蚂蚁安全领域智能化数据治理

高级数据技术专家/霄元(黄国龙)



个人简介

黄团龙高级数据技术专家

- 20年加入蚂蚁,目前是安全大数据团队数据智能资产、内容数据资产负责人
- 0-1主导建设蚂蚁集团内容安全数据资产专项
- 关注大数据研发、数据架构、数据治理等领域





数据治理通用能力

蚂蚁数据治理架构及能力

数据治理 业务案例

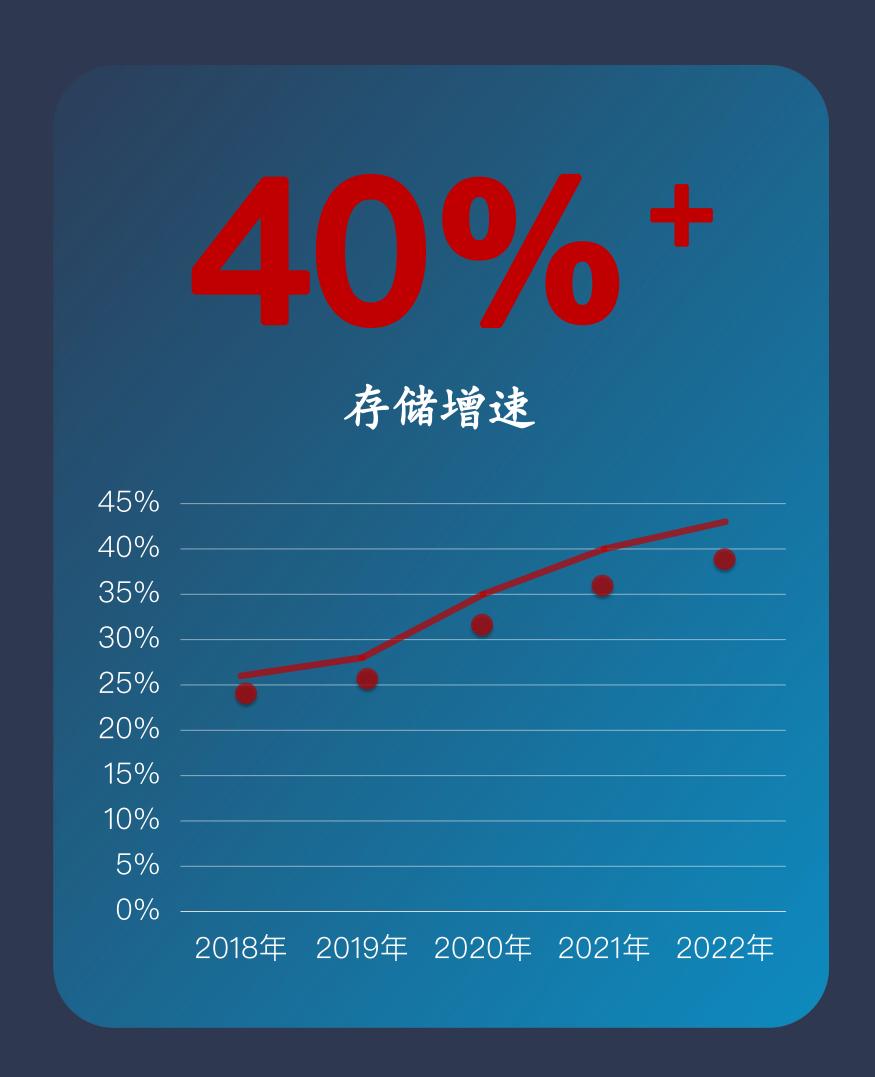
安全领域智能化数据治理实践

数据治理 创新探索

数据治理探索与创新



蚂蚁数据治理架构及能力_面临挑战



成本增速快

- 存储资源:蚂蚁存储达 EB 级别,年增速40%+;
- 计算资源: 计算资源 千KCU/日.

业务需求多

- 新业务资源需求旺盛;
- 人工智能发展快,数据需求呈现爆炸性增长。

成本看不清

- 资源使用细节看不清;
- 成本很难分摊到业务。





蚂蚁数据治理架构及能力_治理思路



能力建设

平台提效 & 技术治理

平台提效:将治理能力产品化服务用户

技术治理: 结合引擎能力升级做技术优化

无效资产治理

重复资产治理

TOP 资产治理

存储治理

计算治理

長群治理



蚂蚁数据治理架构及能力_治理方案

关键能力建设方案:从事前->事中->事后,构建成本治理全链路能力。

触发管控场景 触发治理场景 场景 年度资源预算 模块存储触顶 任务大量变慢 驱动 系统账号开通 项目任务大量变慢 用户加入Project 架构师驱动 管理员驱动 管理员驱动 (事前)规划与准入 管控与监控 (事中) (事后)技术治理 发布管控 查询管控 账号管控 运维管控 项目迭代式 运 生 专项运营治理 架 运营治理 命 维 生命周期 暴力扫描 构 并发限制 补数管理 治理 离职人员资产处置参数不合理 重复表识别 管 周 策略 生命周期长 资产汰换 数据排重 期 控 大表限制 补数监控 血缘要求 汰换任务 热点任务识别 大表暴力扫描 资产下沉推荐 惩 效 管 模 重复采集 过渡埋点 资源混部 理 块 汰换数据 运行监控 无效下线 暴力扫描 Archive压缩 冷数据归档 数据加工(DataPhin) 数据生产 数据同步 数据应用 治理 应用DB → DRC → AntQ explorer 回流任务
应用DB 对象 →加工任务 应用系统 * Rlink

(ods)

odps



应用日志 — SLS



(ods)

蚂蚁数据治理架构及能力_治理架构

任务

并发

资源 管理 门户

资源

治理

领域

资源监控大盘(DRE)

跨集群 带宽

TOP 任务

暴力 扫描 集群 存储

资产治理工作台(个人/团队/业务单元)

资产 健康分 垃圾 资产 处置

治理 活动 分析

资源治理核心领域(专项方案)

风险 拦截 治理 列表 工具

一键

资源成本管理(管理者/DRE)

预算 管理

资源 调拨

关键技术

引擎优化

模型优化

代码优化

资产管理优化

成本 核算

组织文化

治理 委员会

治理达人

红黑榜

治理双周 /月报

资源调拨

任务分时调度

自动化扩缩容

调度并发控制

单元化隔离

业务单元化容灾

集群管理

集群资源混部

资源预算管理

数据项目规划

数据分级存储

数据极限存储

采集治理

无效采集下线

场景化生命周期

日志治理

无效埋点下线

重复采集治理

日志消费管控

业务单元化容灾

消费治理

资源

分析

一键链路退役

废弃报表下线

无效服务下线

消费血缘保鲜

制度规范

标准规范

计存军规

考试培训

(事前-规划&管理) 资源预算管理

基础 平台 能力

业务

管理

预算 分配

预算 分析

预算

发布 管控 拦截

实时

资源 风险 处置

自动化 扩缩容

弹性 分时 调度 健康分

场景化 生命 周期

资产治理中心

自动化 技术 治理

(事后–治理)

业务单 元治理 活动

访问

存储

资源调拨中心(事中-分配&监控)

统一 资源 元数据

预 算

血缘

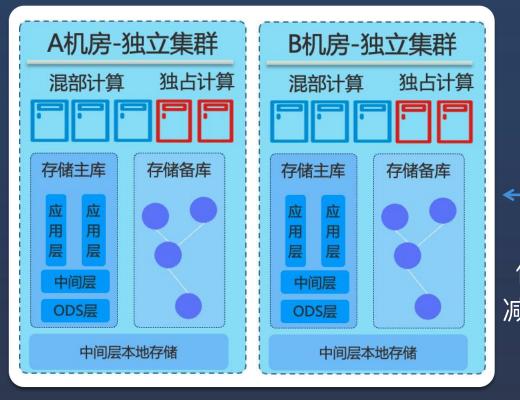
调度

蚂蚁数据治理架构及能力_资源治理领域案例

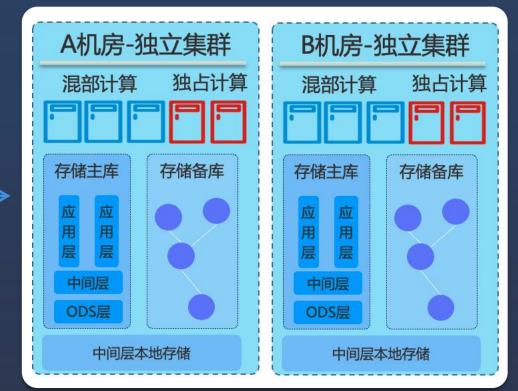
集群资源混部

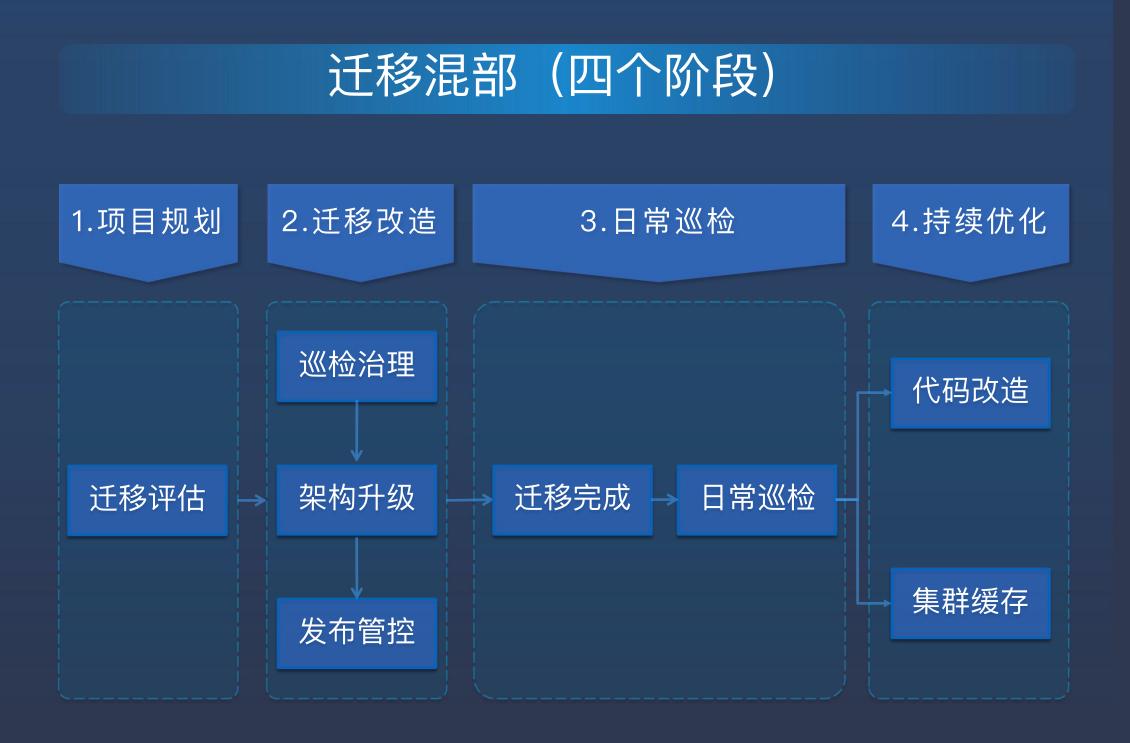
推进在离线混合部署,计算算力会提升10%,机器成本降低25%.

混合部署方案 提升资源利用率, 动态扩容,保障稳定性 存储计算一体 => 存储计算分离 深圳混部



跨城访问 依赖中间层 减少网络开销





开源



蚂蚁数据治理架构及能力_资源治理领域案例

园 关键技术

提升治理自动化率,实现自动识别、归因分析、自动清理,形成常态化管控能力。

引擎优化

参数优化

- Split Size
- 小文件合并
- Reducer instantce
- CPU数
- Dynamic parallelism

调度优化

- 任务归并
- HBO优化
- 集群混部
- 错峰运行
- 冷热分层

模型优化

数仓模型

- 业务领域建模
- 抽象公共层
- 通用应用层
- 配置化指标系统
- 大宽表设计

代码设计

- 全量改采样
- with替代tmp表
- 视图化改造
- 避免数据膨胀
- 执行顺序优化

计存设计

- 渐进计算
 - 累计计算
 - Zorder
 - Shuffle优化
 - Bitmap索引
 - 全改增
 - 极限存储

数据格式

- 重排压缩
- Cube预计算

代码优化

JOIN优化

- Map join
- Hash join
- Skew join
- Dynamic Filter

数据倾斜

- Map端
- Reduce端
- 热点值

UDF优化

- 内置替换
- 提前计算
- 参数调优
- 本地缓存

聚合优化

- Grouping sets
- UDF转UDTF
- Count(distinct)

资产管理优化

生命周期

- Map join
- Hash join
- Skew join
- Dynamic Filter

无效表

- 热点值

• Map端

• Reduce端

计算浪费

- 临时表
- 系统表
- 长周期表
- 回收站优化
- 大字段生命周期

重复资产

- 同链路相似表
- 相似任务节点
- 分区不更新
- 缓慢变化维表

节流

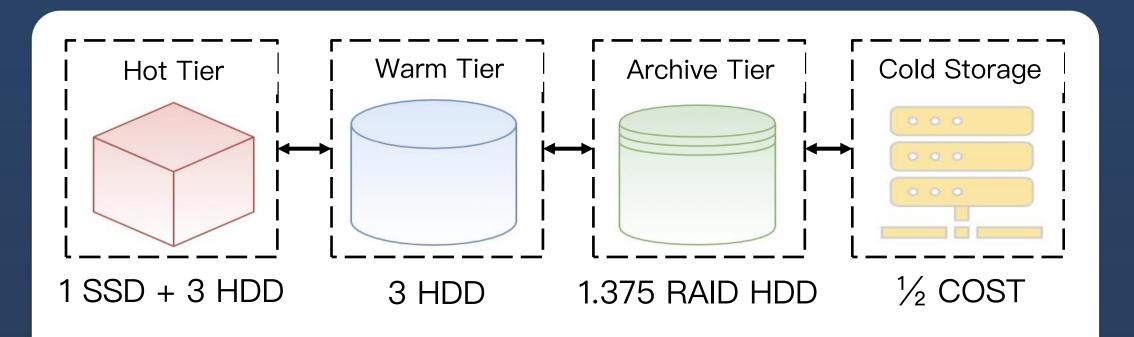


蚂蚁数据治理架构及能力_资源治理领域案例

⇒ 关键技术-实例

冷存体系

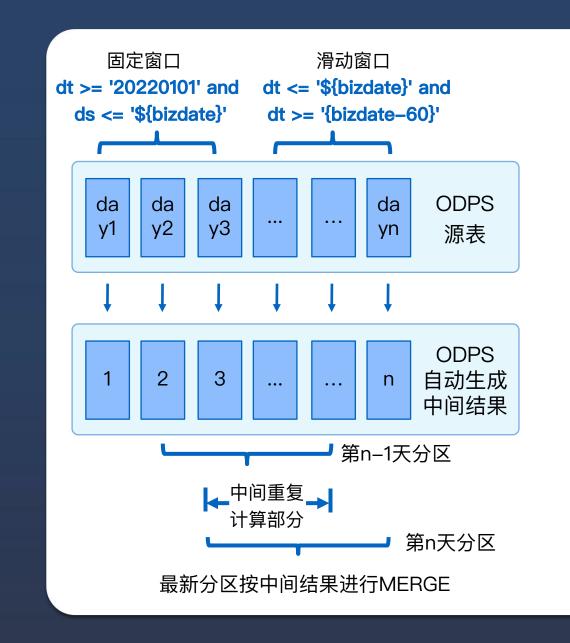
推进在离线混合部署,计算算力会提升10%,机器成本降低25%.



- ➤ Hot Tier: 高频消费的热点数据、优化I/O;
- ➤ Warm Tier: 热数据、读取频率正常;
- ➤ Archive Tier:数据需长期保留,访问频次底;
- ➤ Cold Storage: 长期保留,超低频访问。

渐进计算

设置成渐进计算后,每日计算消耗从795CU降到22CU.



原理:空间换时间,自动生成中间表,避免重复计算,其中中间表可采用hash cluster,提升merge阶段Shuffle效率; Odps支持一键渐进计算、设置一个参数即可。

```
SELECT userid,

SUM( CASE WHEN **** END),

COUNT(CASE WHEN **** END),

MAX(CASE WHEN **** END)

from aseccdm.dwd_sec_evt_______
where dt > '${bizdate_90}'

and dt <= '${bizdate}'

and ((coalesce(biz_no_1,'') <> '')

or (event_name = 'cushcouponnechare')

or (event_name = 'plant, but i')

or (event_name = 'dlant' and length(biz_no) > 0))

group by userid
```





蚂蚁数据治理架构及能力_基础平台能力实例

(事前)发布管控

由平台或业务方事先制定发布管控规则,相关任务上线提交时、平台自动检验结果,如有规则未校验通过,此次上线发布失败。

发布管控规则说明



任务发布状态详情



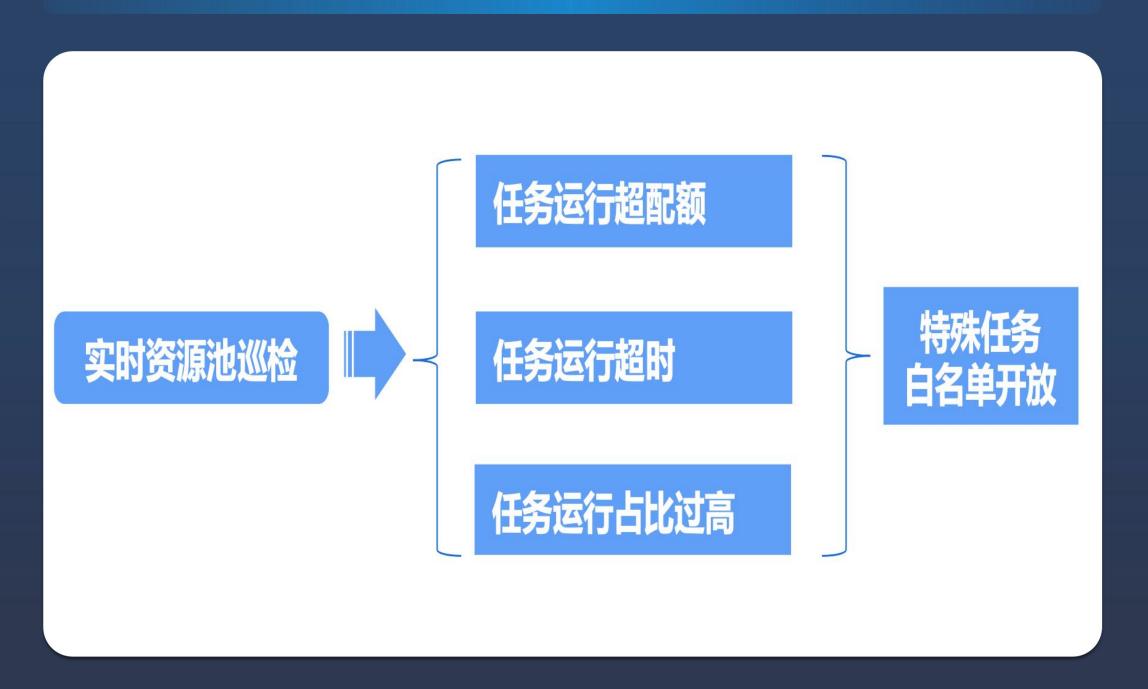


蚂蚁数据治理架构及能力_基础平台能力实例

(事中) 实时巡检

禁止随意提交超大任务导致整个资源池的堵塞和打满,进而造成高昂的成本消耗,同时也兼顾效率,仅对异常使用进行管治。

异常大任务自动查杀方案



异常大任务自动查杀实例

项目 实例

20221125151538929gebd2g6fn375累计资源使 用超告警阈值

- 告警名称 不境禁止20CU以上的 任务
- 项目名称: a
- 作业提交人
- instance_id:

20221125151538929gebd2g6fn375

- 快照时间: 2022-11-25 23:46:07
- 告警时间: 2022-11-25 23:48:23
- cpu累计使用: 20.0(CU), 当前阈值: 20(CU)
- gpu累计使用: 0(CU), 当前阈值: (CU)
- 内存累计使用: 26.0(GB), 当前阈值: (GB)
- 作业自动查杀结果: 作业已被自动查杀

查看logview(代码)

历史告警事件

一键查杀作业

项目名称: 7

- 告警名称: □□□□□□□ 资源消耗超过90%
- 作业提交人:
- instance_id: 20221125125009480gcnliffc2465
- 快照时间: 2022-11-25 20:56:07
- 告警时间: 2022-11-25 20:57:34
- cpu利用率: 395.30%, 当前阈值: 90.00%
- gpu利用率: 0.00%, 当前阈值: 90.00%
- 内存利用率: 123.53%, 当前阈值: 90.00%
- 作业自动查杀结果: 作业已被自动查杀

查看logview(代码)

查杀作业

历史告警事件

阈值有疑问请联系告警添加人:

查杀异常或者优化请加群联系值班同学处理:

11799056

detail



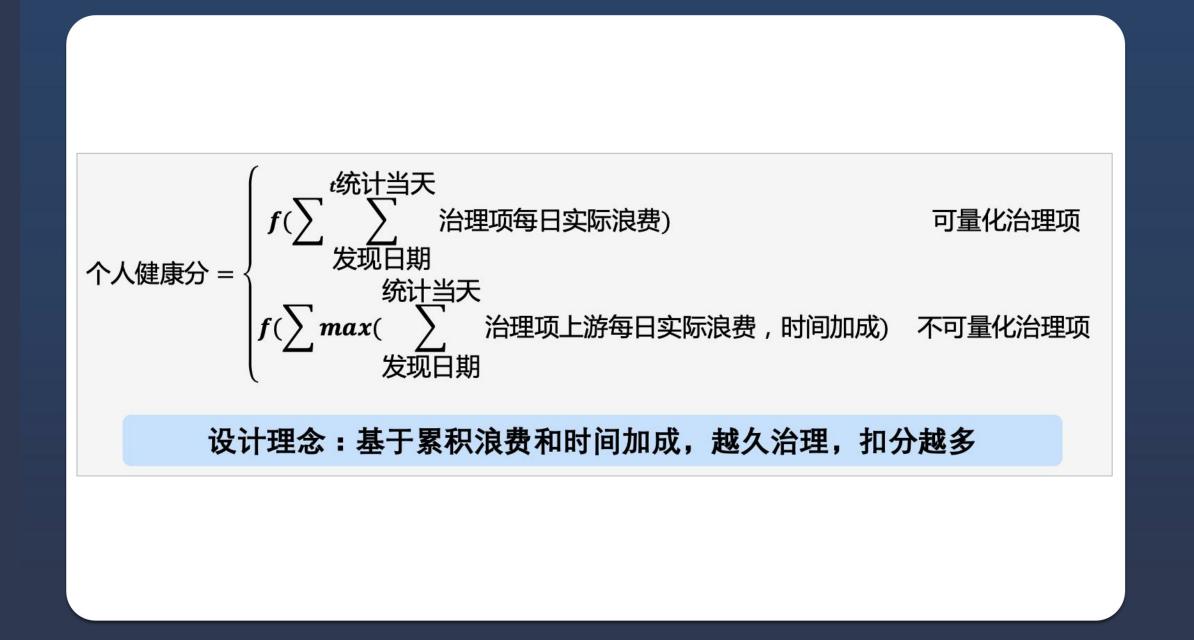


蚂蚁数据治理架构及能力_基础平台能力实例

(事后) 成本健康分

基于累计浪费和时间加成,设计成本健康分算法。通过健康分管理数据平台使用权限。

成本健康分算法



健康分产品运营实例







蚂蚁领域数据治理架构及能力_治理成果

安全领域治理成果总结

已全部达成年度目标,预估节约数据成本 25%+。

治理专项	存储用量(PB)	计算用量(KCU)	表数量(张)
累计治理收益	百PB+	百KCU+	百万+

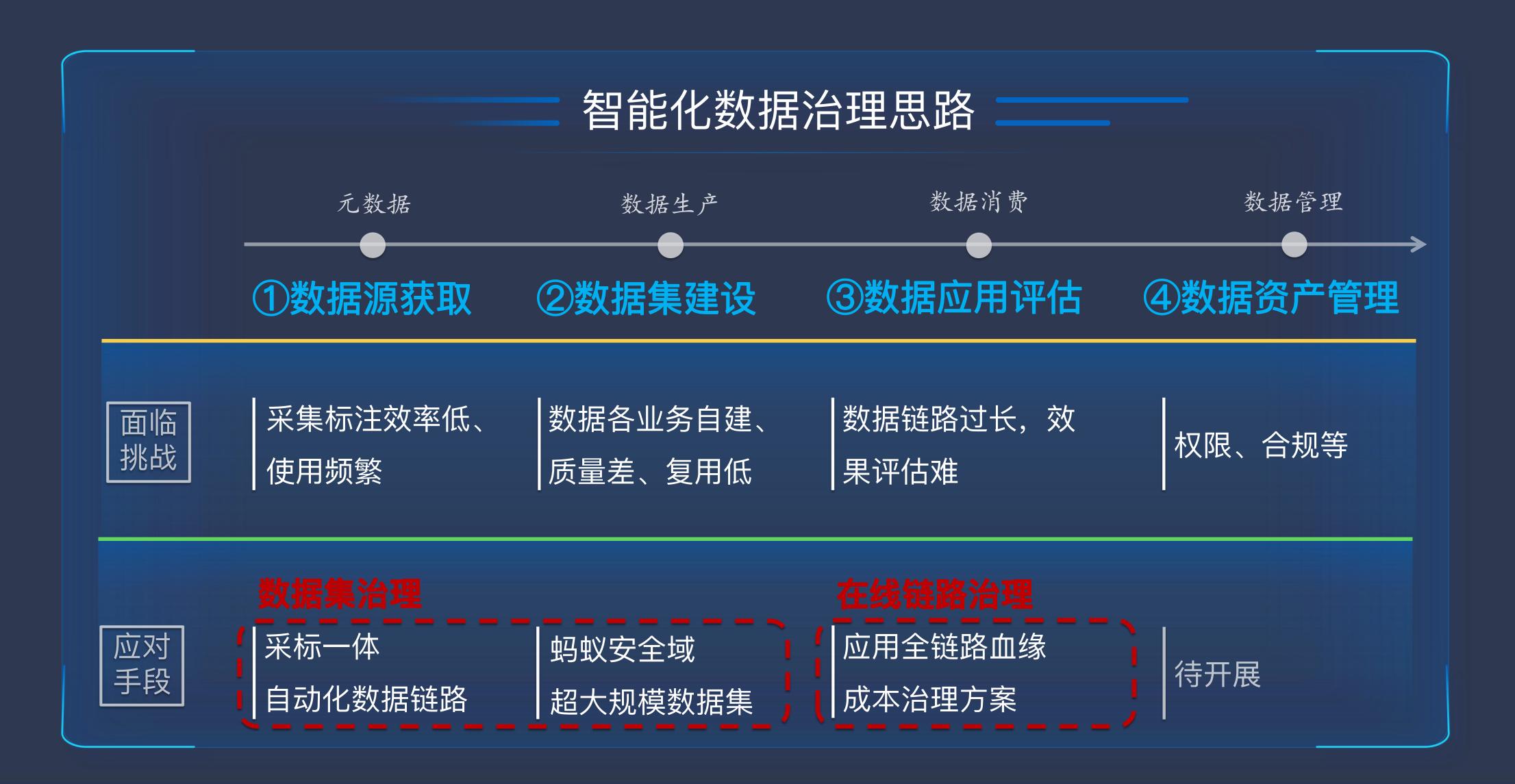


安全领域智能化数据治理实践





安全领域智能化数据治理实践





安全领域智能化数据治理实践



数据集建设阶段, 搭建采集标注自动化数据链路降本增效; 标准化建设蚂蚁安全域超大规模数据集消除数据孤岛……

※ 数据集治理

采标一体自动化数据链路 蚂蚁安全域超大规模数据集

在线运行阶段, 搭建元数据之应用全链路血缘, 助力在线模型策略成本治理优化……

② 在线链路治理

应用全链路血缘 在线模型成本治理方案





安全领域智能化数据治理实践_数据集治理

数据集简介 _____

高品质、多样性、大规模的数据集建设是AI技术应用竞争关键要素之一,

在建设大规模数据集过程中,百PB级别的安全数据资产必将带来高昂成本支

出、及各种质量风险隐患

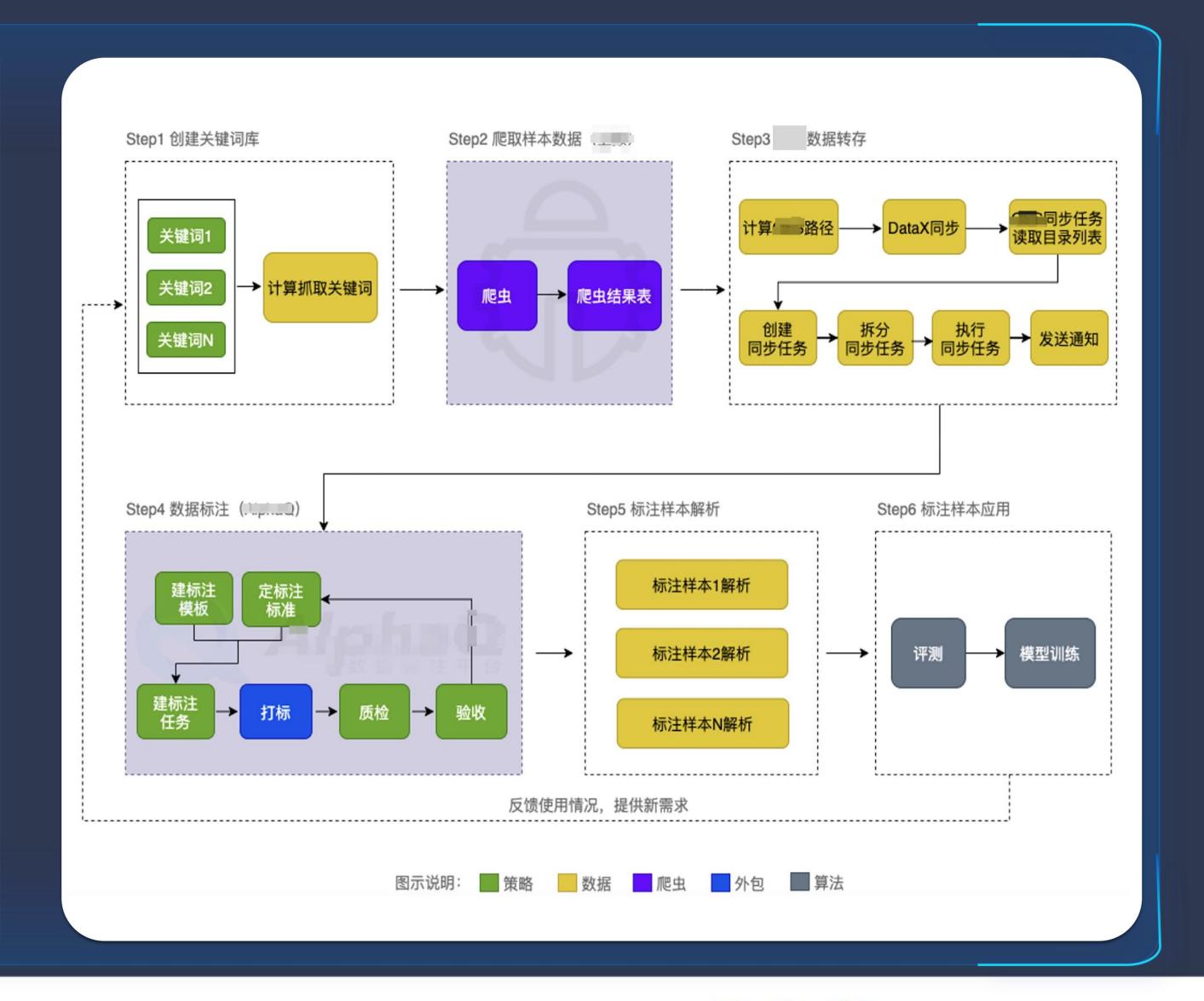
数据集主要有:采集、标注、大规模数据集、训练数据集、评测数据集。



安全领域智能化数据治理实践_数据集治理

数据采集、人工标注是数据集建设关键环节之一、也是首要事项,相关事项涉及合作方众多、且时间也不可控,在人工对接中费时费力。

采标一体化自动化数据链路,从关键 词计算、对接采集、转存及通知、对接标 注、数据ETL全链路实现自动化,降本增 效明显、且数据品质也有保障,大大缩减 人力成本、将原来采集标注2周以上时效降 到5天以内。





安全领域智能化数据治理实践_数据集治理

智能数据建模设计架构

标准数仓建模设计确保数据品质, 通过大规模数据集支撑业务。



大规模数据集分类体系

规范数据集分类体系,消除数据孤岛、共享数据集资源。

-					
类型	来源	数据集	规模(存储/量级)	类型	特点
内容域 —— ——	蚂蚁外	MMC4	存储: 512G / 文本: 23,085 存储: 1.8T / 特征: 23,085	文本图片对	十亿图像语料库Multimodal-C4 fewer-faces数据
		Laion-5b	存储: 24.25T / 量级: 3,769,298,024	文本图片对	非人工注释的多模态图文数据集, 训练图文匹配能力
		语音对话数据集	存储: 496.9G / 量级: 1,982,402.	音频	群聊、对话等场景音频数据
		纸币虚拟币数据集	存储: 164G / 量级: 1,086,253	图片	人民币、美元、比特币等相关图片 样本数据
	蚂蚁内	蚂蚁图片内容源	存储: ***T / 量级: *亿	图片	
	生成式(AIGC)	伪造图片标注数据集	存储: ***M+ / 量级: ***亿+	图片	赌博等图片集
资金域	蚂蚁内	事件特征数据集	存储: ***T / 量级: *亿	文本	
		·	·		

实例: *数据集数据建模设计

数据集实例,整合资金各业务特征、标签形成全域样本集。





安全领域智能化数据治理实践_在线链路治理

在线链路治理简介

良好的治理离不开对数据资产合理评估,通过对当前安全领域的数据资产

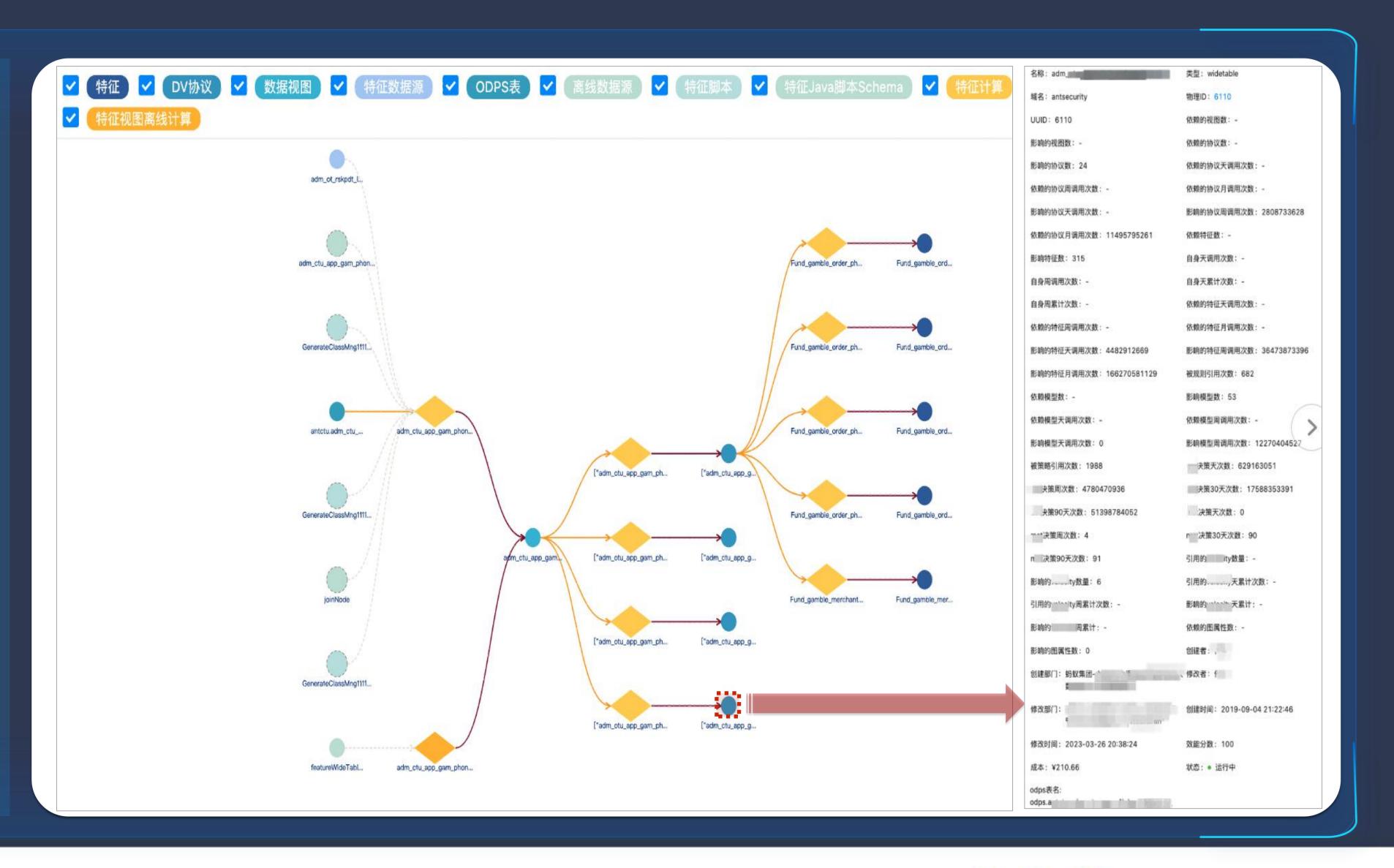
进行了效能评估,产出了数据效能分,量化数据资产在风控系统中发挥的作

用,从而推动无效资产的治理,计算、保障优先级设定,成本优化等。



安全领域智能化数据治理实践_在线链路治理

数据资产类型 包括策略、特征、 模型、协议等90+ 种资产类型,种类 多、数据资产量巨 大、关系链路复 杂,利用我们的二 部图模型,构造了 一张全局的资产大 图,从连接起各个 信息孤岛,打破平 台间的血缘鸿沟。





全链路

血

缘



安全领域智能化数据治理实践_在线链路治理

在线链路治理流程

影响因子定义

影响因子分:血缘静态引用量、 线上流量调用量以及决策日志量 三个层次。

影响因子分层:

- ① 根据数据血缘,计算出该数据资产与其它数据资产在静态血缘上的引用量;
- ② 计算该数据资产线上产生的实际流量,例 如特征的调用量、模型的调用量;
- ③ 计算关联该数据资产的决策日志的量级,作为该资产在风控体系中发挥效能的重要特征。

备注: 为了更合理的评估数据资产长短期的效能情况, 分别从天/周/月/季的维度作为特征。

指标数据

基于上述影响因子特征,统计得到了从不同维度评估数据资产效能的源数据。

字段 描述	被特征 引用次数	被规则 引用次数	被策略 引用次数	特征快照 周调用次数	决策日志 周调用次数	•••••
f_1	45	125	67	3,433	120	-
f_2	115	44	-	1,200	4	-
f_3	-	12	100	12	-	-
f_4	86	-	65	10	-	-
f_5	77	10	4	566	1	-
	-	-	-	-	-	-

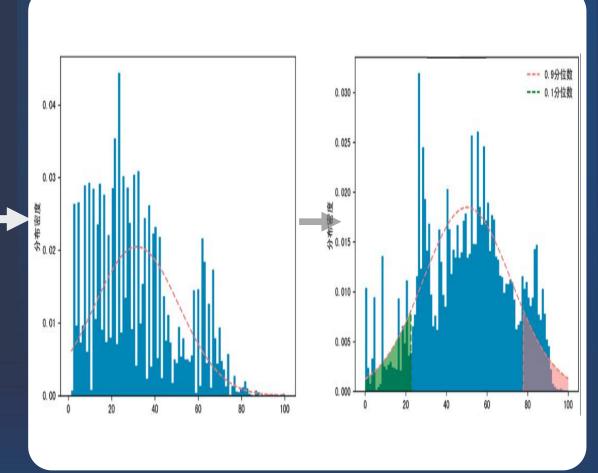
编码器推理

自编码器,对一组特征进行学习,得到有效表征。

encode decode input hidden output

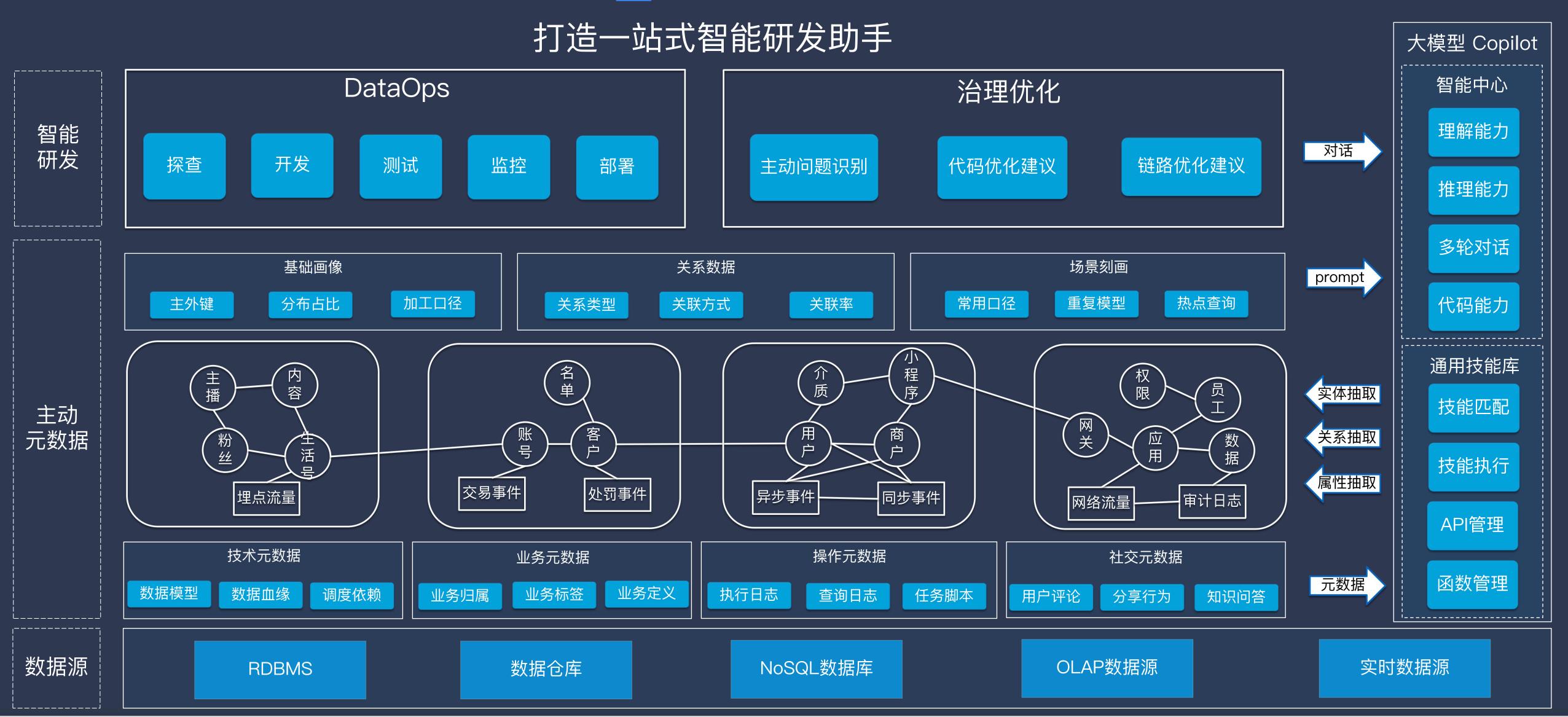
效能分变换

使用回归模型对特征进行回归训练,获各特征重要度;基于特征重要度优化模型和调权重,使效能分产出更合理;通过BOX-COX变换对数据分布进行调整(0~100分正态分布)。





数据治理探索与创新_ ETL AUTOPIPELINE







数据治理探索与创新_创新案例大模型Copilot

小 表 D

小表D:安全大数据一 站式智能研发助手,结合 安全特色,深度整合其他 数据类大模型,以小表D为 切入口,为用户提供丰富 的大模型功能,贯穿用户 整个数据研发生命周期, 在数据分析、任务研发、 任务运维、风险发现等日 常生产环节提供一站式数 据辅助服务,让数据研发 更加智能高效。







想一想,我该如何把这些技术应用在工作实践中?

THANKS



