# 《数据科学R与Pyhon实践》

顾立平

# 来自2012年的"最性感"的未来职业

# 教学与实践

# 如今就业要求...

**难以企及和数据科学家**

**技术**

JMP, SAS, R, Python, Perl, Excel, SQL, Hadoop, Java, JavaScript, loT, Event-Stream

**25%**

**25% 方法论**

数学和统计理论、定量和定性方法、试验（例如A/B测试）、六西格玛质量管理方法、质量提升

**领域**

专业知识包括：工作流、运营和可操作分析等

**25%**

**25% 软技能**

影响力，批判性思维、沟通交流能力，系统性思维，设计思维，可视化思维，设计

# 第四科学（理论、实验、计算、数据）

|  | 演绎（理论、原理的预测） | 归纳（数据的解释与探索） |
|---|---|---|
| 人类社会 | **理论科学**<br>理论构建<br>脑力竞争<br> | **实验科学**<br>经验解释<br>团队组织<br> |
| 数字空间 | **计算科学**<br>自动化<br>博弈胜负<br> | **数据科学**<br>人工智能<br>数据质量<br> |

数据科学R与Python实践所需技能

| | |
|---|---|
| | **数据科学**因为跨学科性质，强调交叉能力：**黑客技能、数学和统计知识**以及学科专业的**领域知识**。 |
| | **黑客技能**有助于处理那些必须获取、清理和操纵的大量电子数据。 |
| | **数学和统计知识**有助于数据科学家选择适当的方法和工具，以便从数据中获得独到见解。 |
| | **领域知识**有助于探索有意义的研究问题及其假设建立，以及验证和解释随后的实验结果。 |
| | **传承的研究方式**有助于在数学和统计知识与领域知识之间进行新的方法遴选、使用和评估。 |
| | **机器学习的方式**有助于结合黑客技能与数学和统计知识，但是缺乏探索有意义的科学问题的科学动机。 |
| | **危险禁区！** 如果没有严谨的科学精神和方法，黑客技能与领域知识相结合，很可能会导致不正确的分析。 |

# 谢谢您

顾立平

邮箱gulp@mail.las.ac.cn

手机/微信15801137067