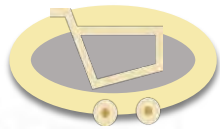


大规模机器学习平台的技术实现

第四范式 胡时伟

2017年 8月5日

AI Works



电商推荐



计算广告



差异化定价



授信审批



内容推荐



实时风控



精准营销



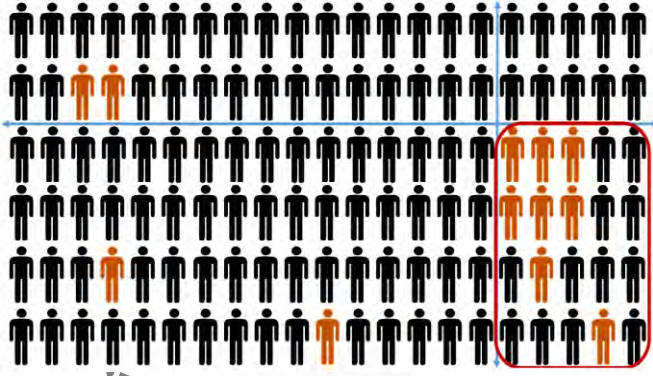
智能产品推荐



智能客服

精细：对个性化和微观业务场景的分析和预测能力要求早已远超传统企业的想象

传统客户触达



传统客户触达：

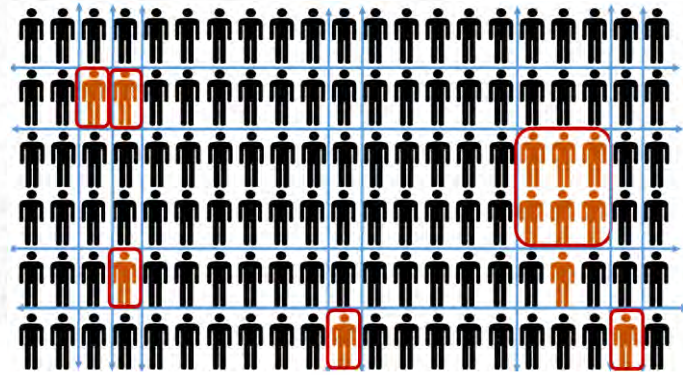
- 用少量特征将用户较为粗放的划分到少量类别中，每个类别中的用户被认为有相似的属性和相同的意愿，丢失了对每个用户的个性化描绘，准确性有限。同时也无法覆盖到部分客群中的个性化用户



AI客户触达

大数据机器学习模型：

- 基于日益丰富的海量数据样本，和千万以上量级数据特征，将用户细分到微观粒度，对每个用户做精细的个性化描述，直接定位到每个有意愿的用户，更精准，更全面



智能：要求企业能够适应不断变化的内外部环境，实现数据价值



传统的决策规则政策迭代周期：
数月、半年甚至一年以上

互联网的决策规则政策迭代周期：
每天、每小时甚至只需每分钟



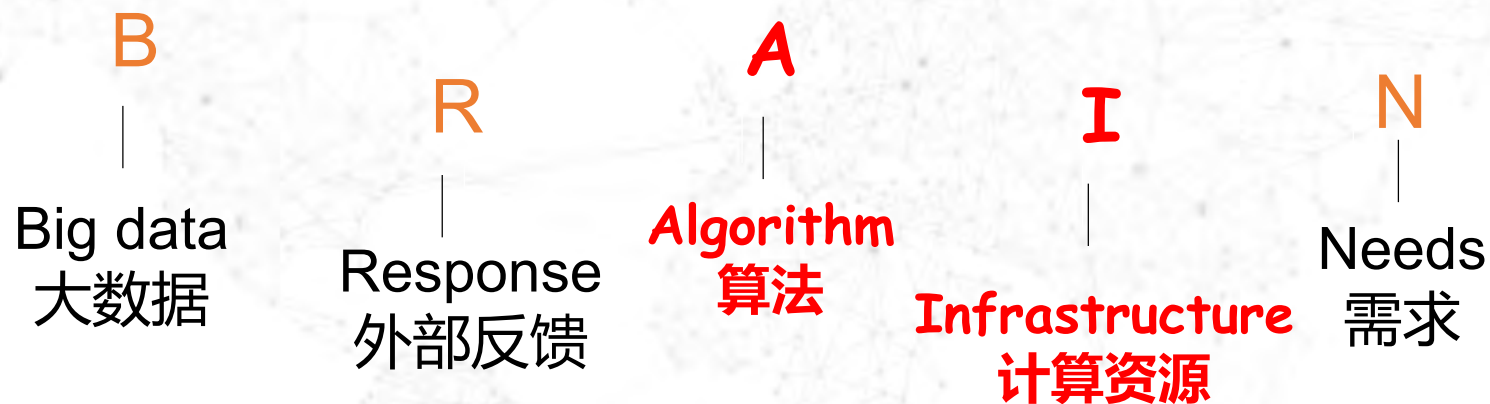
高效：企业需建立实时或准实时的数据采集传输、模型预测和响应决策能力



我们认为的AI是什么？



构建商业AI能力的五个要素



问题:

假设我是一个点餐Pad提供商，已经收集了食客的点餐数据，如何为每一位食客每次就餐提供AI一键点餐，以同时提升餐馆翻台率与食客满意度？

现实与理想的距离

因素	现实	理想
特征工程	建模人员进行少量特征工程探索	庞大的AI团队进行大规模特征工程探索
模型规模	几十到几千维度	几千万到数十亿维度
模型算法	采用神经网络反复炼丹，通过模型变化适应场景	采用大规模机器学习算法，通过特征工程适应场景
模型除错	经常出现穿越、过拟合等问题，线下建模效果很好，上线之后失望	老司机利用经验带领团队排除掉建模过程中的各种风险

如何使数据工程师变成AI专家？

- 特征工程：使数据工程师能够有效探索出足够有效的特征集
- 模型规模：引入一套支撑超高维模型训练的机器学习系统
- 模型算法：使数据工程师能够训练出足够有效的模型
- 模型除错：使数据工程师能够快速了解到模型是否有错误并加以排除

先知平台 – 敏捷AI应用构建平台



Prophet-Web化操作界面构建工业标准AI应用

The screenshot displays the Prophet-Web interface for building an AI application. The main workspace shows a workflow diagram with the following steps:

- Input: bankdata
- Process: 数据拆分 (Data Split)
- Parallel Processes: 4 instances of 特征抽取 (Feature Extraction)
- Parallel Processes: 特征重要性分析 (Feature Importance Analysis) and 自动特征组合 (Automatic Feature Selection, 60%)
- Parallel Processes: 逻辑回归 (Logistic Regression) and 逻辑回归自动调参 (Automatic Hyperparameter Tuning for Logistic Regression)
- Parallel Processes: 模型预测 (Model Prediction)
- Final Step: 模型评估 (Model Evaluation)

The left sidebar contains a navigation menu with categories like 公共数据, 项目数据, 数据源, 数据源, 样本数据, 项目模型, 数据处理, 特征工程, 分类算法, 聚类分析, 自定义脚本, 模型预测, and 模型评估.

The right sidebar shows configuration options for the '自动特征组合' (Automatic Feature Selection) step, including: 输入源 (Input Source), 训练集 (Training Set), 验证集 (Validation Set), 特征选择 (Feature Selection), 评分指标 (Evaluation Metric), 学习率 alpha系数 (Learning Rate alpha coefficient), 桶大小 (Bucket Size), 随机调参次数 (Number of Random Hyperparameter Tuning Iterations), 特征预排序配置 (Feature Pre-sorting Configuration), 特征池评测配置 (Feature Pool Evaluation Configuration), and 特征池大小 (Feature Pool Size).

The bottom status bar includes controls for 另存为 (Save As), 保存 (Save), 终止 (Stop), 异常校验 (Error Check), 删除计划 (Delete Plan), 运行历史 (Run History), 评估对比 (Evaluation Comparison), and a '运行中 00:06:59' (Running 00:06:59) indicator.

图形机器学习操作界面 - Lamma

The screenshot displays the Lamma graphical machine learning interface. On the left is a sidebar with a search bar and a list of categories: 项目数据, 项目模型, 数据拆分, 特征工程, 分类算法, 自定义脚本, 模型检测, and 模型评估. The main workspace contains a Directed Acyclic Graph (DAG) workflow starting with a 'data_in' node, followed by '数据拆分算子' and 'KQ算子'. The DAG branches into four '特征抽取算子' nodes, which lead to '逻辑回归', 'XGBoost', and two '模型训练' nodes. The workflow concludes with '模型评估算子' and '模型评估' nodes. A red box labeled '算子区' points to the sidebar. Another red box labeled 'DAG操作区' points to the central workflow diagram. A third red box labeled '参数配置区' points to the 'console编辑' panel on the right, which shows a list of parameters: 1 label=label(col_21), 2 a=log(col_10), and 3 b=Discrete(col_2). At the bottom, a red box labeled '计划操作区' points to a toolbar with icons for '删除', '保存', '停止', '清除', and '运行历史'.

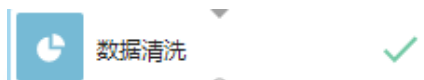
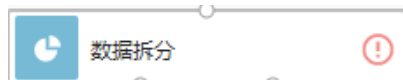
算子区

DAG操作区

参数配置区

计划操作区

Prophet – 简化数据预处理过程



数据拆分

拆分方式 ⓘ

按比例拆分数据

按规则拆分数据

先排序后拆分数据

随机拆分



随机种子 ⓘ

0

分层拆分 ⓘ



字段选择

输出结果1存储类型

压缩文件 (Parquet)

输出结果2存储类型

压缩文件 (Parquet)

▶ 备注

输入数据大小: [349.138672]KB

输出数据大小: [349.062500]KB

输入数据文件个数及每个文件大小: [5]

[_SUCCESS]: [0]B

[_common_metadata]: [1.809570]KB

[_metadata]: [6.844727]KB

[part-r-00000-08a2c1ad-78d3-4ca7-936d-6c16d0ab4e89.gz.parquet]: [170.694336]KB

[part-r-00001-08a2c1ad-78d3-4ca7-936d-6c16d0ab4e89.gz.parquet]: [169.790039]KB

输出数据文件个数及每个文件大小: [5]

[_SUCCESS]: [0]B

[_common_metadata]: [1.809570]KB

[_metadata]: [6.831055]KB

[part-r-00000-fc028169-6a84-4617-b580-ae9b26585687.gz.parquet]: [170.664063]KB

[part-r-00001-fc028169-6a84-4617-b580-ae9b26585687.gz.parquet]: [169.757813]KB

输入数据行数: [41188]

输出数据行数: [41188]

输入数据列数: [21]

输出数据列数: [21]

函数签名: [replace(job, "-", "")_1501837500085], 清洗配置: [job=replace(job, "-", "")]

总共数据条数: [41188]

成功处理数据条数: [10675]

匹配失效条数: [30513]

示例:

1, technician

2, management

3, services

4, retired

5, management

8, admin.

10, housemaid

11, management

12, management

13, services

Prophet- 简化特征工程

离散特征编码

特征组合函数

console编辑 ⓘ

☰ 验证 ⚙ 生成配置



暂无字段名和类型信息
可点击右上角图标查看
系统推荐字段名和类型

```
7 discrete_feature_36_5=discrete(emp_var_rate) # emp_var_rate
8 discrete_feature_36_6=discrete(month) # month
9 discrete_feature_36_7=discrete(contact) # contact
10 discrete_feature_36_8=discrete(poutcome) # poutcome
11 discrete_feature_36_9=discrete(previous) # previous
12 discrete_feature_36_10=discrete(pdays) # pdays
13 discrete_feature_36_11=discrete(job) # job
14 discrete_feature_36_12=discrete(default) # default
15 discrete_feature_36_13=discrete(marital) # marital
16 discrete_feature_36_14=discrete(age) # age
17 discrete_feature_36_15=discrete(loan) # loan
18 discrete_feature_36_16=discrete(education) # education
19 discrete_feature_36_17=discrete(day_of_week) # day_of_week
20 discrete_feature_36_18=discrete(housing) # housing
21 discrete_feature_36_19=discrete(campaign) # campaign
22 discrete_feature_36_20=discrete(combine(duration,pdays)) # duration pdays
23 discrete_feature_36_21=discrete(combine(duration,previous)) # duration previous
24 discrete_feature_36_22=discrete(combine(duration,poutcome)) # duration poutcome
25 discrete_feature_36_23=discrete(combine(default,duration)) # duration default
26 discrete_feature_36_24=discrete(combine(loan,duration)) # duration loan
27 discrete_feature_36_25=discrete(combine(contact,duration)) # duration contact
28 discrete_feature_36_26=discrete(combine(duration,pdays,previous)) # duration previous pdays
29 discrete_feature_36_27=discrete(combine(education,month)) # education month
30 discrete_feature_36_28=discrete(combine(campaign,nr_employed)) # campaign nr_employed
31 discrete_feature_36_29=discrete(combine(pdays,nr_employed)) # nr_employed pdays
32 discrete_feature_36_30=discrete(combine(duration,poutcome)) # duration poutcome pdays
33 discrete_feature_36_31=discrete(combine(loan,duration,pdays)) # duration loan pdays
34 discrete_feature_36_32=discrete(combine(marital,nr_employed)) # marital nr_employed
35 discrete_feature_36_33=discrete(combine(marital,education,month)) # marital education month
36 discrete_feature_36_34=discrete(combine(duration,education,previous,poutcome)) # duration previous poutcome pda
```

保存

确定

Prophet– 简化特征工程

支持两种特征编码方法：连续值特征和离散值特征

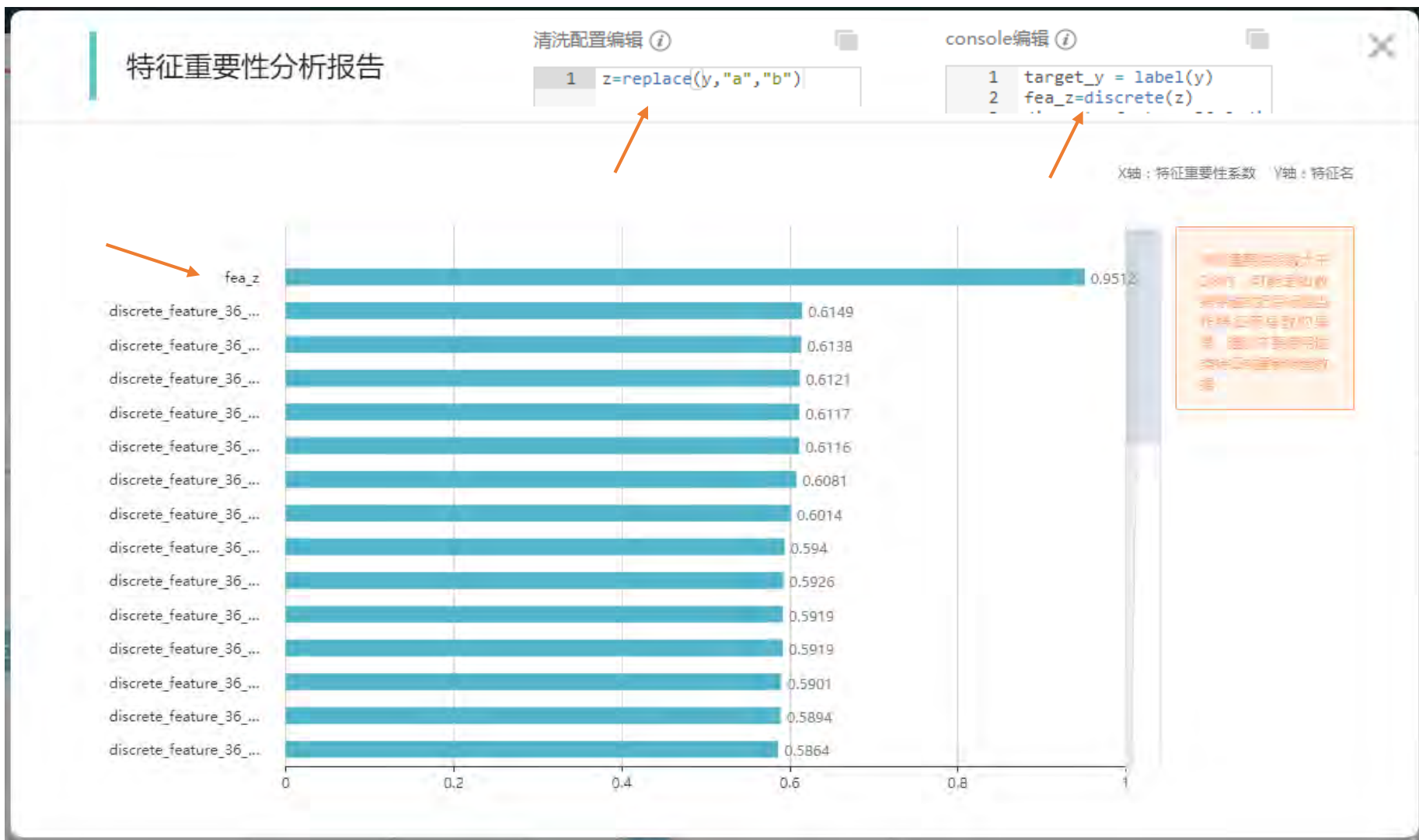
支持多种高维特征处理方法：

- Log/Floor/Lineartrans等数值处理
- Year/Hour/Minute/Second/Datediff/Timediff日期处理
- Eliminatechar/Split/Mapping/SplitbyKey等字符串处理
- Combine (组合) /Wordseg (切词)/Top (排序)

特征处理支持嵌套，例如：

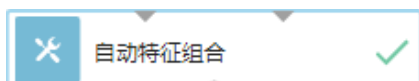
- `Y=top(int(splitbykey(age, ";", ";", ";"), "2")`

Prophet – 特征重要性分析，防止穿越



Prophet – 自动特征组合

`discrete_feature_36_26=discrete(combine(duration,pdays,previous)) # duration previous pdays`



自动特征组合

输入源 ①

- ① 训练集
- ② 验证集

特征选择 ①

- ④ 选择全部特征
- 自定义

评分指标 ①

slotwise_auc

学习率 alpha系数 ①

0.1

桶大小 ①

70 100 1000 10000 100000

随机调参次数 ①

0

特征预排序配置

样本采样率 ①

0.25

迭代轮数 ①

0

特征组合分析报告

特征数 35 显示统计评分后的输入特征配置

特征池序	特征抽取配置	slotwise_auc
1	discrete_feature_36_20=discrete(combine(duration,pdays)) # duration pdays	0.864894
2	discrete_feature_36_21=discrete(combine(duration,previous)) # duration previous	0.882880
3	discrete_feature_36_22=discrete(combine(duration,poutcome)) # duration poutcome	0.892774
4	discrete_feature_36_23=discrete(combine(default,duration)) # duration default	0.897527
5	discrete_feature_36_24=discrete(combine(loan,duration)) # duration loan	0.902156
6	discrete_feature_36_25=discrete(combine(contact,duration)) # duration contact	0.905734
7	discrete_feature_36_26=discrete(combine(duration,pdays,previous)) # duration previous pdays	0.909546
8	discrete_feature_36_27=discrete(combine(education,month)) # education month	0.910162
9	discrete_feature_36_28=discrete(combine(campaign,nr_employed)) # campaign nr_employed	0.910919
10	discrete_feature_36_29=discrete(combine(pdays,nr_employed)) # nr_employed pdays	0.912560
12	discrete_feature_36_30=discrete(combine(duration,pdays,poutcome)) # duration poutcome pdays	0.914066
14	discrete_feature_36_31=discrete(combine(loan,duration,pdays)) # duration loan pdays	0.915714
16	discrete_feature_36_32=discrete(combine(marital,nr_employed)) # marital nr_employed	0.915436
17	discrete_feature_36_33=discrete(combine(marital,education,month)) # marital education month	0.915305

discrete_feature_36_33=discrete(combine(marital,education,month)) # marital education month

Prophet – 自动参数探索

逻辑回归参数分析报告

调参方法： RandomSearch

调参次数： 1

算法模式： 速度优先

初始化值

最大训练轮数	学习率	L1正则项系数	L2正则项系数
4	0.5	0	0

调参运行记录

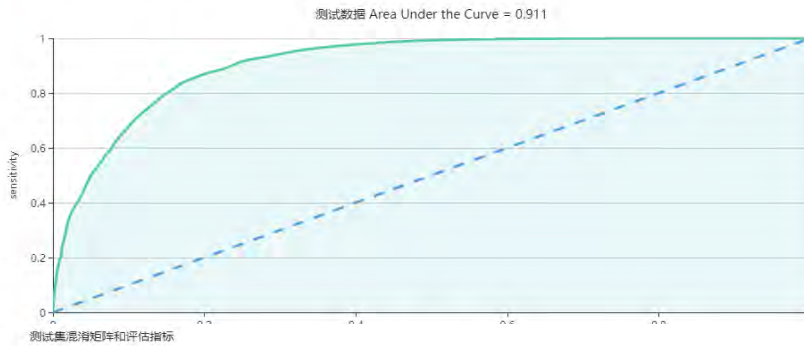
验证AUC	最大训练轮数	学习率	L1正则项系数	L2正则项系数	运行状态
0.854299843	6	4.768562353958775	29.92301754889432	15.335804343727975	成功

Prophet – 模型评估报告

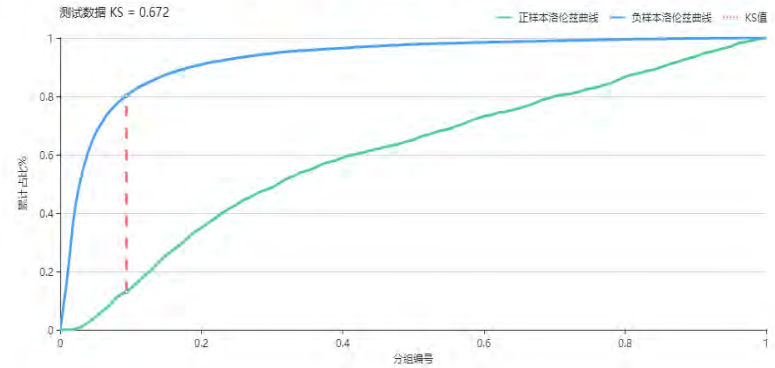
模型评估报告

测试数据评估结果

ROC图 Precision/Recall图 Lift图 K-S图 Gain图

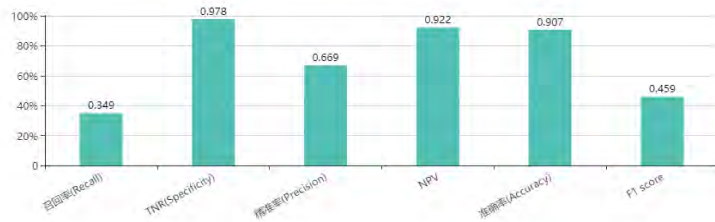


ROC图 Precision/Recall图 Lift图 K-S图 Gain图



正负例判定阈值

	模型预测 (正)	模型预测 (负)	总数
实际值 (正)	486 3.94%	905 7.34%	1391 11.28%
实际值 (负)	241 1.95%	10698 86.76%	10939 88.72%
总数	727 5.9%	11603 94.1%	



分段统计

	样本占比	真正例	伪正例	伪负例	真负例	准确率	召回	Specificity	本组AUC	累积AUC
1]	0.009	89	18	0	0	0.832	1.000	0.000	0.000	0.000
2]	0.011	97	34	0	0	0.740	1.000	0.000	0.000	0.000
3]	0.012	92	55	0	0	0.626	1.000	0.000	0.001	0.001
4]	0.012	95	56	0	0	0.629	1.000	0.000	0.001	0.002
5]	0.015	113	78	0	0	0.592	1.000	0.000	0.002	0.005
6]	0.018	0	0	85	140	0.622	0.000	1.000	0.005	0.010
7]	0.028	0	0	142	200	0.585	0.000	1.000	0.008	0.018
8]	0.049	0	0	193	414	0.682	0.000	1.000	0.022	0.040
9]	0.107	0	0	283	1037	0.786	0.000	1.000	0.073	0.113
10]	0.739	0	0	202	8907	0.978	0.000	1.000	0.798	0.911

Prophet – 预估服务发布

集群资源监控

API访问接口

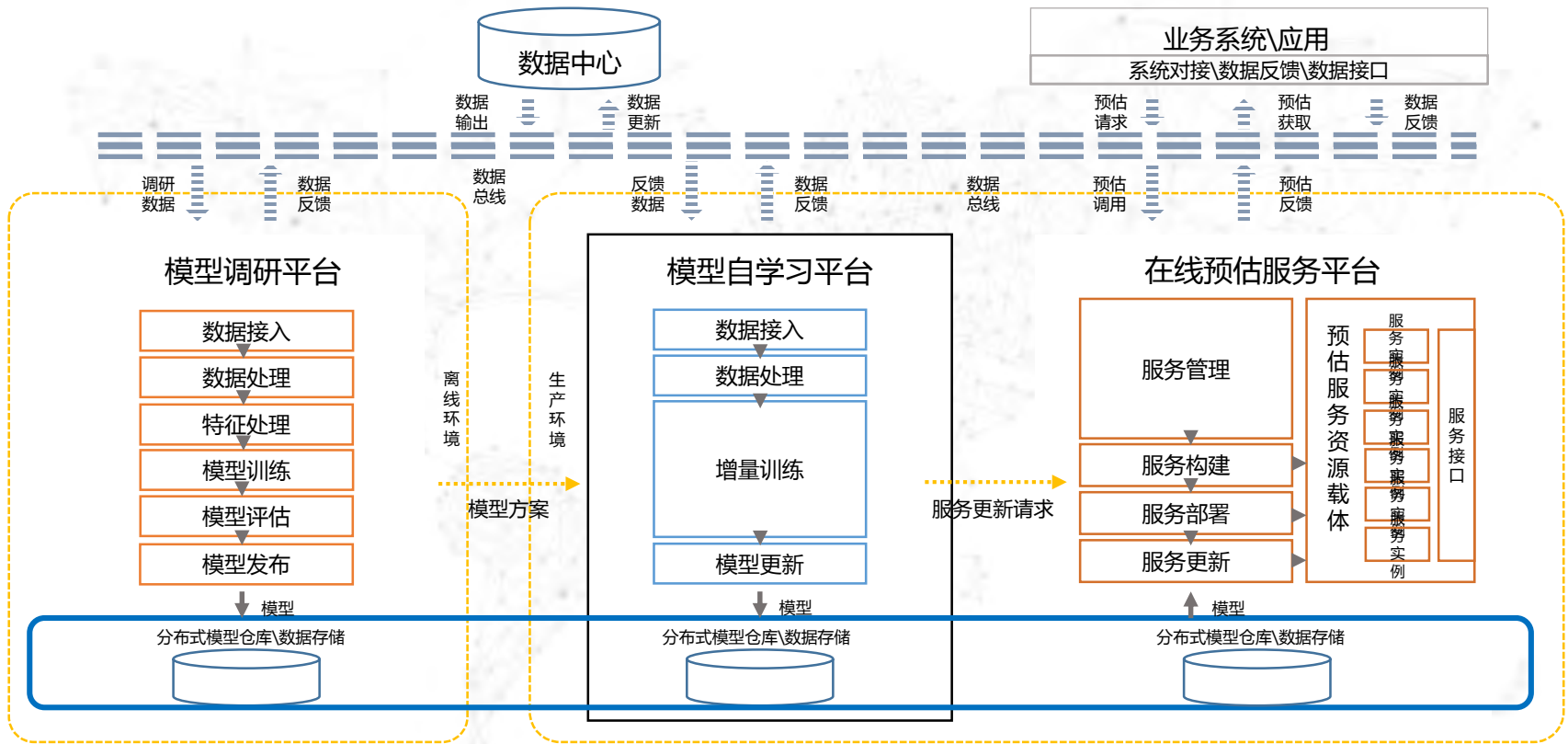
测试

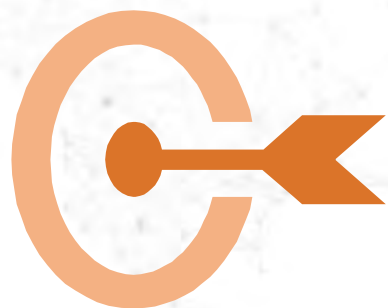
字段名	数据类型	输入值
duration	Int	210
emp_var_rate	Double	1.4
month	String	aug
contact	String	cellular
poutcome	String	nonexistent
previous	Int	0
pdays	Int	999
job	String	blue-collar



测试

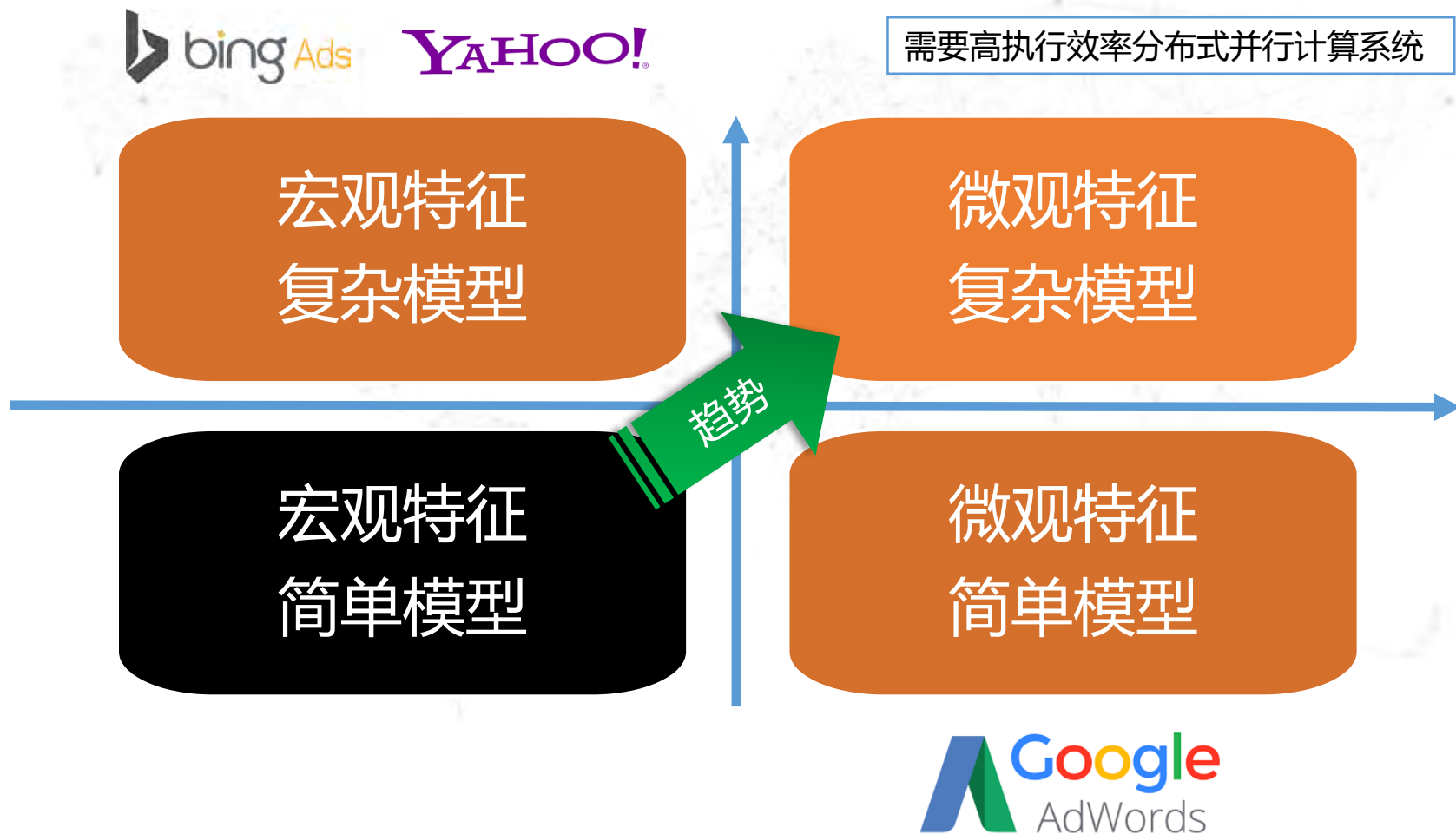
先知：平台化的机器学习架构支撑机器学习全流程平台





平台技术分享

机器学习算法在工业应用中的4个象限



高维模型：计算能力是第一生产力

线性模型：

- Logistic Regression
- SVM
- GMM
- LDA

非线性模型：

- GBDT
- DNN
- Kernel Methods



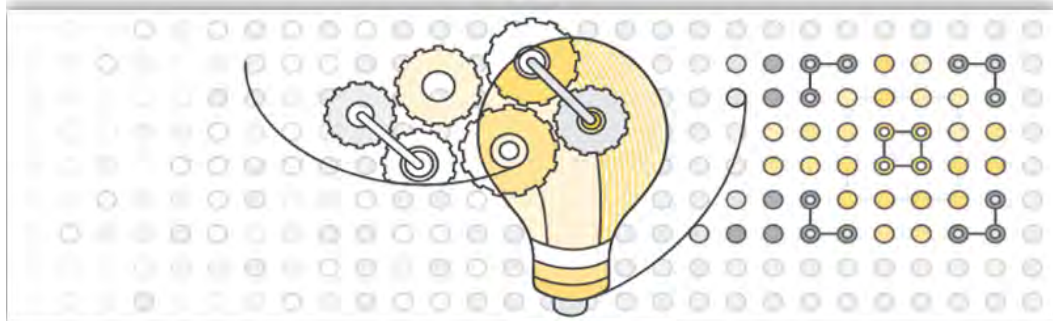
计算能力是第一生产力

VC维：

- 能被模型完全区分的最大数据集大小

PAC学习：

- 给定VC维下的算法泛化能力界



Theorem 2.9 (upper bound on sample complexity, [Blumer et al., 1989])

Let H and F be two function classes such that $F \subseteq H$ and let A an algorithm that derives a function $h \in H$ consistent with m training examples. Then, $\exists c_0$ such that $\forall f \in F, \forall D$ distribution, $\forall \epsilon > 0$ and $\delta < 1$ if

$$m > \frac{c_0}{\epsilon} \left(VC(H) \times \ln \frac{1}{\epsilon} + \frac{1}{\delta} \right)$$

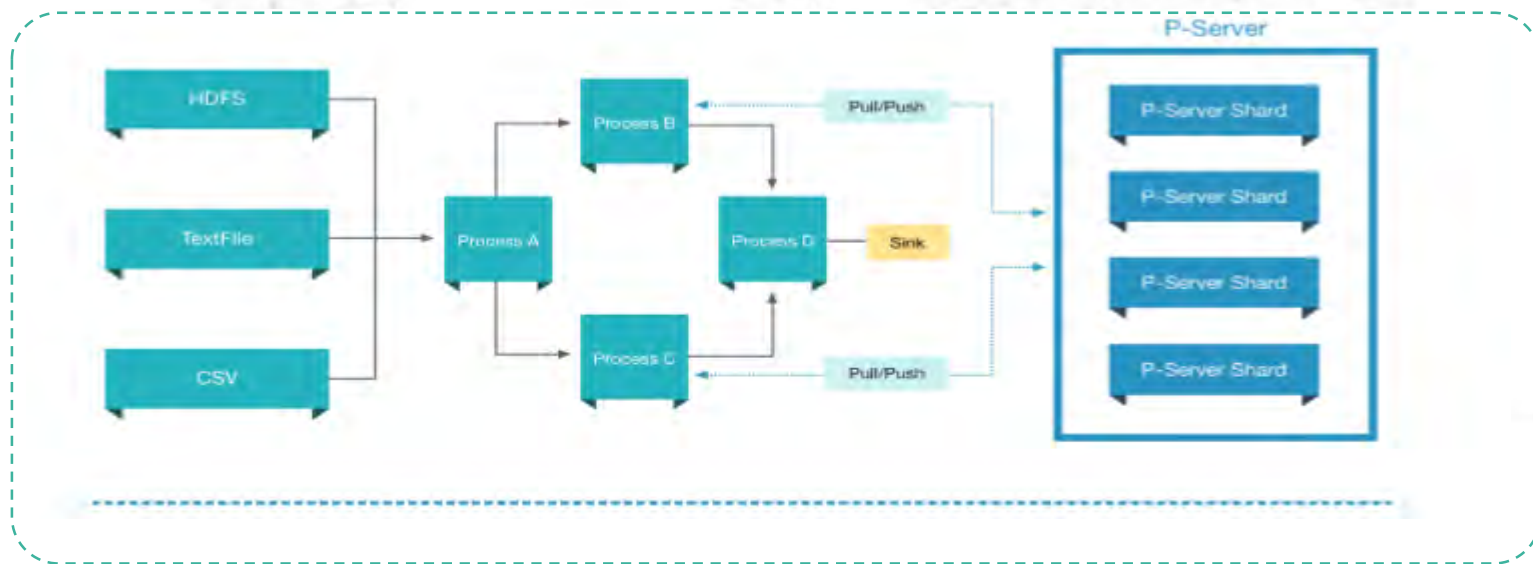
then with a probability $1 - \delta$,

$$error_D(h) \leq \epsilon,$$

where $VC(H)$ is the VC dimension of H and $error_D(h)$ is the error of h according to the data distribution D .

大规模机器学习框架GDBT

- GDBT是一个分布式机器学习框架，由C++ 14编写
- 机器学习过程抽象，隐藏分布式细节
- 兼容MPI/Yarn/Mesos/K8s等运行环境
- 实现了高性能的Parameter Server架构
- 优化多模型共同训练过程，大幅提升AutoML相关算法的性能



标准算法优化：HDLR（第四范式）Vs. 传统LR

	传统逻辑回归算法	大规模离散逻辑回归算法（第四范式）
特征维度	几十到几千	几千万到几十亿，甚至上万亿
数据兼容性	需要使用高饱和度数据	可以直接使用低饱和度稀疏数据，例如互联网数据。
样本数据抽样	样本数据抽样，只使用抽样出来的样本建模	无需抽样，使用全量样本建模

从上表中我们可以看到，虽然两者都叫逻辑回归，但在特征维度、数据兼容性和样本数据是否需要抽样上存在着巨大差异。

标准算法优化：GBDT（第四范式）Vs. 传统决策树

	传统决策树算法 (如Cart,C4.5算法)	传统集成学习决策树算法 (如SAS上随机森林算法等)	GBDT / HE-TreeNet (第四范式)
树的数量	单棵树	多棵树	多棵树
模型准确度	树过深容易过拟合， 刻画准确和过拟合难以兼得	用很多棵简单的树迭代，不容易 过拟合	用很多棵简单的树迭代，不容易 过拟合
样本数量	几百万	几百万到几千万	数亿甚至几百亿
输入特征	数千	数千到数万	没有限制，由平台节点规模而定
离散特征 使用能力	无法处理大规模离散特征	无法处理大规模离散特征	通过HE-TreeNet实现对 大规模离散特征的处理和使用

从上表中我们可以看到，先知平台大规模机器学习建模的GBDT算法和传统集成学习决策树算法都具有多棵树，可以兼顾模型准确度的要求和防止模型过拟合的要求，而从支持建模样本数量和输入特征数量上，都比传统集成学习决策树算法大大提升。

算法样例：HE-Treenet

高维离散嵌入树网络 (Hyper-dimension Ensemble Tree Net)

- 摘要

决策树做模型训练时，如果遇到高维离散特征，会建立非常复杂的树，这会导致训练变得很慢，同时很有可能会造成模型的过拟合。这时就需要通过某种方法把**离散特征转成连续特征**再做训练，HE-TreeNet便是解决这个问题的一种实现。

- 解决的问题

实际情况中，连续特征和离散特征同时存在；
最大程度的利用数据价值。

- 适用场景

数据多，样本间有时序关系，时间跨度大；
短时间内数据分布变化不大，近期样本分布
对label影响更大；
离散特征多，需要用连续值模型。

高性能特征工程框架 - 1

- 特征处理：将原始二维表转化为高维稀疏特征矩阵
 - 无筛选、大规模：高维机器学习特定面对的问题
- 目标：支持使用某种领域特定中间语言描述特征处理过程
- 技术选型：为兼容Spark/在线程序，选用JVM上的语言
- 性能问题
 - 特征处理占用机器学习计算过程的50%以上时间
 - 类与方法、字符串、堆内存分配
 - 运行图：公共子表达式、递归调用、动态类型判断

高性能特征工程框架-2

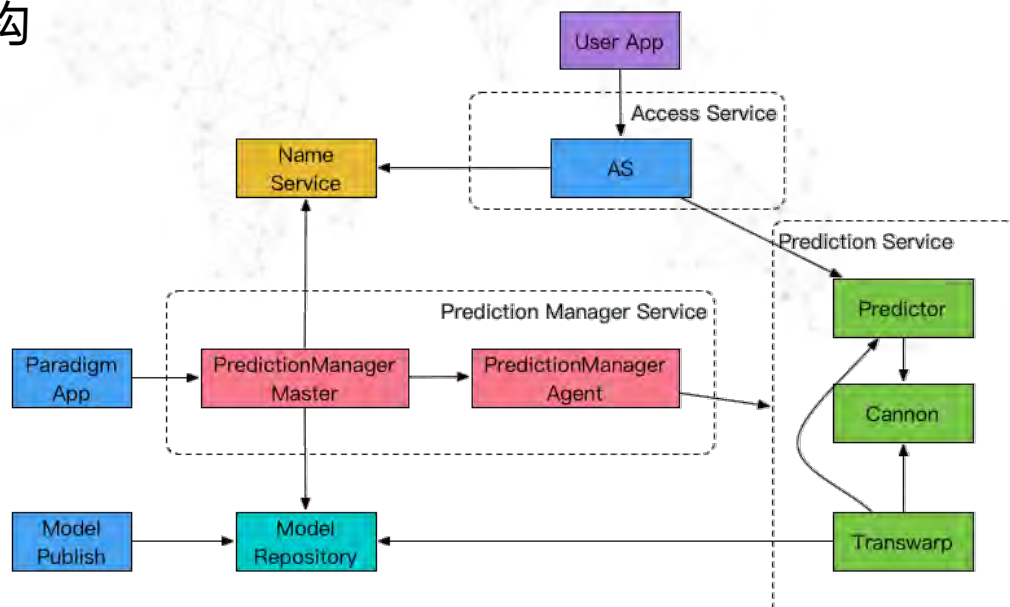
- 运行优化：AST遍历运行 – 运行时代码生成（脚本–DAG–代码）
- 死代码与公共子表达式消除
- 引入强类型系统与类型推断
- 通过使用原生类型数组 `int[]` `double[]` `char[]` 预先分配 显式管理内存
- Cache

Rtidb – 面向机器学习的特征蓄水池

- 解决核心问题：统计与序列特征的存储与查询
 - 反欺诈中查询账户最近一个月的交易记录，并衍生大量特征
 - 基于流式处理的预计算 VS 实时查询与特征衍生
 - 千万特征千条历史 实时交易反欺诈 c4.xlarge TP99 20ms QPS 400+
- 局部有序存储，满足高性能时序数据读取要求
- 并发读写友好，读写互不影响
- 支持TTL，内存回收时对读数据无影响
- 支持高级时序数据结构：CountWindow、TimeWindow、Session
- 高度定制的序列化协议

线上服务支撑组件

- 模型仓库：
 - 线下DAG图到线上DAG图的自动转换
 - 线上模型的滚动更新
- Cannon：分布式在线模型存储与访问
- 基于Kubernetes的高可用架构



workflow

- 区别于Tensorflow/GDBT的高层任务执行系统
- 支持通过图形界面进行定义与执行
- 可以从任意节点开始执行，支持中断、恢复
- 支持中间计算结果持久化（运行过的算子不再运行）
- 支持全局异常校验与推断（Schema-aware系统）



任务调度器

- 背景：工作流的实际执行计划包括多种任务
 - Python | Spark-ETL | GDBT | Tensorflow | Spark-MLLib
 - 执行环境：Local/Yarn/Kubernetes
- 多租户支持 – Quota、身份（User Mapping）、沙箱与安全
- 动态资源调度 – 如何给第一个上班的人分配资源
- 智能调参 – Cost Based VS Model Based
 - 使用先知来优化先知

适用于混合云的部署架构

- 先知整体支持容器化部署
- k7s = k8s – network – docker
- 实现一个符合容器标准的轻量级容器
- 基于IP/Port架构，避免Flannel等虚拟网络组件对企业内网架构的影响
- Prophet on Ingress = 基于域名的服务请求转发

总结

- 计算是第一生产力
 - 分布式、框架、算法优化
- 平台的目的是使计算生产力易于获取
 - GUI、DSL、调度、通用组件
 - AutoML



案例参考

第四范式先知平台成功案例（金融领域）

解决方案	典型业务场景	业务效果
风险控制	欺诈预防与侦测：利用机器学习建模技术，对潜在欺诈风险主体进行数据建模，发现欺诈意图，从而在欺诈交易发生前进行阻截，并通过自学习迭代自动发现新型欺诈手段和模式	某股份制银行信用卡中心 召回80%欺诈交易的前提下，准确率从1%提升至2.68% Visa渠道的准确率从1%提升至7.62%
精准营销	信用卡分期产品营销：预测客户对于分期产品（交易分期、取现分期、账单分析）的响应率，对高响应率的客户进行精准营销	某股份行信用卡中心：交易分期，响应率提升68%，收入提升61% 某股份行信用卡中心：取现分期，响应率提升22%，收入提升22%
	汽车贷款分期营销：在千万微信公众号客户中，挖掘近期有购车意向的客户，通过微信营销购车分期业务	某股份制银行信用卡中心：响应率提升200%~300%
个性化推荐	理财产品个性化推荐：预测客户对于不同理财产品的偏好，进行精准的产品推荐	某股份制银行：按不同资产段，响应率提升2倍~11倍，成交金额提升50%~500%
	内容分发个性化推荐：提供千人千面的效果和个性化体验	某新闻客户端产品：产品第一版本的模型比基线提高20%
智能客服	知识点个性化推荐：基于客户个性化属性及特征，预测并展示特定客户个性化问题；	项目进行中，目标为根据用户的行为进行实时预估，给出最有针对性和时效性的个性化知识推荐
	智能动态IVR菜单：通过对历史电话呼入纪录与处理数据建立机器学习模型，实现IVR渠道客户个性化动态菜单，精准命中客户来电需求	某股份制银行信用卡中心 原始基于规则的动态IVR菜单Top@5 准确率为25% 基于人工智能动态IVR菜单Top@3 准确率为53%
差异化定价	差异化定价：预测客户对于不同价格的响应率，实现精准的定价，实现总利润的提升	某股份制银行信用卡中心 客均手续费提升23%

第四范式先知平台成功案例（互联网领域）



某国Top1的新闻App 推荐，优化点击率 提升**34%**



某知识分享领域Top3 App 音频推荐，优化听完率 提升**43%**



某秀场类直播Top3 App 主播推荐，优化收看时长 提升**21%**



某国内最大的UGC社区 内容推荐，优化点击率 提升**93%**

运营小编专
家经验规则

用户
喜欢

用户
无感

机器学习
个性化推荐

机器学习模
型推荐

用户
喜欢

用户
无感



扫码免费体验先知平台：Prophet.4paradigm.com



THANKS

第四范式 胡时伟