



数据挖掘与商务分析

课程导论

主讲教师：肖升生

xiao.shengsheng@shufe.edu.cn



关于授课教师

■ Instructor: 肖升生

- Email: xiao.shengsheng@shufe.edu.cn
- Tel: 021-65904410-837
- Office: #837

■ Research areas:

- (1) BA & Data Mining
- (2) Digital Economics

■ Faculty website:

<https://de.sufe.edu.cn/18/4c/c12089a202828/page.htm>



讲授提纲

- 01 数据类型与价值使用**
- 02 数据挖掘、AI大模型与商务智能**
- 03 跨行业的数据挖掘流程**
- 04 课程内容与设计**
- 05 课程学习材料**



讲授提纲

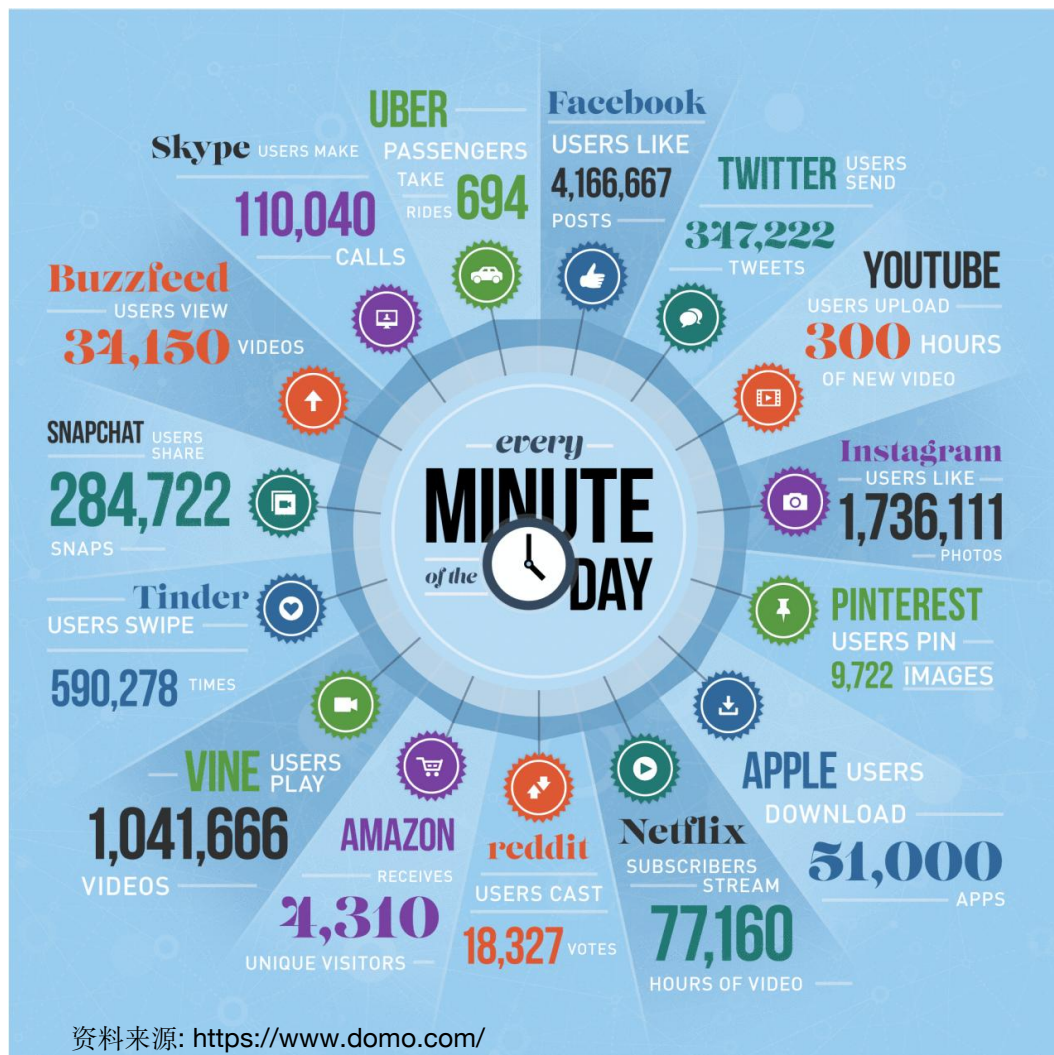
- 01** 数据类型与价值使用
- 02** 数据挖掘、AI大模型与商务智能
- 03** 跨行业的数据挖掘流程
- 04** 课程内容与设计
- 05** 课程学习材料



数据类型与量级

■ 数据类型:

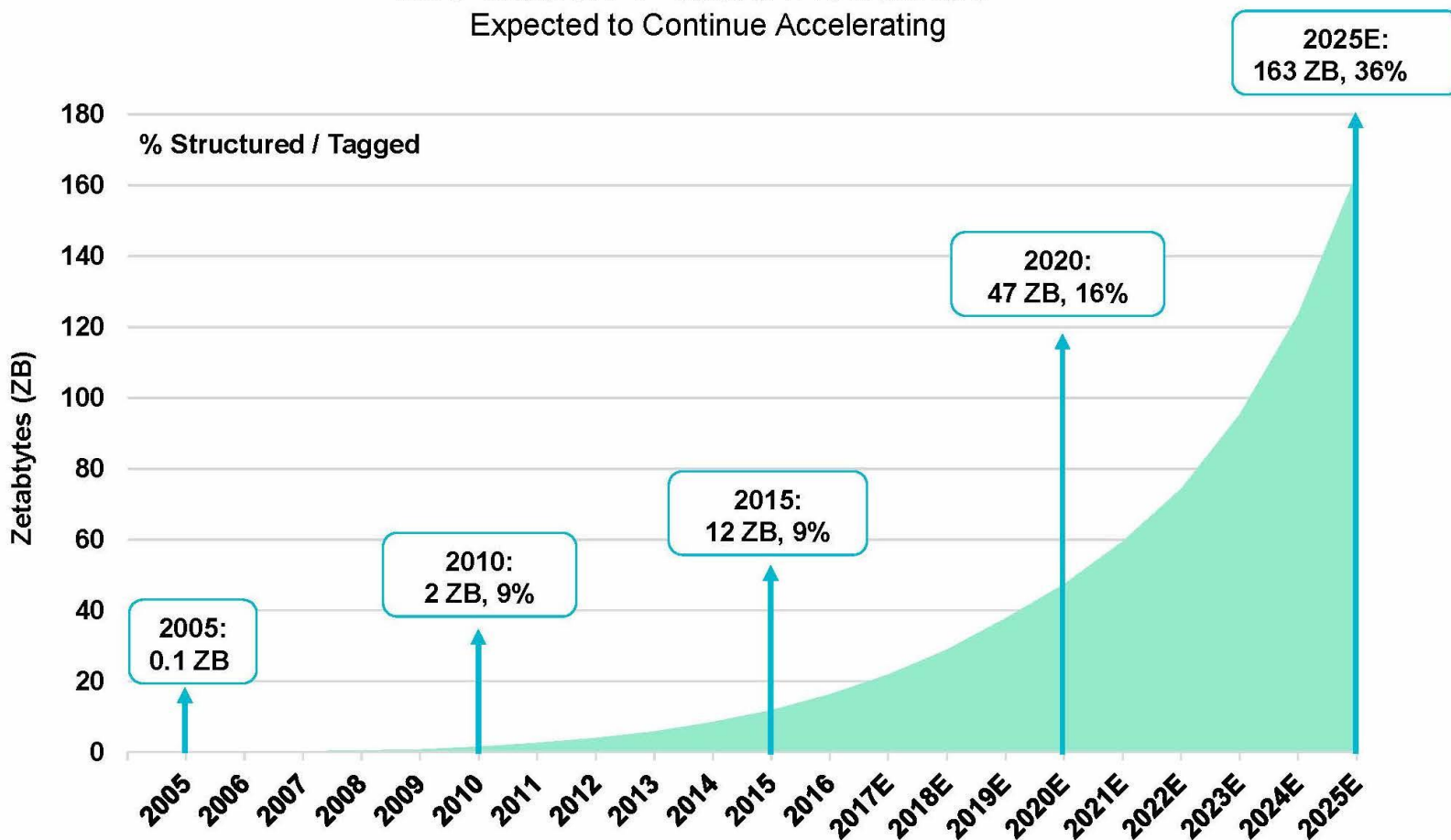
- 数值
- 文本
- 位置
- 声音
- 视频
- ...





数据的快速增长

Information Created Worldwide =
Expected to Continue Accelerating



Source: IDC DataAge 2025 Study, sponsored by Seagate (3/17)
Note: 1 petabyte = 1MM gigabytes, 1 zeta byte = 1MM petabytes

Bit、Byte、KB、MB、GB、TB、PB、EB、ZB、YB



数据的利用率低

- 数据被称为了新“石油”和新资产
- 但跟石油类似，数据的价值需要提炼
- 现状：“Data Rich but Information Poor”
 - 大量的信息隐藏在海量的数据背后
 - 绝大部分的数据都没有被分析和使用
 - 有用信息的挖掘需要耗费大量人力和物力



讲授提纲

- 01** 数据类型与价值使用
- 02** 数据挖掘、AI大模型与商务智能
- 03** 跨行业的数据挖掘流程
- 04** 课程内容与设计
- 05** 课程学习材料



什么是商务智能

商务智能是利用数字智能技术从大量数据中提取信息，转化为可指导决策的知识和洞察力。



零售业库存优化

某大型零售商利用BI分析销售数据，精准预测需求，减少库存积压，提高资金周转率。



银行业风险控制

银行通过BI实时监控交易模式，有效识别欺诈行为，降低信贷风险，提升客户信任度。



制造业生产效率

制造企业运用BI分析生产线数据，优化资源配置，减少浪费，显著提升生产效率和产品质量。



商务智能的几个发展阶段

商业智能初期阶段 (2005-2013)

可视化数据分析阶段 (2013-2016)

浅层决策智能阶段 (2016-2018)

多维决策智能阶段 (2018-至今)

传统商业智能

新型商业智能

企业上线适应自身业务的应用系统，类似于ERP、CRM、OA、HIS等。市场仍主要被SAP、Oracle、IBM等老牌巨头占领。主要用户群体集中于大型企业，且相对封闭。



可视化数据分析产品出现，企业项目中原有的商业智能初期产品逐步下线，此消彼长，可视化的数据分析产品集中进入市场，国内外厂商处于快速成长期。同时，随着IT基础设施逐步完善，更多企业用户拥抱商业智能。



大数据、人工智能技术的发展支持商业智能进入数据挖掘的浅层决策智能阶段。云服务的普及推广支持商业智能解决方案的云端部署，吸引更多的中小企业用户。行业进入新型商业智能阶段。



伴随2018年人工智能技术的全面商业化落地，集合AI、大数据、云服务、RPA、运筹学等技术的新型商业智能开始为企业客户提供多维决策的智能服务。融合技术、打磨场景、优化解决方案的部署成本是现阶段商业智能企业的发展重点。



商务智能与行业变革：汽车行业

机械化汽车

信息技术向汽车设计、生产制造等环节渗透
提高了生产效率

汽车雏形 → 单件少量生产 → 大规模生产

美国
T型车+流水线
1903-1927



美国、欧洲
自动化生产线+精益生产
1947-1980s



英国、法国、美国
蒸汽汽车
1705-1834

德国
内燃机四轮车
1876-1886

日本
多样化+准时化+精益生产
1970-1976

来源：阿里研究院

机电化汽车

数字控制和互联网技术向汽车产品和服务渗透
提升了汽车性能和舒适度，创造更高产品价值

机电一体 → 车联网

电子控制式喇叭、微处理器控制的
ABS/ESP/安全气囊
1970-1982



微电脑控制的
车辆集中电控、
GPS定位/离线导航/移动出行服
务
1982-2000s

车载无线电对讲
1990s

车机互联
Carplay、Android Auto
2013以来

智能化汽车

云计算、人工智能技术与汽车产业深度融合
重新定义汽车，重塑汽车产业

云网端AI一体

自动驾驶：自动泊车、
智能巡航、自动驾驶
2018以来

智能座舱：智能交互、
智能仪表、360°影像
2018



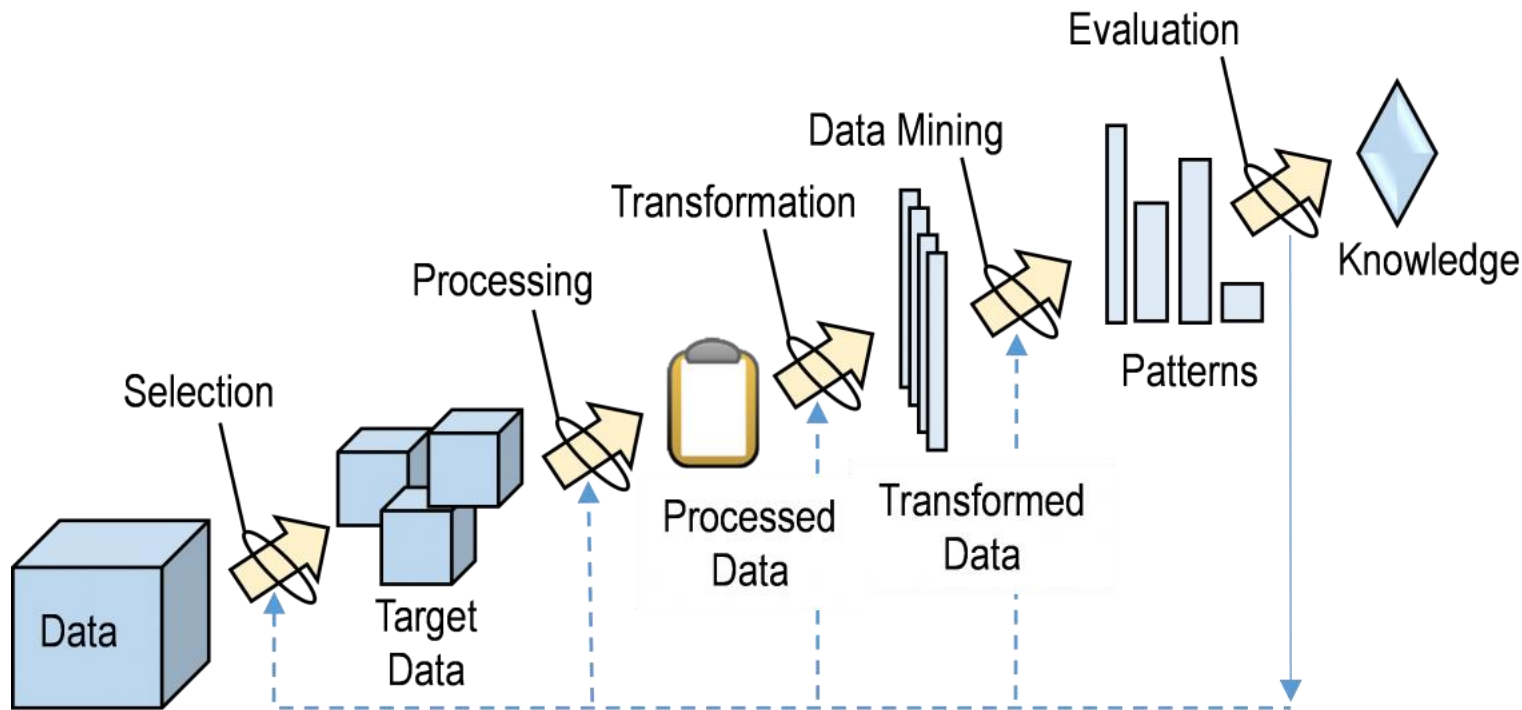
订阅式软件服务：信
息娱乐、系统升级
2020以来

智能导航：路况实时交
互、路线动态优化
2020以来



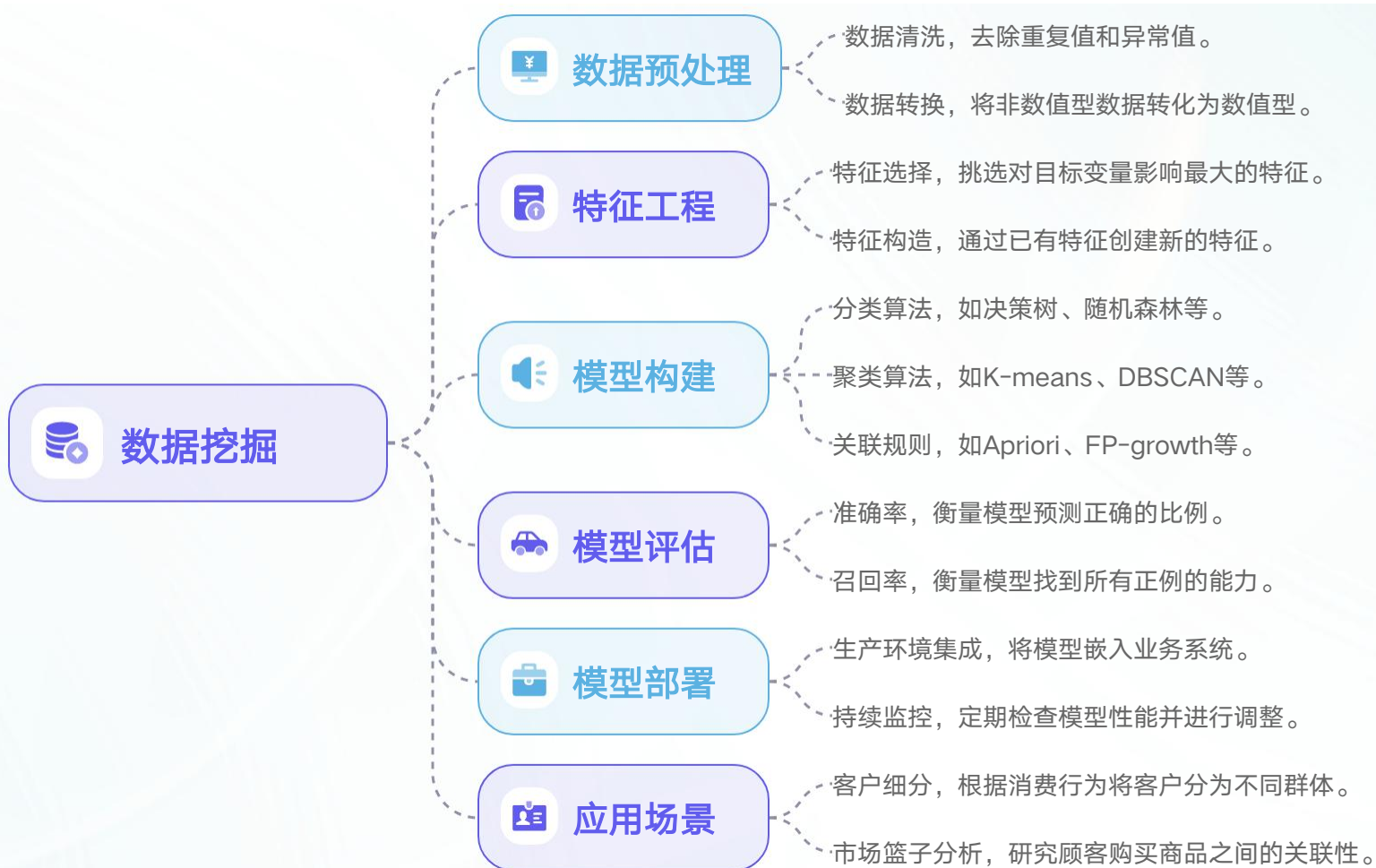
什么是数据挖掘

- 数据挖掘(Data Mining): 从大量的数据中使用智能化的方法自动地发现有用信息的过程





数据挖掘的关键环节





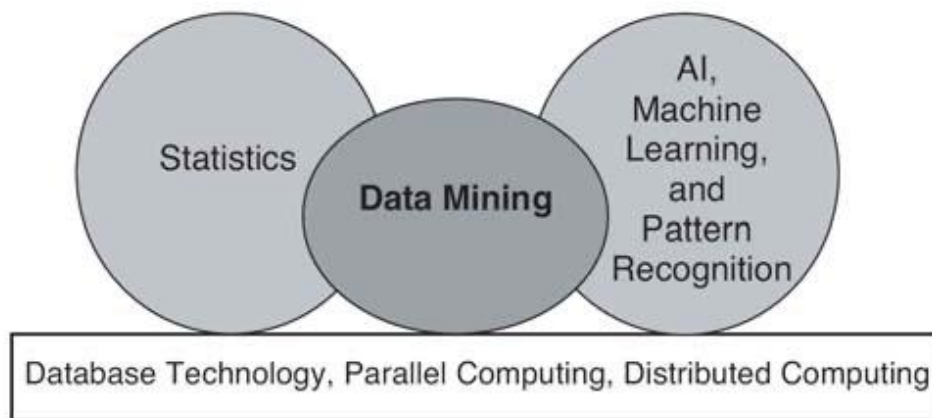
数据挖掘要解决的问题

■ 待解决的问题

- 高维度
- 异构性
- 方法的可伸缩性
- 分布式数据存储

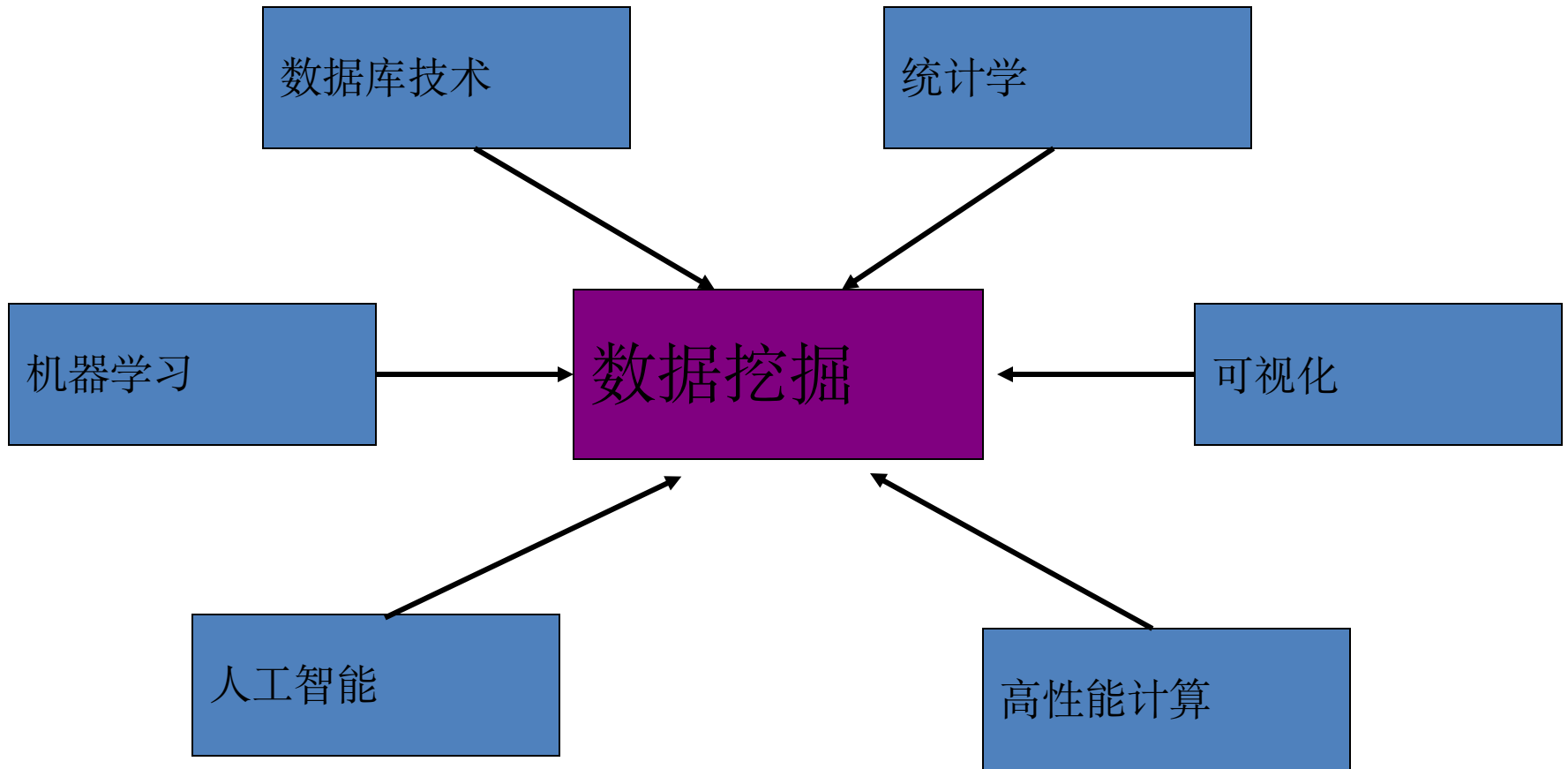
■ 借鉴的方法论来源

- 统计学
- 人工智能
- 机器学习
- 数据库技术
- 分布式计算





数据挖掘是多学科融合产物





数据挖掘基本任务

■ 数据挖掘的基本任务:

- 预测: 根据已有属性值预测特定属性值
- 描述: 概括数据中潜在的关系模式

■ 数据挖掘的基本内容:

- 分类分析 [预测性]
- 聚类分析 [描述性]
- 关联规则分析 [描述性]



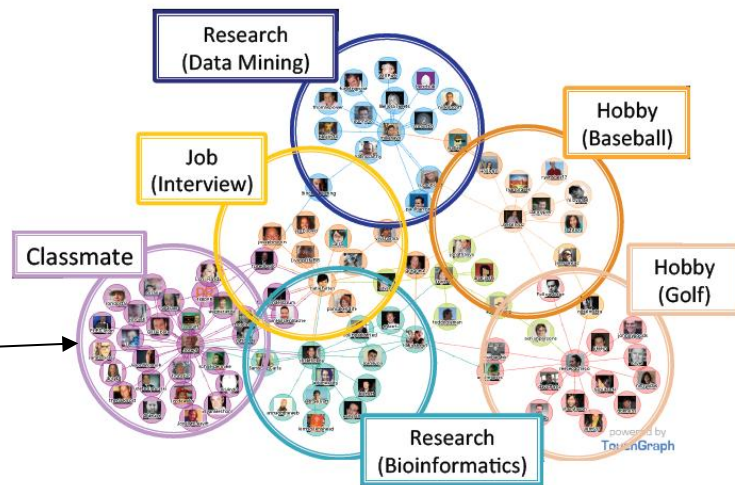
数据挖掘与商业应用

数据挖掘与商业应用



文本的大量使用

People locating (e.g., safety family/ child tracking, friend finder)	at&t FamilyMap	Pink Map
Location check-in/ sharing on social community applications	foursquare	loopt
City/regional guide, neighborhood service search	Google maps	LON lonely planet



个体间联系的网络化

带有位置信息的智能终端化



数据挖掘的应用领域





AI大模型

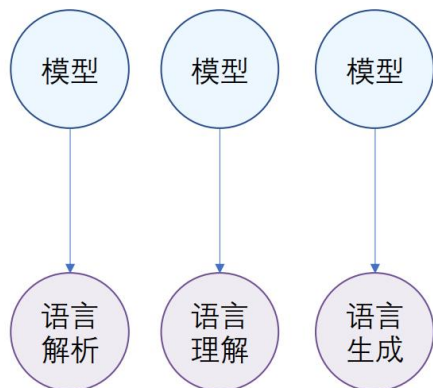
AI大模型：指基于海量数据和强大算力训练出的、拥有数十亿甚至万亿参数的深度学习模型（如统DeepSeek, Chat-GPT等）。它们具备强大的模式识别、自然语言处理、生成和推理能力。



	基础模型				ChatBot	其他应用		
国外	Google LaMDA PaLM PaLM-E	Google DeepMind T5 Imagen Flan Gopher Chinchilla Gato	Meta LLaMA MMS OPT-175B LIMA-65B	OpenAI GPT-4 DALL-E2 CodeX	BigScience Bloom T0 BloomZ	stability.ai Stable Diffusion StableLM	Bard BingChat ChatGPT Claude	Notion AI Cedille AI Copilot Colab Copilot
	Stanford University Stanford Alpaca	databricks Dolly 2.0	AI21 studio Jurassic-1 Jumbo	AI Claude	GPT-J 6B	LMSYS ORG vicuna-13b		
国内	基础模型				ChatBot	其他应用		
	BAAI 悟道 Baidu 文心 达摩院 通义 华为云 盘古 国家超级计算天津中心 天河天元大模型	idea 二郎神 inspur 浪潮 源1.0 JD.COM 言犀 MINIMAX 开放平台 基础模型 科大讯飞 iFLYTEK 星火	澜舟科技 孟子 商汤 日日新 腾讯 混元 网易伏羲 玉言 云从科技 CLOUDWALK 自研大模型		ChatGLM ChatJD 从容 MOSS 商汤 SenseChat 天工 讯飞星火 文心一言 360 智脑	钉钉 斜杠 WPS AI wondershare 万兴科技 学而思网校 MathGPT HAOMO. 雪湖·海若 METASOTA 写作猫	出门问问 序列猴子 EMOTIBOT FRIDAY 有道 youdao 子曰 达观数据 曹植 知乎 知海图AI 小冰	



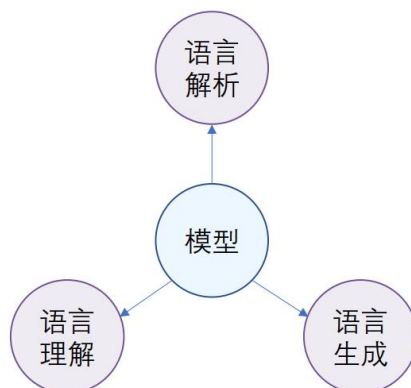
大模型与新学习范式



过去

为每个任务训练独立的模型

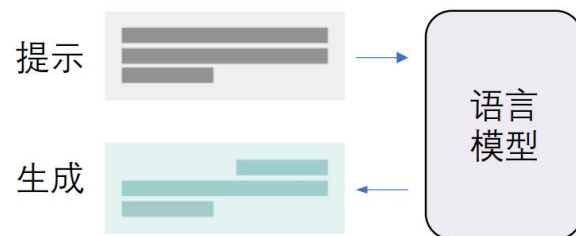
个体化训练



不久之前

中心节点完成预训练，用户在此基础上面向任务微调

中心化训练 + 个体化微调

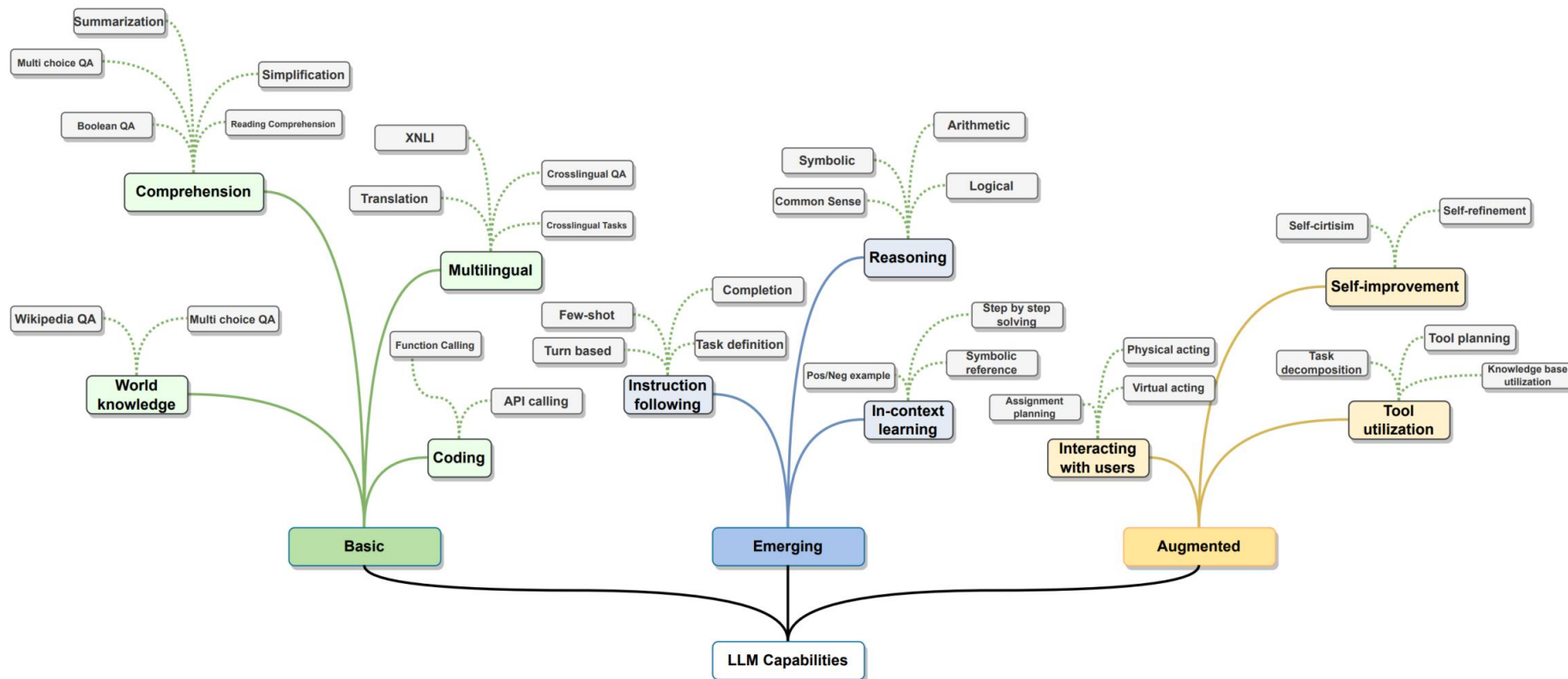


现在（大规模语言模型）

- ▶ 提示学习
 - ▶ 上下文学习
 - ▶ 思维链提示
- ▶ 轻量化微调



AI大模型的能力版图



Shervin Minaee, Tomas Mikolov, Narjes Nikzad, Meysam Chenaghlu, Richard Socher, Xavier Amatriain, Jianfeng Gao. Large Language Models: A Survey. <https://arxiv.org/pdf/2402.06196.pdf>



思考

- 有了AI大模型，数据挖掘的学习还有必要吗？



讲授提纲

- 01** 数据类型与价值使用
- 02** 什么是数据挖掘
- 03** 跨行业的数据挖掘流程
- 04** 课程内容与设计
- 05** 课程学习材料



跨行业的数据挖掘流程



Cross-industry standard process for data mining (**CRISP-DM**)



项目案例

企业画像构建

- 财务指标维度
- 融资风险维度
- 经营风险维度
- 法律诉讼维度

多维度数据收集和处理

- 企业经营数据
- 舆情数据：抓取与企业相关的舆论情绪、社会评价
- 司法诉讼数据：企业法律纠纷情况

融资利率确定

- 基于成本加成的融资利率确定模型
- 融资利率模型开发与落地部署

01 02 03

↓

供应商企业风险评级及供应链金融利率确定



项目案例：业务理解与数据处理

结构化数据—财务信息

项目	2022-03-31	2021-12-31	2020-12-31	2019-12-31	2018-12-31	2017-12-31	2016-12-31	2015-12-31	2014-12-31
营业收入	1,065,917.17	6,222,614.00	5,400,624.69	4,362,368.02	3,726,624.10	3,152,643.77	2,882,846.88	2,680,317.46	2,542,181.14
营业成本	11.16	15.20	15.97	17.06	18.21	9.36	7.56	5.43	8.17
营业利润	41,820.53	296,070.50	279,542.26	279,163.10	253,697.80	232,666.76	197,789.63	182,266.46	161,106.30
净利润	30,733.08	242,679.19	230,677.54	218,103.45	199,750.99	183,305.42	167,655.49	150,430.49	141,522.76
总资产	27,134.70	226,581.50	210,280.45	200,874.55	184,927.26	165,902.03	150,932.34	131,838.32	110,522.41
净资产	20,151.48	243,724.47	212,983.07	178,330.09	131,100.29	134,935.18	106,617.32	106,837.33	64,225.25

结构化数据—财务信息

序号	披露日期	处罚日期	处罚类型	违规事由	文号	处理人	法律依据	原文
1	2022-05-11	-	监管关注	其他	-	上海证券交易所	-	-
2	2016-05-18	-	整改通知	未按规定履行职责	金地银改(2016) 6	上海市地方税务局金山区分局	《中华人民共和国税收征收管理法》	-

非结构化数据—法律纠纷

案由	(2024)辽0106民诉前调3094号	案由	租赁合同纠纷
案件类型	民事	当事人	原告: 沈阳双普交通科技有限公司 被告: 1. 李* 2. 黄** 3. 陈** 第三人: 支付宝(中国)网络技术有限公司
法院	沈阳市铁西区人民法院	承办部门	-
承办法官	-	法官助理	-
案件状态	结案	立案日期	2024-02-27
开庭时间	-	结束时间	-

非结构化数据—年度报告

章节	页码
第一节公司概况	5
第二节会计数据、经营情况和管理层分析	6
第三节重大事件	21
第四节股份变动、融资和利润分配	23
第五节行业信息	27
第六节公司治理	28
第七节财务会计报告	33
附件会计信息调整及差异情况	151



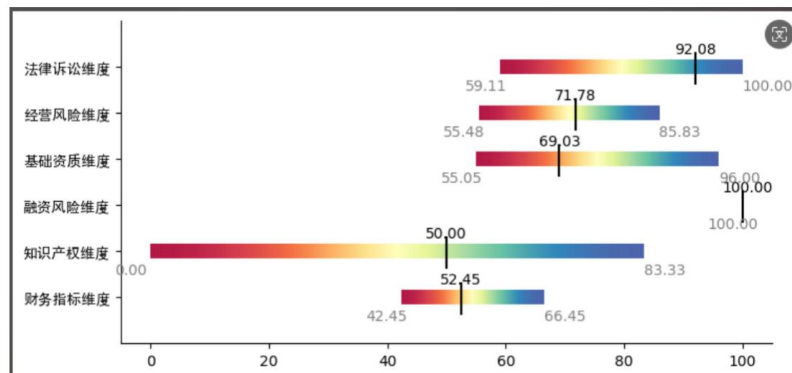
项目案例：成果交付



企业画像各维度得分详情

风险预警	风险标题	更新时间
	【破产清算风险提示】对方(原告)王**与上海天德建设(集团)有限公司破产清算特别债权人第一次债权人会议决定书	2022-02-21
	【买卖合同纠纷】对方(原告)上海浦东机场建设发展有限公司与上海天德建设(集团)有限公司买卖合同诉讼一审案件民事裁定书	2022-02-10
	【建设工程合同纠纷】上海特精建筑工程有限公司与上海特精实业有限公司建设工程施工合同纠纷民事二审案件民事裁定书	2022-01-28
	【股权转让合同纠纷】对方(原告)林**、吴**与股权转让合同纠纷二审民事裁定书	2022-01-25
	【股权转让纠纷】上海天德建设(集团)有限公司与股权转让合同纠纷二审民事裁定书	2022-01-17
	【买卖合同纠纷】对方(原告)东莞市国康实业有限公司与上海天德建设(集团)有限公司买卖合同诉讼一审案件民事裁定书	2022-01-04
	【买卖合同纠纷】对方(原告)东莞市国康实业有限公司与上海天德建设(集团)有限公司买卖合同诉讼一审案件民事裁定书	2022-01-04
	【股权转让合同纠纷】(270000.0元)上海天德建设(集团)有限公司与上海国康科技有限公司股权转让合同纠纷诉讼一审案件民事判决书	2021-12-21
	【提供劳务者受害责任纠纷】对方(原告)陆**与415412541541515号提供劳务者受害责任纠纷诉讼一审案件民事裁定书	2021-12-02
	【买卖合同纠纷】(68080.0元)宣州市高塘镇顺泰水泥制品构件厂、上海天德建设(集团)有限公司买卖合同诉讼一审民事判决书	2021-11-29

指标预警信息列表详情



企业各维度得分情况在同行业中所处位置



项目案例：成果交付

供应链金融应用DEMO

机械有限公司 工程机械II

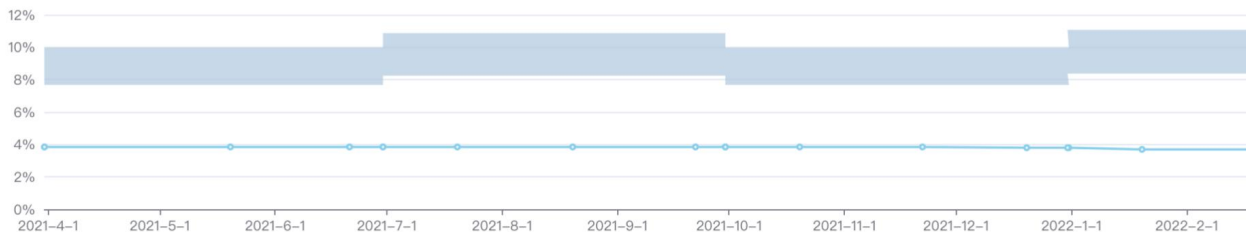
供应商画像

利率区间

风险预警

利率区间:8.25%-10.86% 1

LPR 金融机构加权平均贷款利率 2



3

时间年限 ?

时间年限
1年

违约损失率最小值:0.2 ?

违约损失率最大值:0.4 ?



项目案例：成果交付

供应链金融应用DEMO

时间年限 ?
时间年限
1年

4

违约损失率最小值:0.2 ?

违约损失率最大值:0.4 ?

5

[供应商相关信息]

股权价值 ?	14329236024
股权价值年化波动率% ?	50.01
流动负债 ?	3973385416
非流动负债 ?	951716763

数据测试

[资金使用信息]

财务成本 ?	0.043
期望利润率 ?	0.045
违约风险补偿 ?	0.021

数据测试

[测试输出记录]

测试[2022/5/20 14:17:32]	信用评级: A	利率: 9.73%	利率区间: 8.25% - 8.25%	违约概率: 3.05%
测试[2022/5/20 14:17:33]	信用评级: A	利率: 9.73%	利率区间: 8.25% - 8.25%	违约概率: 3.05%

6

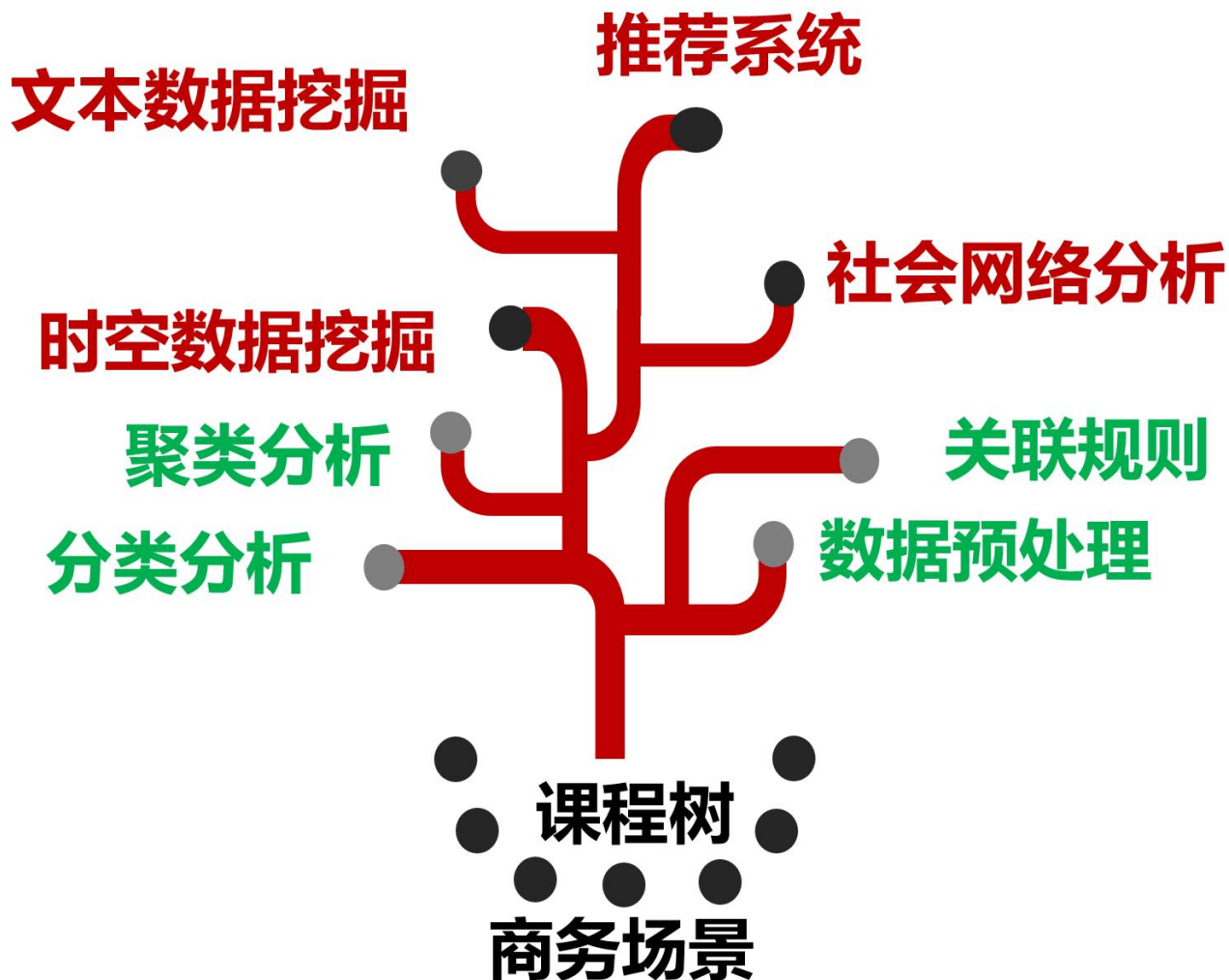


讲授提纲

- 01** 数据类型与价值使用
- 02** 什么是数据挖掘
- 03** 跨行业的数据挖掘流程
- 04** 课程内容与设计
- 05** 课程学习材料

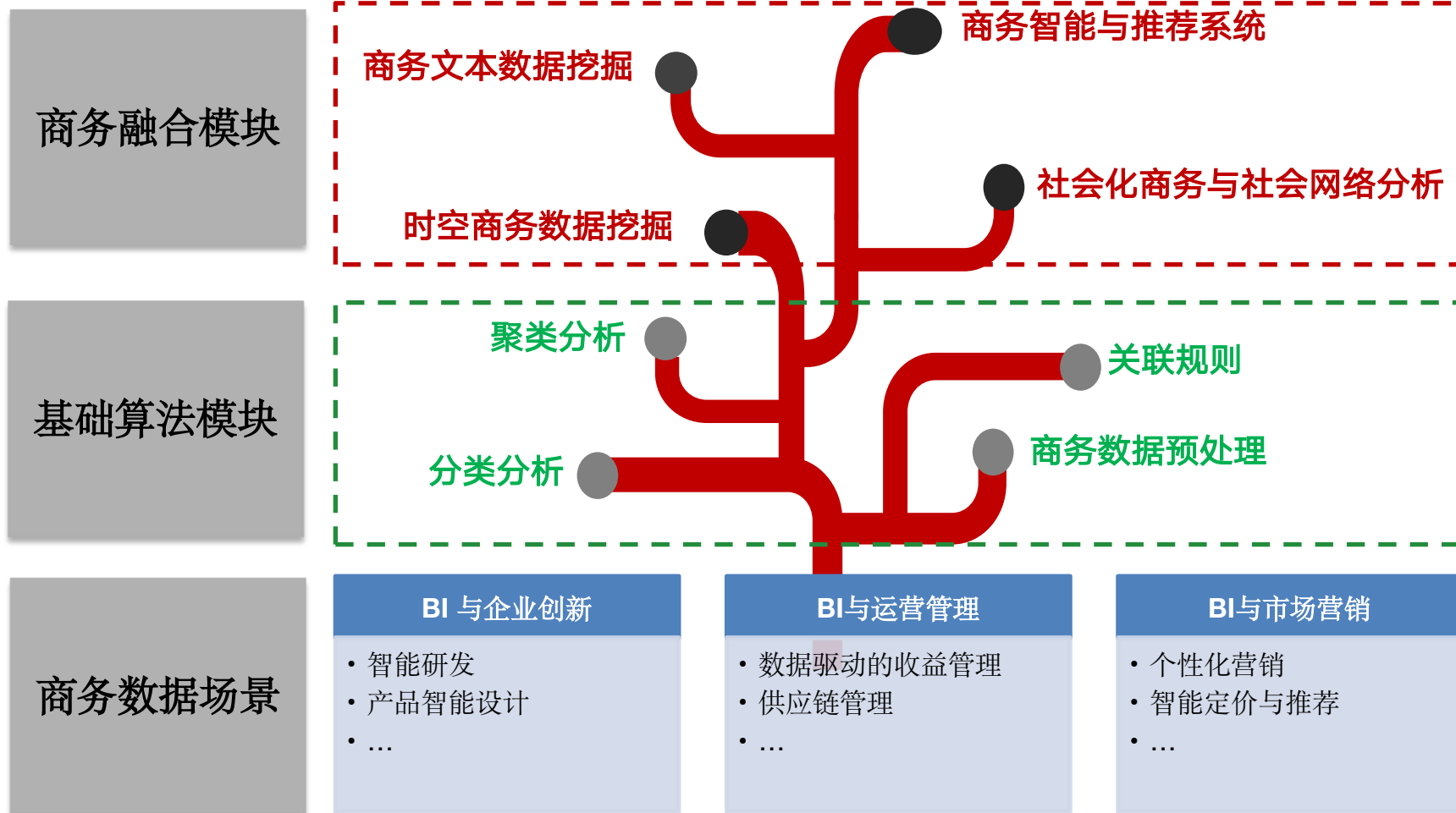


课程内容





课程具体设计





先修课程要求

■ 理论课程（建议）：

- 概率论与数理统计
- 高等数学

■ 编程课程（建议）

- 具有Python/R 编辑基础

■ 已修/正修如下课程的可不选修本课程：

- 《数据挖掘》课程
- 《机器学习》课程



课程考核计划

- 考勤及课堂表现(10%):
 - 随机点名
 - 课堂表现
- 随堂测测验 (20%)
- 个人作业(30%)
 - 数据分析实践
 - 数据分析方法原理练习
- 期末Project (40%)
 - 个人/团队均可
 - 总人数 ≤ 5 人
 - 不能直接使用其他课程的期末项目



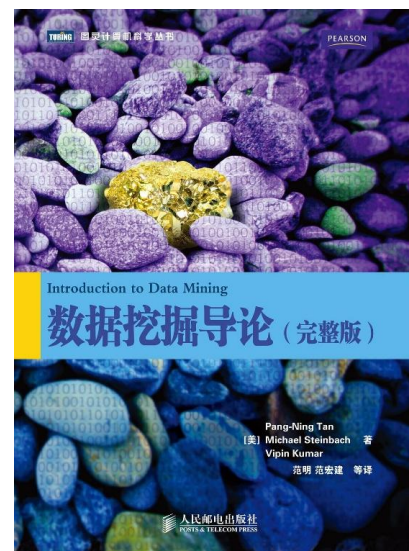
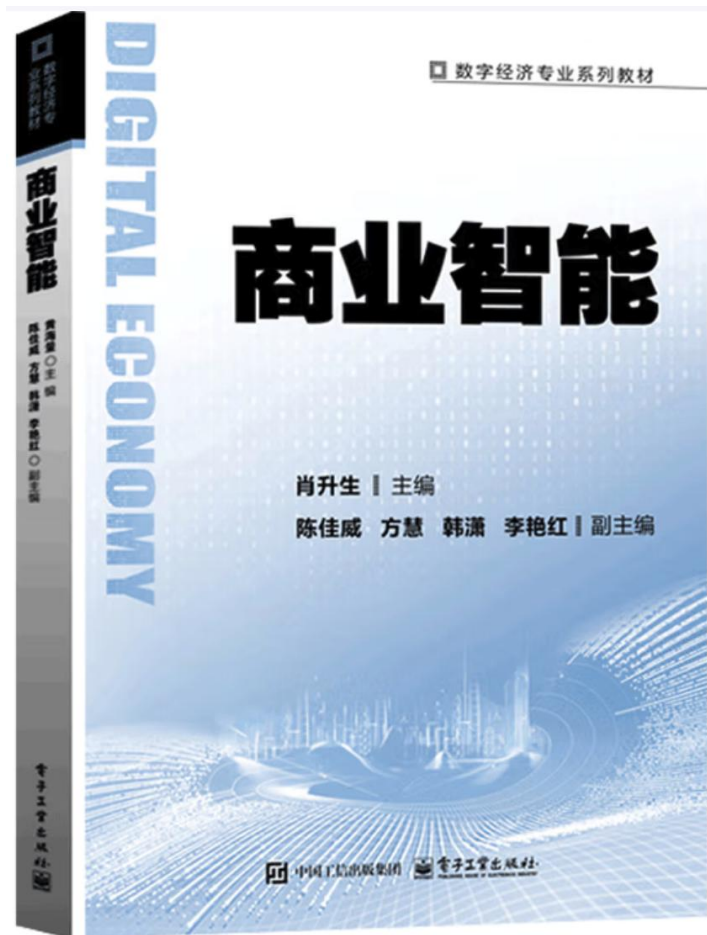
讲授提纲

- 01** 数据类型与价值使用
- 02** 什么是数据挖掘
- 03** 跨行业的数据挖掘流程
- 04** 课程内容与设计
- 05** 课程学习材料



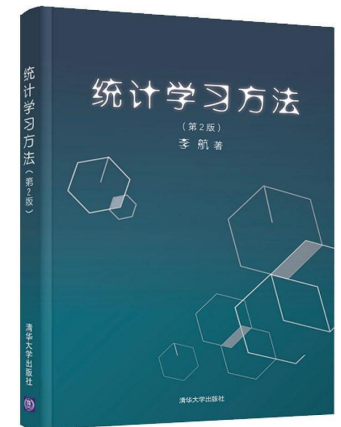
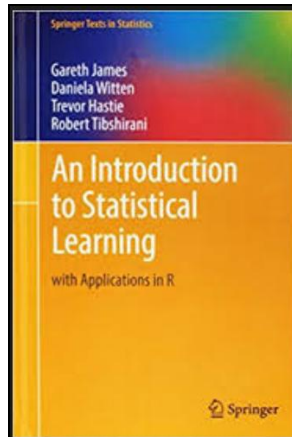
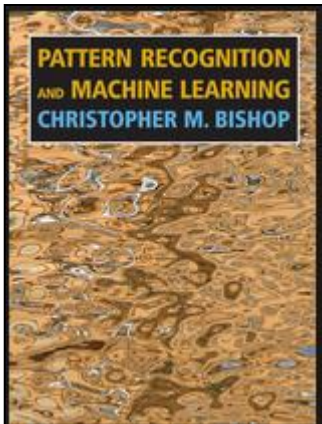
课程参考教材

■ 相关教材





课外阅读材料





更多学习资料

■ 理论学习

- ACM SIGKDD Conference on Knowledge Discovery and Data Mining
- International Conference on Machine Learning
- International Conference on Data Mining
- IEEE Transactions on Knowledge and Data Engineering

■ 实践学习

- 天池大赛: <https://tianchi.aliyun.com/>
- Kaggle: <https://www.kaggle.com/>



数据挖掘与商务分析



400年前发明了显微镜，改变了测量的标准，人类研究物体的细微程度从此不同。

大数据分析带来的变革，就像400年前的显微镜一样，我们能够掌握事件、行为的精细程度，也将从此进入全新的境界。

—— Erik Brynjolfsson