

机器学习在企业级场景中的实践与探讨

Enterprise Machine Learning Practices

Jerry Zhu 朱辉
IBM机器学习全球研发总经理
IBM中国开发中心认知计算实验室总经理

大多公司在这儿

数据驱动

洞察驱动

数字化创新

结果

文化转型
打破信息孤岛
发现：“是什么？”
理解：“为什么？”

预测
优化
自动化
协作

新业务模型
颠覆技术
实时决策

能力

自服务
报表
商业智能

模型
可视化
应用

嵌入
排程
集成

驱动因素

节省成本、
现代化

竞争

市场领导力

数据价值

提升数据价值的关键因素

访问数据



不同来源和结构的数据
获取和存储、提供统一
访问接口

治理数据



组织、管控数据以提
升数据质量，安全性
和可信性

获取洞察



从数据中快速准确的学习和
获取洞察与智能，帮助业务
转型

访问

治理

分析

混合数据管理

统一的数据管控和集成

数据科学和认知商业



Write Once, Access Anywhere

with a common access layer to promote application independence



Prepare, Publish, Protect

your data to drive insights while mitigating compliance risks



Describe, Predict, Prescribe

to understand the current, predict the future and change the outcome

融入 **机器学习**
无缝集成 **云上和云下**
基于 **开源技术**

人工智能

像人类一样的推理能力

- Watson
- 无人驾驶
- 围棋比赛
- 更多...

机器学习

不需要显示编程的学习能力

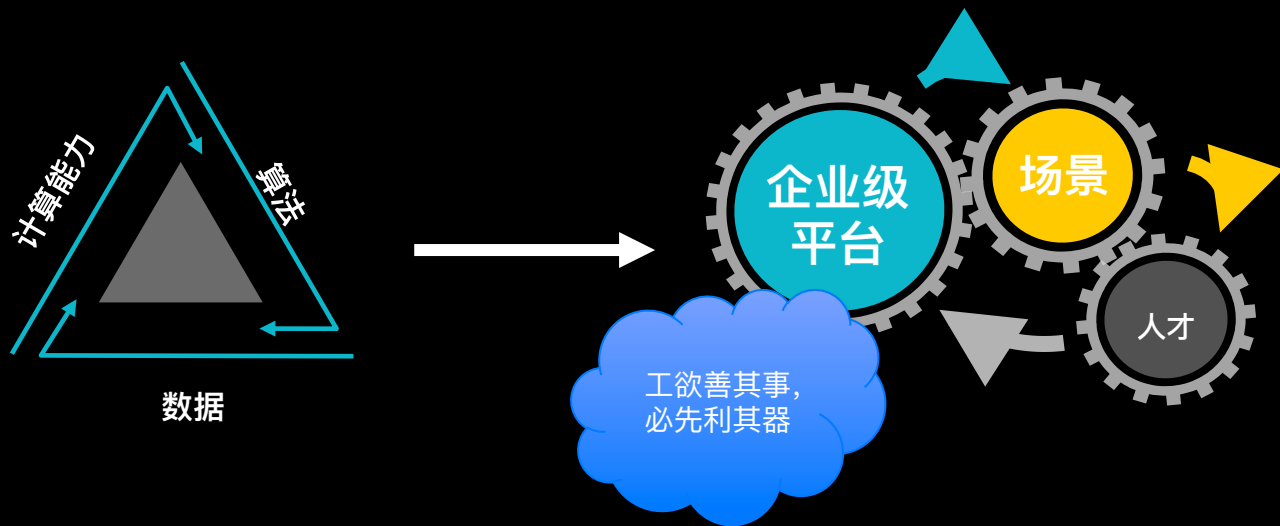
- 欺诈检测
- 推荐系统
- IT自动化运维
- 更多...

深度学习

自主学习新数据集的能力

- 智能健康
- 虚拟客服
- 机器翻译
- 更多...

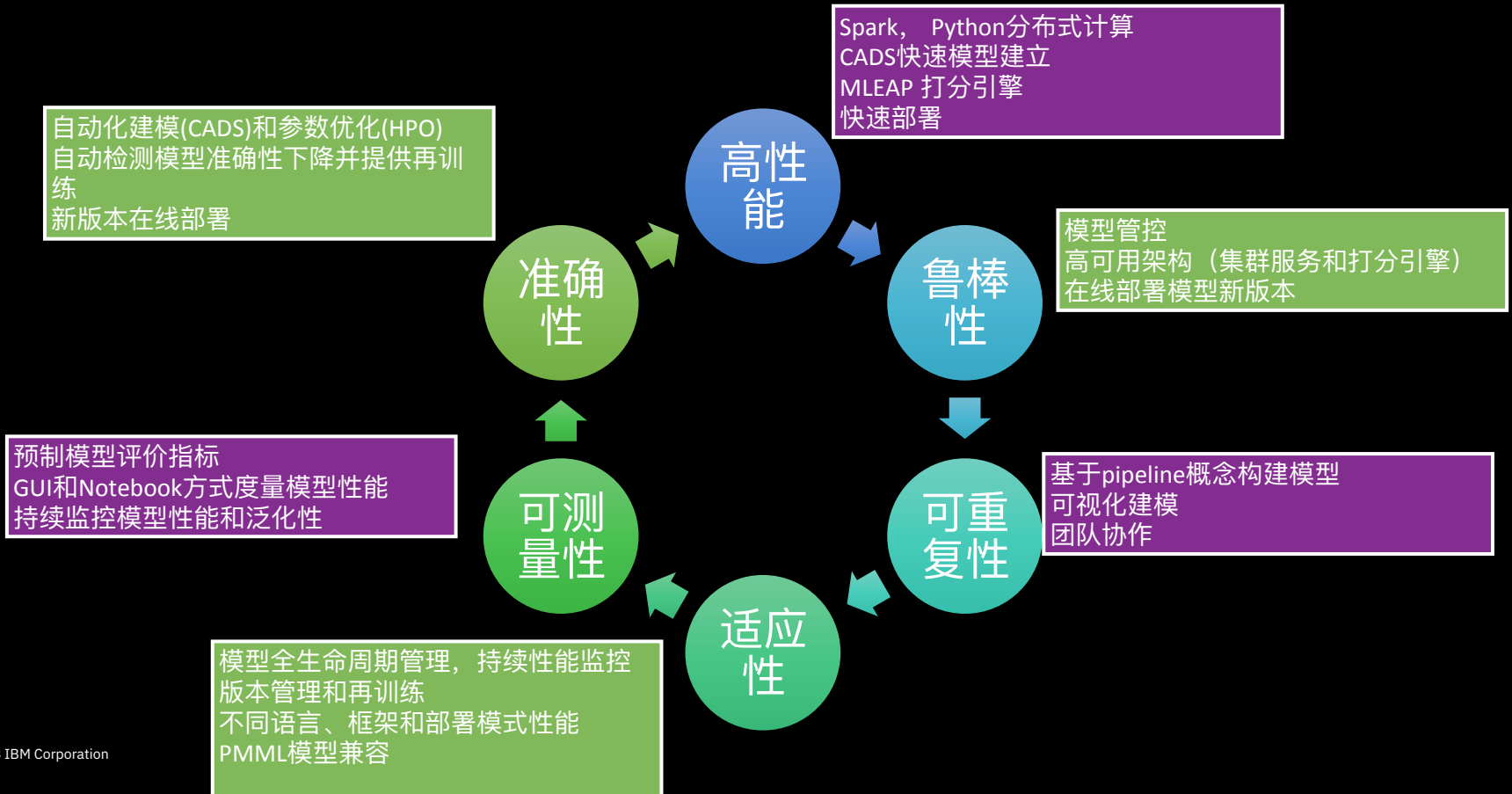
企业级机器学习的工程化挑战



常规机器学习流程及挑战

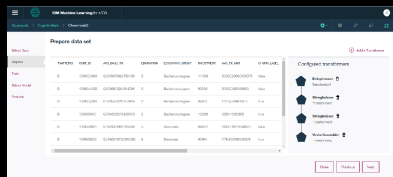


企业级机器学习考量点 – 模型可操作性 (Model Operationalization)

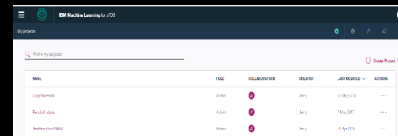


机器学习全生命周期管理与协作

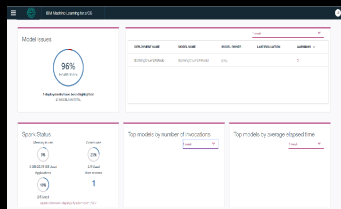
通过一个监控仪表盘管理企业内部所有部署的模型



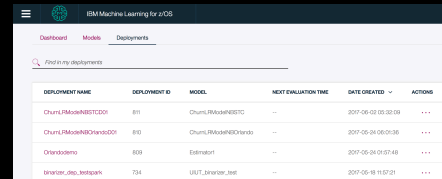
业务分析师



数据科学家



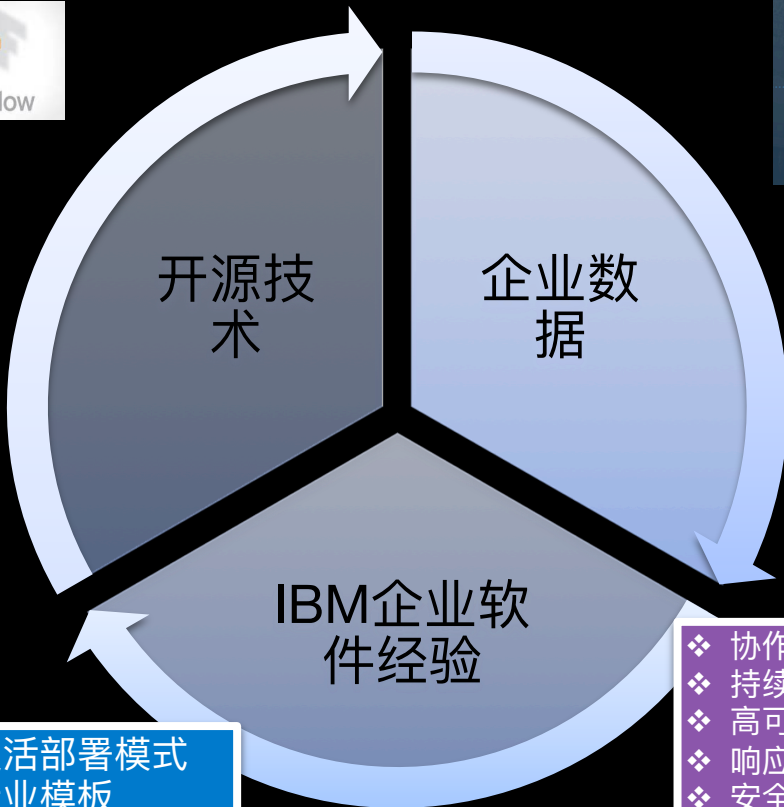
应用开发工程师



机器学习工程师

IBM 企业级机器学习平台

结合最新开源软件、企业数据和IBM特有的企业软件经验



全球

90%

以上的数据无法通过互联网获取



统一治理

覆盖所有数据的更佳洞察及合规



混合数据管理

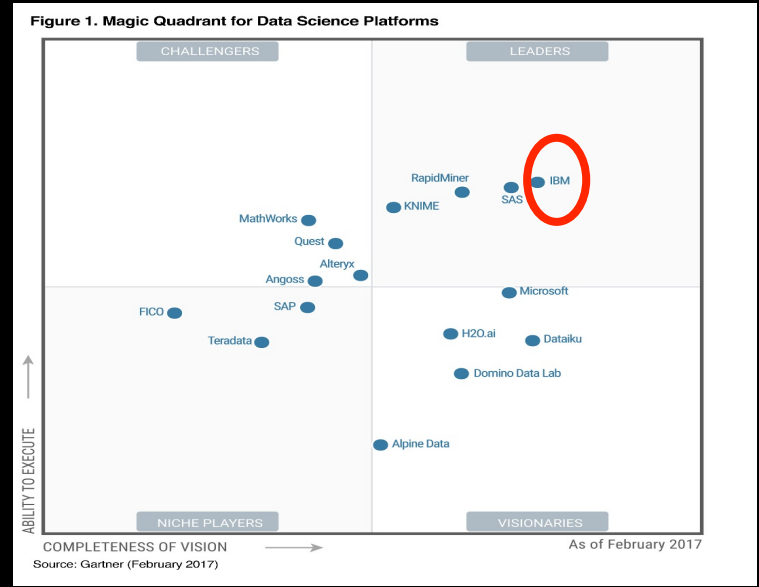
统一数据及内容上云之路

- ❖ 灵活部署模式
- ❖ 行业模板

- ❖ 协作和生命周期管理
- ❖ 持续监控和模型优化
- ❖ 高可用性
- ❖ 响应时间，高吞吐
- ❖ 安全性

IBM 数据科学平台: Gartner 和 Forrester 报告双料冠军

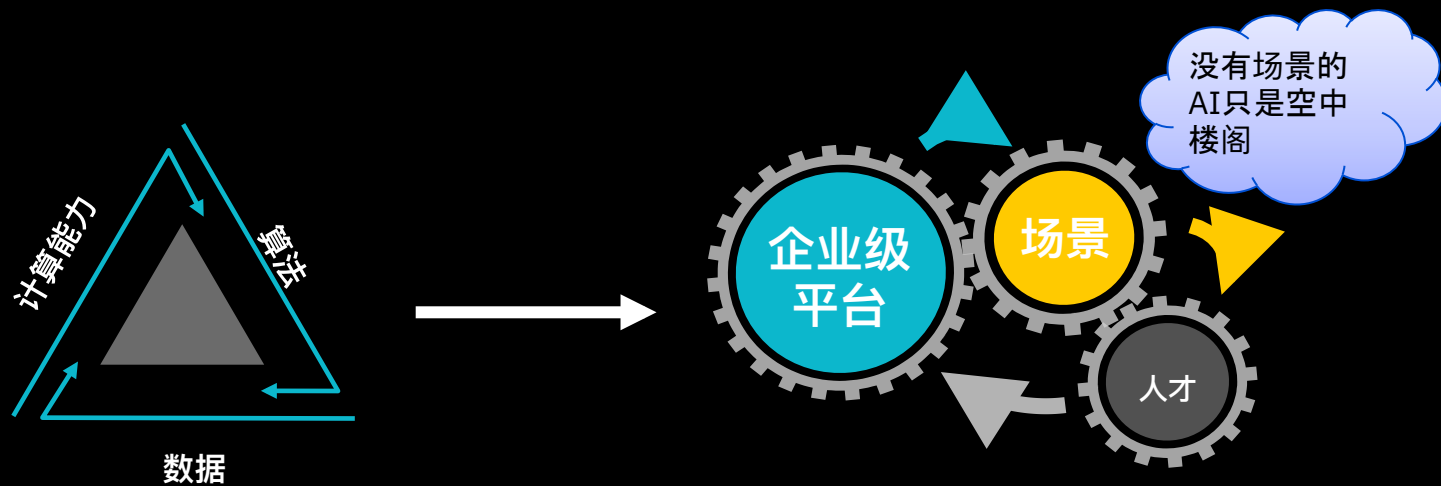
IBM被2017 Gartner 数据科学魔力象
限评为最具远见和执行力的领导地位



The Forrester Wave™: Predictive Analytics
And Machine Learning Solutions, Q1 2017
© 2016 IBM Corporation

Forrester: IBM拥抱开源, 发布的Data Science
Experience 云服务能快速创建基于Spark集群的开源
Jupyter 和 Rstudio notebooks

企业级机器学习的工程化挑战



人工智能助力交通管控 (1)

短时交通预测

长时交通预测

拥堵扩散规律

交通异常检测

拥堵趋势预测

拥堵模式分析

交通预测分析算法及模型

机器学习平台

(数据探索、特征工程、模型选择、自动化建模、一键部署、持续学习)

IBM混合数据管理平台

手机信令

出租GPS

车联网

轨道卡

天气数据

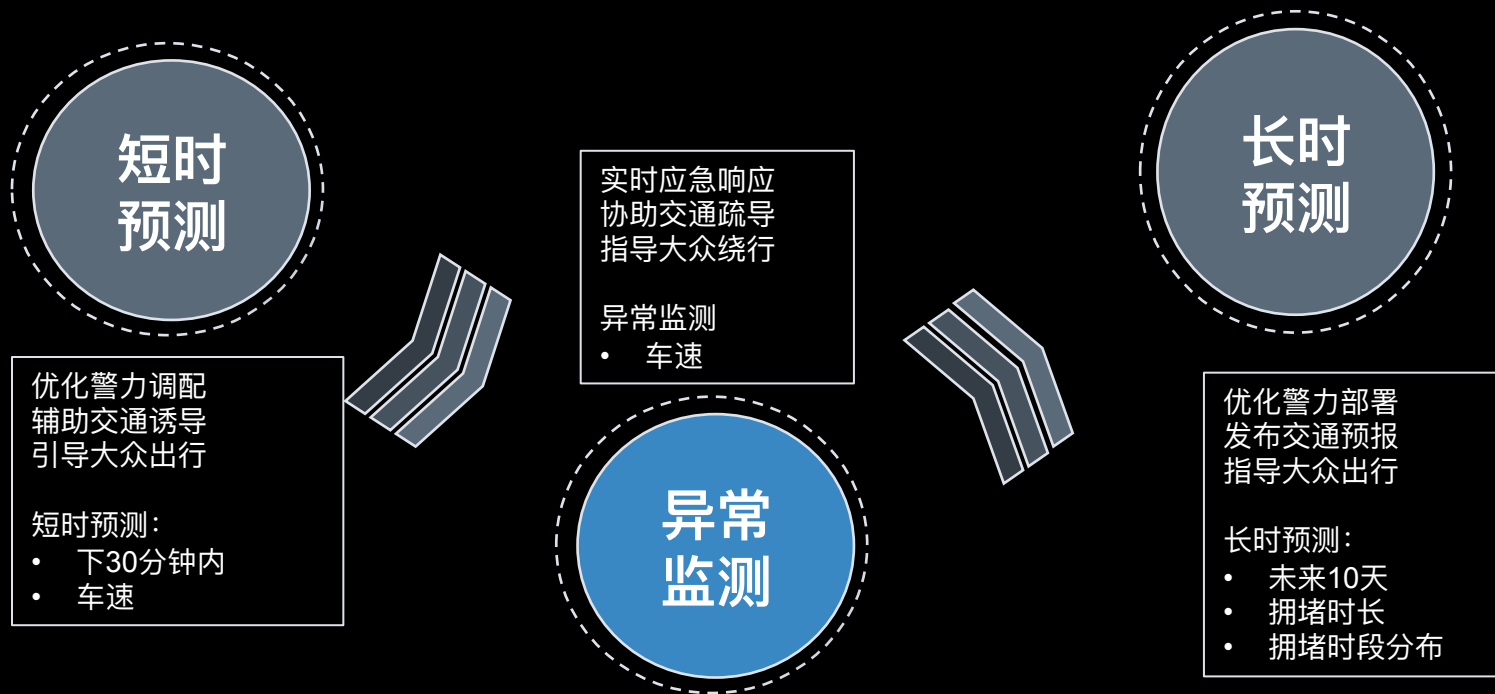
路网规划

交通事件

RFID/卡口

基础硬件环境

人工智能助力交通管控 (2)



人工智能助力交通管控 (3)

[异常监控预警](#)

[实时路况预报](#)

[7日路况预报](#)

[常发拥堵查询](#)

[拥堵趋势查询](#)

[拥堵对比查询](#)

[拥堵事件查询](#)

[拥堵模式查询](#)

[路网堵点查询](#)



人工智能助力交通管控

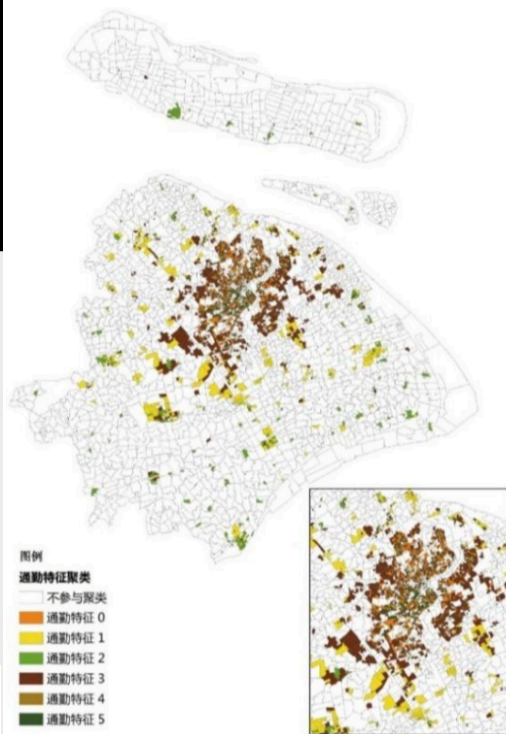
人工智能助力交通管控 (4)



机器学习辅助城市规划

- 融合来自政府和企业的普查、交通、通信、房价等数百项数据，构造多维城市运行指标。
- 利用机器学习算法，对城市进行多角度的现状评估、风险预警，决策模拟。
- 应用1:
- 根据城市地块的就业类、居住类几十项定性指标，研究各类地块在城市中的空间分布。
- 应用2:
- 从就业、居住功能角度分类描述特定地块。
- 应用3:
- 模拟政策对地块指标带来的数量变化，并预测政策是否会使地块发生质的变化。

- 通勤特征0：中等距离少量通勤波动型
- 通勤特征1：长距离极少量通勤型
- 通勤特征2：短距离极少量通勤型
- 通勤特征3：中等距离极少量通勤型
- 通勤特征4：中等距离大量通勤型
- 通勤特征5：中等距离中等量通勤波动型



基于通勤特征的居住空间
聚类结果空间分布

IBM 机器学习平台提供一个数据科学家的协作平台，能帮助我们快速探索数据并创建模型，并将高效模型服务到同衡的客户。

-- 李栋博士 清华同衡创新研究院 常务副主任

机器学习预测租房价格

• 业务场景

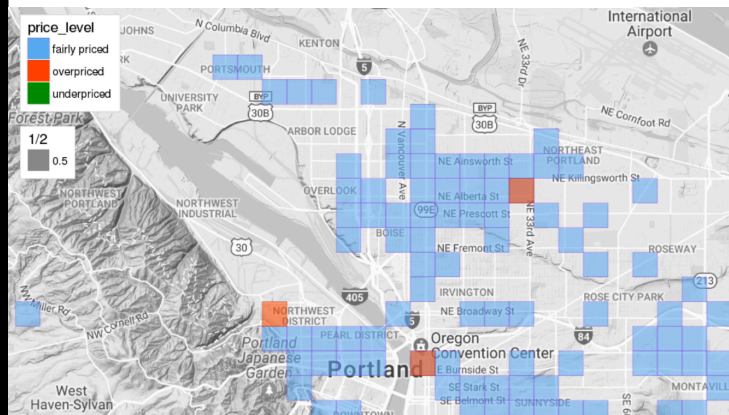
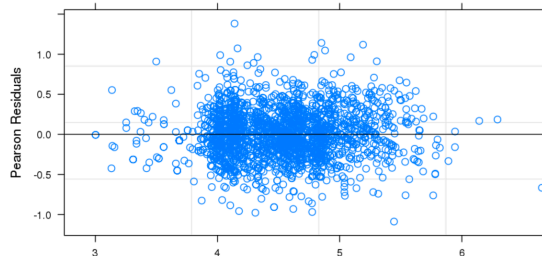
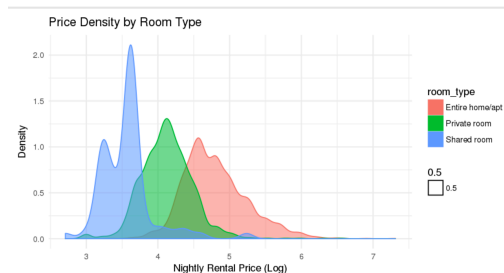
- 去到一个新的城市，怎么从Airbnb上选择性价比高的房子
- 出租房屋时，如何才能定合适的价格呢？

• 机器学习方法

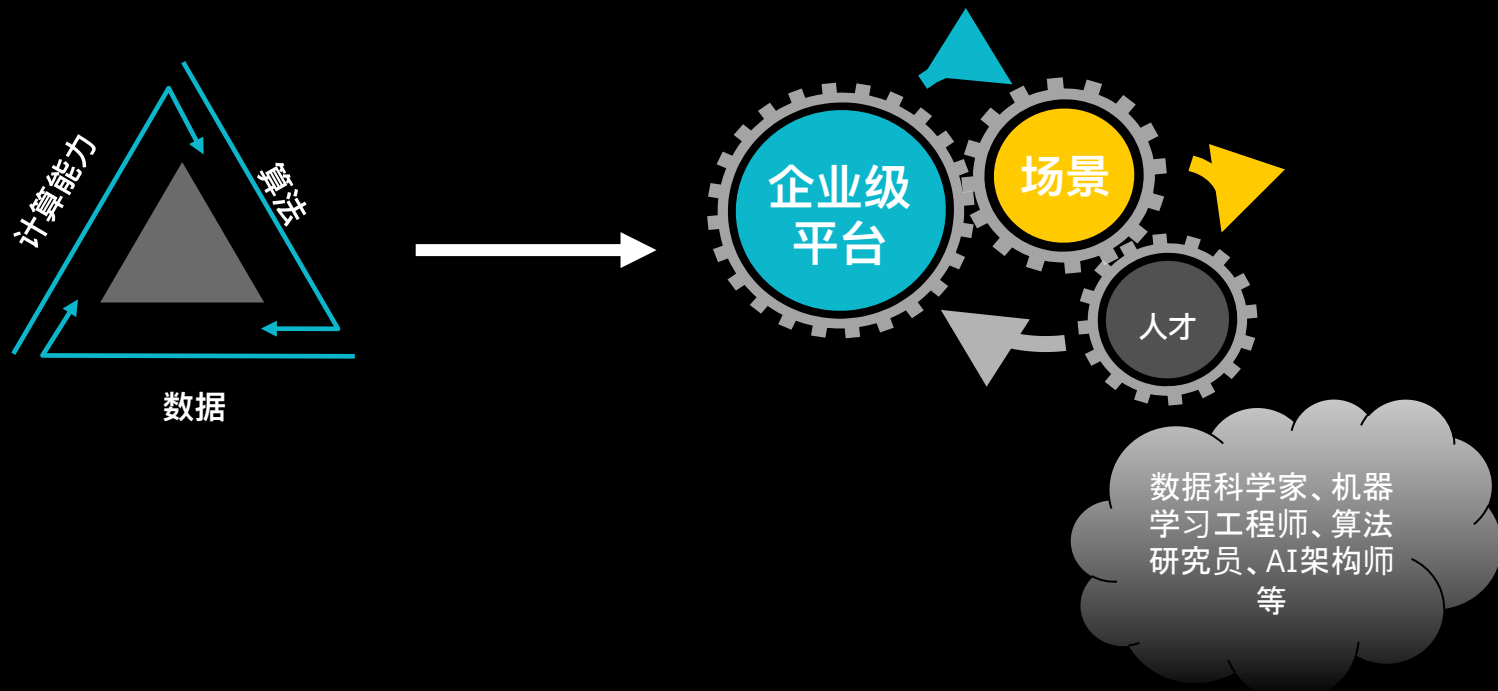
- 数据获取-- 从Airbnb 上抓取出租房屋信息和价格
- 数据清洗- 原始数据来自于互联网，需要进行异常值处理、缺失值填充等
- 数据探索- 基于区域、房屋类型、价格、可住人数等进行分析
- 特征工程- 数据分组、分箱、特征选择
- 模型选择和超参数调优 - RF, Xgboost, Mixed-effects Model 等

• 业务价值

- 分析发现被低估价值的出租房
- 在Airbnb上找到性价比高的租房
- 为自己的房子出租前定合适的价格



企业级机器学习的工程化挑战



2

IBM机器学习中心

IBM Machine Learning Hub

学习 创新 应用
Learn Innovate Apply

- ❖ 你的数据+场景+数据科学家 AND
- ❖ IBM数据科学家+IBM机器学习平台

Community

Notebooks

| | | | |
|---|--|---|--|
| <p>NOTEBOOK Learn basics about notebooks and Spark</p> <p>AUTHOR DATE IBM May 30, 2016</p> <p>TOPIC SOURCE Enrollment External</p> | <p>NOTEBOOK Analyze open data sets with DataFrames</p> <p>AUTHOR DATE IBM Aug 29, 2016</p> <p>TOPIC SOURCE Society External</p> | <p>NOTEBOOK Analyze energy consumption in buildings</p> <p>AUTHOR DATE IBM Sep 22, 2016</p> <p>TOPIC SOURCE Science & Technology External</p> | <p>NOTEBOOK Visualize car data with Brunel</p> <p>AUTHOR DATE IBM Sep 12, 2016</p> <p>TOPIC SOURCE Transportation Self-Contained</p> |
| <p>NOTEBOOK Use R to load data and run SQL queries</p> <p>AUTHOR DATE IBM May 15, 2016</p> <p>TOPIC SOURCE Transportation Self-Contained</p> | <p>NOTEBOOK Use Python to load data and run SQL queries</p> <p>AUTHOR DATE IBM May 15, 2016</p> <p>TOPIC SOURCE Transportation Self-Contained</p> | <p>NOTEBOOK Maximize oil company profits by using decision optimization</p> <p>AUTHOR DATE IBM-DO Jun 02, 2016</p> <p>TOPIC SOURCE Economy & Business Self-Contained</p> | <p>NOTEBOOK Use DO to help a sports league schedule its games</p> <p>AUTHOR DATE IBM-DO May 15, 2016</p> <p>TOPIC SOURCE Leisure Self-Contained</p> |

3

产品自带案例
和应用

1



学习路径 课程 徽章 活动 博客

数据科学基础

当蝴蝶展她的翅膀，会发生什么？他会飞到别的别处的花
造新秩序！

认知学堂

完成向**AI时代**的转型不仅仅
涉及到技术...

...还有同样重要的**思想**
和文化



IBM