397	根据人名对徽章进行分类	43	100.00%
384	助教教学表现分类	37	100.00%
383	中文手机评论情感倾向判定	135	65.93%
377	Tictactoe游戏结果分类	54	98.15%
354	鸢尾属花的分类	58	87.93%
353	泰坦尼克号幸存者分类	39	100.00%
351	城市地表覆盖分类	32	100.00%
350	LED显示屏显示数字的辨别	22	100.00%
349	贷款违约预测	54	87.04%
346	红酒品质的分类	41	78.05%

### 383. 中文手机评论情感倾向判定

时间限制:100s 内存限制:1024MB

#### 题目描述

网购在我们的日常生活中起到了越来越重要的作用，当我们购买某一产品时，往往第一反应是参考相应的评论。每条评论中都有其情感倾向，无论是好评、差评还是中评。我们通过分析评论，有更大的概率做出较优的决策。

本题目提供了一份从京东爬取的某款手机评论数据，每一条评论按照其包含的情绪类别已被标注为好评(1)、差评(-1)或中评(0)，要求构建分类模型，根据用户评论，预测用户评论的情绪类别。

注意：用户评论为原始文本形式，需要用户自行完成中文分词等预处理步骤。

#### 开始答题

Python  
R

显示范例

```
1 class Solution(MLWorker):
2     # 请在下面区域作答 #
3     def train(self, dataframe_trainx, dataframe_trainy=None):
4
5
6     def predictValue(self, model, dataframe_testx):
7
8
9
```

清空

提交

返回列表

# 大数据专业建设的关键因素四

## 大数据教学实训平台 - 数据嗨客

### 评测报告：平台自动分析、评价数据分析与建模效果

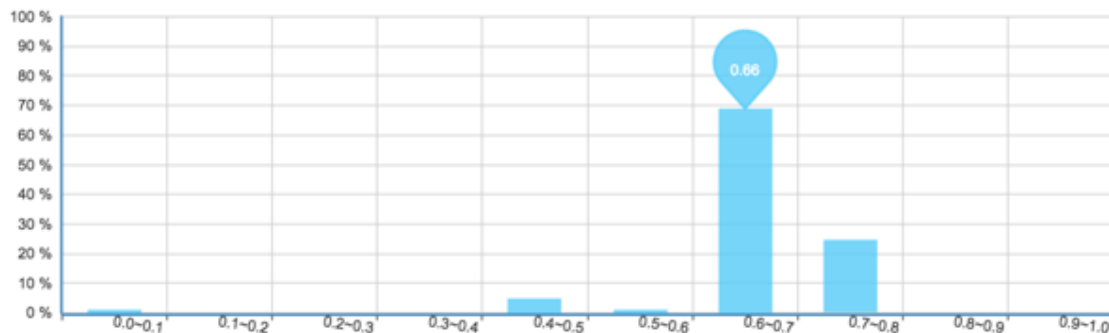
#### 模型评价

准确率 Accuracy	精确度 Precision	召回率 Recall	F1值 F1_score	对数损失 Log_loss	马修斯相关性系数 Matthews_corrcoef	ROC曲线下面积 ROC_AUC
0.66	0.63	0.66	0.64	--	--	--

分类 Classification	精确度 Precision	召回率 Recall	F1值 F1_score
-1	0.71	0.66	0.68
0	0.25	0.16	0.20
1	0.71	0.85	0.77

所有结果分布图

准确率



#### 提交代码

```
3 class Solution(MLWorker):
4
5     def train(self, dataframe_trainx, dataframe_trainy=None):
6
7         from sklearn.tree import DecisionTreeClassifier
8         from sklearn.feature_extraction.text import TfidfVectorizer
9         import jieba
10
11         x = dataframe_trainx[list(dataframe_trainx.columns.values)[0]].map(lambda xx:
12 list_text = list(x)
13 self.tfidf = TfidfVectorizer(max_df=0.95, stop_words=stopwords).fit(list_text)
14 array_trainx = self.tfidf.transform(list_text)
15 array_trainy = dataframe_trainy.values.ravel()
16 # 训练模型
17 model = DecisionTreeClassifier().fit(array_trainx, array_trainy)
```

#### 编译信息

```
1 /tmp/p_sandbox/112300.py:13:80: E501 line too long (93 > 79 characters)
2 /tmp/p_sandbox/112300.py:19:1: E302 expected 2 blank lines, found 0
3 /tmp/p_sandbox/112300.py:26:70: E231 missing whitespace after ':'
4 /tmp/p_sandbox/112300.py:28:5: E301 expected 1 blank line, found 0
5 /tmp/p_sandbox/112300.py:37:1: E302 expected 2 blank lines, found 0
6 /tmp/p_sandbox/112300.py:50:80: E501 line too long (112 > 79 characters)
7 /tmp/p_sandbox/112300.py:52:80: E501 line too long (86 > 79 characters)
8 /tmp/p_sandbox/112300.py:64:80: E501 line too long (110 > 79 characters)
9 /tmp/p_sandbox/112300.py:72:5: E303 too many blank lines (2)
10 /tmp/p_sandbox/112300.py:73:80: E501 line too long (90 > 79 characters)
11 /tmp/p_sandbox/112300.py:75:80: E501 line too long (99 > 79 characters)
12 /tmp/p_sandbox/112300.py:91:4: W292 no newline at end of file
13 0.06 seconds elapsed
14 0 directories per second (0 total)
15 16 files per second (1 total)
```

# 大数据专业建设的关键因素四

## 大数据教学实训平台 - 数据嗨客

数据探索：提供数据集分析的云环境

线上教学：组织管理在线教学

The screenshot shows the Data HackData platform interface. At the top, there is a navigation bar with '学习' (Learn), '练习' (Practice), '教室' (Classroom), and '排名' (Ranking). Below the navigation bar is a menu with 'Edit', 'View', 'Insert', 'Cell', 'Kernel', and 'Help'. The main area is a code editor with the following code:

```
In [49]: data.iloc[0,0] = 29
import matplotlib inline
import seaborn as sns
sns.heatmap(data.corr())

#3 测试集和训练集划分 交叉验证

#4 模型训练
from sklearn import linear_model

#5 模型预测

#6 效果评估
```

Below the code editor is a heatmap visualization showing the correlation matrix for various features. The features listed on the y-axis are: age, salary, education, graduated.years, marriageStatus, hasChild, hasHouse, houseLoan, carLoan, hasCar, workYears, officeScale, months, borrowType, credit, status, gender\_0, and gender\_1. The heatmap shows a strong positive correlation (red) between 'age' and 'salary', and a strong negative correlation (blue) between 'gender\_0' and 'gender\_1'. A color scale on the right ranges from -0.8 (blue) to 0.8 (red).

The screenshot shows the Data HackData platform interface displaying a list of online courses. The navigation bar and menu are the same as in the previous screenshot. The main area shows a grid of course cards. Each card includes a title, a description, the instructor's name, and the number of learners.

- 162 普林大数据人员内部培训**  
普林大数据人员内部培训  
教师: welyunlei  
5人学习
- 161 国家信息中心大数据培训**  
大数据实践与工具, 大数据行业案...  
教师: 欧高炎  
13人学习
- 160 第一期“数据科学与大数...**  
大数据受到国家和企业的高度重视...  
教师: 欧高炎  
16人学习
- 153 数据科学导引**  
《数据科学导引》是北京大学数据...  
教师: 鄂维南  
76人学习
- 144 浙江物联网大数据培训**  
浙江省物联网产业协会联合北京大...  
教师: 欧高炎  
29人学习
- 135 大数据分析的模型与算法**  
主要演讲人: 鄂维南 (北京大数据...  
教师: 欧高炎  
74人学习



# 大数据专业建设的关键因素四

## 大数据教学实训平台 - 数据嗨客

作业评测：自动生成班级作业测评报表，辅助教学

### 文本分析实战作业

截止时间：2016/12/12 23:00

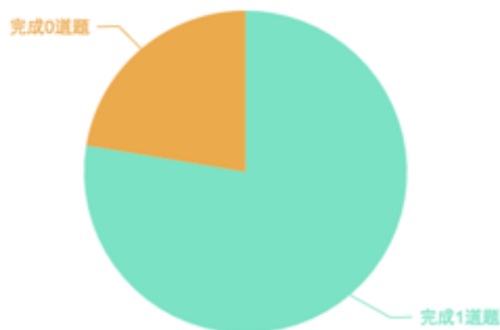
编辑

最晚时间：2016/12/12 23:00

作业要求：本次作业给定一份中文手机评论数据，要求构建文本分类模型。评价指标为F1值。给分标准为：未通过 0分  $F1 < 0.5$  1分  $0.5 \leq F1 < 0.6$  2分  $0.6 \leq F1 < 0.7$  3分  $0.7 \leq F1 < 0.8$  4分  $0.8 \leq F1$  5分

### 作业完成情况

完成1道题  
完成0道题



### 文本分析实战作业

题目：中文手机评论情感倾向判定

用户	提交情况	提交次数	最后提交时间	状态	查看代码
wenzaiwen	未提交	0	-	成功0次	查看
jingweiw	未提交	0	-	成功0次	查看
龚佃选	未提交	0	-	成功0次	查看
王希舜	未提交	0	-	成功0次	查看
董彬	未提交	0	-	成功0次	查看
刘艳云	未提交	0	-	成功0次	查看
余冰	按时	39	2016-12-11 21:51	成功33次	查看
孟少帅	未提交	0	-	成功0次	查看
陈潇漪	未提交	0	-	成功0次	查看
任惠霞	按时	11	2016-11-20 12:52	成功1次	查看
温见培	按时	5	2016-12-07 15:14	成功0次	查看
梁泽	按时	12	2016-12-09 01:58	成功4次	查看
赵星楠	按时	26	2016-12-10 17:18	成功10次	查看
陈龙	按时	25	2016-12-11 21:55	成功23次	查看

# 大数据专业建设的关键因素四

## 大数据教学实训平台 - 数据嗨客

行业竞赛：帮助企业解决大数据的实际问题

The screenshot shows the '数据嗨客大数据竞赛' (DataHiKe Big Data Competition) website. The header includes navigation tabs: 练习 (Practice), 学习 (Learn), 教室 (Classroom), 排名 (Ranking), and 竞赛 (Competition). The main banner features the title '数据嗨客大数据竞赛' and a sub-header '打造国际高端算法竞赛，通过开放数据源和计算资源，让选手用算法解决社会或业务问题，数据引领新生代力量！' Below the banner, there is a list of five competitions:











竞赛状态	竞赛名称	进行中	报名参赛
进行中	第五期数据嗨客大数据竞赛	进行中	报名参赛
进行中	第四期数据嗨客大数据竞赛	进行中	报名参赛
进行中	第三期数据嗨客大数据竞赛	进行中	报名参赛
已结束	第二期数据嗨客大数据竞赛	已结束	报名参赛
已结束	第一期数据嗨客大数据竞赛	已结束	报名参赛



# 大数据专业建设的关键因素四

## 大数据教学实训平台 - 数据嗨客

能力鉴定：基于大数据的人才能力评价

排名	用户	做题数	提交次数	通过率
1	 welyunlei	204	2555	89.82%
2	 姜晓东	130	717	52.16%
3	 马艳艳	111	1164	68.27%
4	 欧奕炎	88	387	74.16%
5	 杨嘉琪	72	178	82.58%
6	 jiekang	82	205	80.98%
7	 高佳佳	70	290	80.34%
8	 李雪鹏	47	62	82.26%
9	 张兆强	47	146	87.67%
10	 马璇	35	84	64.29%

1 2 3 4 5 6 下一页 尾页 257条 26页

- 讲义和案例的浏览记录
- 练习题目的效果排名
- 参与竞赛的排名情况
- 语言能力鉴定（Python和R等）
- 数据分析历史提交代码
- 题目和竞赛完成数量



大数据能力鉴定报告

# 大数据专业建设的关键因素五

## 大数据教育体制内外合作



为了给大数据专业人才教育提供支持，北京大数据研究院、博雅大数据学院联合新华三集团，并以大数据教育联盟为平台共同推出大数据人才培养计划，主要包含以下三方面内容：

1

大数据学科建设咨询服务：为全国院校提供《数据科学与大数据技术》、《大数据技术与应用》学科的人才培养体系建设、课程体系设计、师资培养、教材课件开发等专业服务。

2

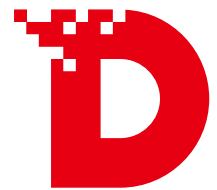
大数据实训环境搭建：为培养学生实践及创新能力，在传统线下教学的基础上，增加大数据教育线上实训平台-数据嗨客，以及大数据实验室，使学生上课时可进行结合真实案例的学习与练习、可操作真实设备等。在课程体系规划中，将实训部分作为重点，指导各院校建立大数据创新应用实训中心，搭建真实完整的机房环境，提供软硬件一体的实训加科研解决方案，根据真实项目实施流程开发实训课程，让本科及研究生具备创新开发能力、高职类学生具备企业级工程师的素养。

3

共同推出国家权威认证，为众多学员进入大数据领域提供了更有分量的筹码。

# 大数据专业建设的关键因素五

## 大数据教育体制内外合作



Powered by DataEngine



### 第三篇 大数据技术实战演练

- 大数据平台组件操作演练
- 使用真实数据模拟演练
- 应用开发实战

### 第二篇 大数据关键技术

- 大规模并行计算技术介绍
- Hadoop生态介绍
- 常用数据挖掘

### 第一篇 大数据基础知识

- 大数据产生背景
- 大数据技术介绍
- 大数据发展前景



**大数据教育联盟**  
Big Data Education Alliance

大数据教育联盟是北京大数据研究院、北京大学元培学院发起并联合众多国内外高等院校、政府主管部门、科研机构和企业事业单位共同组建的非盈利性组织。联盟理事长由鄂维南院士担任。

联盟致力于贯彻和落实国家大数据战略，开展大数据专业人才认证标准研究、数据科学体系建设与推广、学术与人才交流沟通等工作，切实推进大数据“产学研用”的无缝结合，促进商学互动，为国家大数据战略创造良好的生态体系，为大数据产业输送专业的人才。

**理事长：**鄂维南院士

**副理事长单位：**北大、清华、人大、复旦、中山大学、中南大学、贵州大学、对外经济贸易大学、工信部、教育部、新华三、微软、京东等为副理事长单位，理事单位包括百余所高校机构和企业。

# 大数据教育联盟成立仪式 暨大数据人才培养高峰论坛



为了推动大数据行业的发展，帮助更多高校开展大数据学科的建设，联盟定于**2017年5月23日**上午9点在**北京大学**英杰交流中心阳光厅举办“大数据教育联盟成立仪式暨大数据人才培养高峰论坛”。活动首先将与所有参会嘉宾共同见证大数据教育联盟揭牌成立的重要历史时刻，随后的大数据人才培养高峰论坛将由复旦大学、中国人民大学、中南大学、贵州大学、北京信息科技大学等成功申报“数据科学与大数据技术”学科的院校代表为大家分享大数据学科建设及人才培养方案的经验。



**大数据教育联盟**  
Big Data Education Alliance

主办单位：北京大数据研究院，北京大学元培学院  
承办单位：博雅大数据学院  
时间：2017年5月23日 上午9点  
地点：北京大学英杰交流中心阳光厅



# 大数据教育联盟成立仪式 暨大数据人才培养高峰论坛



**大数据教育联盟**  
Big Data Education Alliance



## 议程安排

时间	议题	演讲嘉宾
8:30-9:00	签到	
9:00-9:05	主持人开场致辞	王涛 联盟副秘书长
9:05-9:15	联盟理事长致辞	鄂维南院士 北京大数据研究院院长
9:15-9:40	联盟副理事长单位、常务理事代表致辞	北京大学、清华大学、对外经贸大学、中山大学、微软加速器、新华三集团、京东金融等
9:40-10:00	成立仪式	
大数据人才培养高峰论坛		
10:00-10:20	复旦大学 大数据人才培养模式分享	高卫国 大数据学院院长助理
10:20-10:40	中国人民大学 大数据人才培养模式分享	赵彦云 统计学院院长
10:40-11:00	中南大学 大数据人才培养经验介绍	邹北骥 信息科学与工程学院院长
11:00-11:20	贵州大学 大数据人才培养工作介绍	谢泉 大数据与信息工程学院院长
11:20-11:40	北京信息科技大学 大数据人才培养工作介绍	王兴芬 教务处处长
11:40-12:00	大数据教育平台建设规划	欧高炎 联盟秘书长





北京大数据研究院公众号



大数据教育联盟公众号



THANK YOU!