

# 中国电子云&重庆大学

## 专业综合设计

云计算+大数据综合实战项目  
——用户精细化运营系统



中国电子云  
CECLOUD  
项目文档

2023. 7. 10

# 一、项目计划

## 1. 项目背景

随着互联网的发展和用户饱和度的提升，获客成本不断增加，因此，任何一家公司都需要做好用户的精细化运营，实现用户价值最大化。本项目适用于主营业务在线上进行的公司，如电商、教育、旅游、外卖等等。这类公司需要对用户的线上访问行为、消费行为、业务操作行为进行统计分析和数据挖掘，以支撑公司的业务运营、精准画像营销和个性化推荐等，从而提高业务转化率和改善公司运营效果。

本项目主要包含三大模块：数据仓库、用户画像系统和推荐系统。数据仓库作为数据的管理和运算中心，含有数据存档和各种统计、运算任务的核心平台；用户画像系统深入分析用户并打上各种规范标签，以支撑精细化运营；推荐系统根据用户画像及物品相似度推荐不同的物品，以改善用户体验并增加销量。该运营系统底层采用了云计算技术，搭建部署大数据平台环境，用来支撑海量数据的离线加实时计算，并结合多种算法模型进行用户画像分析，为推荐系统提供数据支撑。

## 2. 项目目标

- (1) 理解项目的业务背景、整体架构体系、技术及算法选型；
- (2) 掌握大数据平台环境搭建和组件安装配置；
- (3) 掌握数据仓库的设计和数准备，包括数仓建模、数据采集、数据清洗等；
- (4) 掌握埋点日志分析，包括流量分析、用户分析、留存分析、活跃分析、事件分析、路径分析、转化分析、广告分析等；
- (5) 掌握业务数据分析，包括销量分析、复购分析、订单分析、购物车分析等；
- (6) 掌握任务脚本开发和任务调度系统 Azkaban 的使用；
- (7) 掌握 Atlas 元数据管理系统，管理数据资产；
- (8) 熟悉机器学习与算法，包括相似度计算、KNN 算法、NaiveBayes 算法、朴素贝叶斯等；
- (9) 掌握用户画像分析方法，调用 SparkMLlib 算法库，完成品牌偏好、评论分类、流失预测、性别预测、风险预测等模型标签分析；
- (10) 掌握 OLAP 运营分析平台的技术框架选型及开发流程；

- (11) 掌握数据可视化方法，实现画像数据查询分析、用户分群查询分析、自定义标签查询分析等；
- (12) 学会数据分析报告撰写方法。

### 3. 项目计划

- (1) 理解项目的业务背景、整体架构体系、技术及算法选型；

时间		模块	课程主题	概述	工作内容
Day01	08:30-09:15	开班典礼	实训开班	开班典礼、领导致辞、纪律要求、学员分组、破冰环节	任务 1：了解项目背景、技术架构、业务流程、最终目标
	09:25-10:10	项目介绍	项目架构解析	项目业务背景介绍、实现功能、整体架构、技术及算法选型、15 天课程安排	任务 2：熟悉用户日志数据源
	10:20-11:05				任务 3：搭建部署 Hadoop 集群，并完成配置（至少 3 节点）
	11:15-12:00	数据介绍	数据源拆解	App/Web 埋点日志整体结构、事件类型、字段类型、UTM 广告跟踪、业务数据介绍	任务 4：搭建部署 Zookeeper 集群，并完成配置（至少 3 节点）
	13:30-14:15	大数据平台环境搭建	平台组件安装配置	完成 Hadoop 完全分布式、Spark、Hive、Kafka、Flume、Sqoop、HBase 等软件的安装配置。	任务 5：搭建部署 Kafka 集群，并完成配置（至少 3 节点）
	14:25-15:10				任务 6：搭建部署 Flume 集群，并完成配置（至少 2 节点）
	15:20-16:05				任务 7：编写 Flume 拦截器
	16:15-17:00				任务 8：生成模拟数据，观察 Kafka 消费者是否能消费到数据

Day02	08:30-09:15				任务 9：熟悉电商业 务数据表结构
	09:25-10:10				任务 10：安装部署 MySQL，并完成配置 (至少 2 节点)
	10:20-11:05				任务 11：生成业务模 拟数据
	11:15-12:00				任务 12：搭建部署 Maxwell，并完成配 置
	13:30-14:15				任务 13：测试业务数 据采集通道，观察 kafka 能否消费到数 据
	14:25-15:10				任务 14：编写 Flume 拦截器，生成模拟日 志数据采集到 kafka
	15:20-16:05				任务 15：搭建部署同 步工具 DataX，并完 成配置，将数据全量 同步到 hdfs
	16:15-17:00				任务 16：搭建部署 Flume，编写拦截 器，使用 Maxwell 工 具将业务数据增量同 步到 hdfs
Day03	08:30-09:15	数据仓库 -数据准 备	数据采集	APP 埋点日志预处理-流 程分析、代码实现、程 序打包部署上线运行、 数据加载到数仓 ODS 层；数据迁移工具 Sqoop 将用户信息表数	任务 17：搭建部署 Hive 集群，并完成配 置(至少 3 节点)
	09:25-10:10				任务 18：理解数据仓 库、数仓建模相关概 念
	10:20-11:05				任务 19：数据仓库 ODS 层开发

Day04				据、订单信息表数据导入 ODS	
	11:15-12:00		日志预处理	清洗过滤、数据解析、数据集成、数据修正、保存结果、预处理开发实现	任务 20：数据仓库 DIM 层开发
	13:30-14:15	数据仓库-埋点日志分析	数据加载	ODS 层数据加载、数据建模、模型，表结构、加载数据	任务 21：数据仓库 DWD 层开发
	14:25-15:10		明细构建	DWD 层建模及开发	任务 22：HiveSQL 编码实现——当日各渠道独立访客数
	15:20-16:05		流量分析	流量情况报表分析-方案设计、会话聚合表、用户聚合表、全局汇总表开发，多维度分析	任务 23：HiveSQL 编码实现——当日各渠道会话总数
	16:15-17:00				任务 24：HiveSQL 编码实现——当日各渠道会话平均浏览页面数
	08:30-09:15				任务 25：HiveSQL 编码实现——当日各渠道会话平均停留时长
	09:25-10:10		用户分析	日新增用户分析、日活跃用户分析	任务 26：HiveSQL 编码实现——当日各渠道跳出率
	10:20-11:05				任务 27：HiveSQL 编码实现——当日各小时新访客数
	11:15-12:00				任务 28：HiveSQL 编码实现——当日各小时页面浏览数
	13:30-14:15				任务 29：HiveSQL 编码实现——当日各

					小时新访客数
	14:25-15:10		留存分析	新用户留存、用户新鲜度分析	任务 30 : HiveSQL 编码实现——各类访客数量、页面浏览数
	15:20-16:05				任务 31 : HiveSQL 编码实现——各类访客页面浏览数
16:15-17:00	任务 32 : HiveSQL 编码实现——各类访客平均在线时长				
Day05	08:30-09:15		活跃分析	用户活跃成分分析、用户访问间隔分析、活跃用户留存分析	任务 33 : HiveSQL 编码实现——各类访客平均访问页面数
	09:25-10:10				任务 34 : HiveSQL 编码实现——当日各关键词评分
	10:20-11:05				任务 35 : HiveSQL 编码实现——当日回流用户数
	11:15-12:00				任务 36 : HiveSQL 编码实现——当日新增用户数
	13:30-14:15				任务 37 : HiveSQL 编码实现——当日活跃用户数
	14:25-15:10		App 分析	app 版本用户分布、app 版本升级分析	任务 38 : HiveSQL 编码实现——当日首页浏览人数
	15:20-16:05	任务 39 : HiveSQL 编码实现——当日商品详情页浏览人数			
	16:15-17:00	事件分析	交互事件分析主题-需求分析、整体建模方案设	任务 40 : HiveSQL 编码实现——当日加	

				计、代码实现	购人数
Day06	08:30-09:15		路径分析	访问路径分析需求解析、整体建模、sql 实现、访问路径概况统计	任务 41 : HiveSQL 编码实现——当日下单人数
	09:25-10:10				任务 42 : HiveSQL 编码实现——当日支付成功人数
	10:20-11:05				任务 43 : HiveSQL 编码实现——当日新增下单人数
	11:15-12:00				任务 44 : HiveSQL 编码实现——当日新增支付成功人数
	13:30-14:15		转化分析	转化漏斗分析、离线常规、在线定制	任务 45 : HiveSQL 编码实现——最近 7/30 日截至当前各品牌复购率
	14:25-15:10				任务 46 : HiveSQL 编码实现——当日各品牌订单数量、人数、金额
	15:20-16:05		广告分析	站内广告分析、站外广告分析	任务 47 : HiveSQL 编码实现——当日各品牌退单数量、人数
	16:15-17:00				任务 48 : HiveSQL 编码实现——当日各品类订单数量、人数、金额
Day07	08:30-09:15	数据仓库-业务数据分析	业务模型	业务表数据模型梳理、整体介绍、需求罗列、数据表导入策略	任务 49 : HiveSQL 编码实现——当日各品类退单数量、人数
	09:25-10:10				任务 50 : HiveSQL 编码实现——当日各 SPU 订单数量、人

					数、金额
	10:20-11:05		销量分析	成交额 GMV 分析	任务 51 : HiveSQL 编码实现——当日订单总额
	11:15-12:00		复购分析	复购率分析	任务 52 : HiveSQL 编码实现——当日订单数量
	13:30-14:15		订单分析	用户订单画像标签表计算、用户退拒行为画像标签表开发	任务 53 : HiveSQL 编码实现——当日退单人数
	14:25-15:10				任务 54 : HiveSQL 编码实现——当日各省份订单数量
	15:20-16:05				任务 55 : HiveSQL 编码实现——当日各省份订单金额
	16:15-17:00				任务 56 : HiveSQL 编码实现——当日优惠券金额
	08:30-09:15		购物车分析	用户购物偏好画像标签表开发	任务 57 : HiveSQL 编码实现——当日优惠券补贴率
	09:25-10:10				任务 58 : HiveSQL 编码实现——当日活动补贴率
	10:20-11:05	数据仓库-任务脚本	任务调度脚本开发	Spark 任务调度脚本开发、Sql 任务调度脚本开发、Sqoop 任务调度脚本开发	任务 59 : 编写 Spark 任务调度脚本
Day08	11:15-12:00				任务 60 : 编写 SQL 任务调度脚本
	13:30-14:15				任务 61 : 编写 Sqoop 任务调度脚本
	14:25-15:10		脚本中向外发送告	短信的方式、邮件方式	任务 62 : 配置脚本告



			警		警规则
	15:20-16:05	数据仓库-任务调度	任务调度系统 Azkaban	Azkaban 在项目中任务实战调度、项目 Hive 脚本调度、配置参数传递给任务的 Shell 脚本。	任务 63：部署并配置 Azkaban 工具
	16:15-17:00				任务 64：配置参数传递任务
Day09	08:30-09:15				任务 65：部署并配置 Atlas 元数据管理工具
	09:25-10:10				任务 66：元数据的注入
	10:20-11:05	数据治理	Atlas 元数据管理	Atlas 的整体概念解析、数仓元数据管理系统、元数据的注入、Atlas 的功能、安装部署和启动、配置连接 Kafka、配置连接 Hive、adminUI 功能学习-人工录入元数据实体、元数据搜索、分类管理和术语标签管理、工具调用演示、代码调用演示。	任务 67：配置连接 Kafka
	11:15-12:00				任务 68：配置连接 Hive
	13:30-14:15				任务 69：人工录入元数据实体、元数据搜索、分类管理和术语标签管理
	14:25-15:10				任务 70：工具调用、代码调用实现
	15:20-16:05	机器学习与算法	项目中涉及的算法模型	相似度计算案例代码实现、KNN 算法手撕代码、NaiveBayes 算法手撕代码、SparkMllib 算法库介绍-基础编程接口-Vector 向量使用、SparkMllib 算法库-朴素贝叶斯算法调用示例	任务 71：KNN 算法编码实现
	16:15-17:00				任务 72：NaiveBayes 算法编码实现
Day10	08:30-09:15				任务 73：朴素贝叶斯算法编码实现
	09:25-10:10				任务 74：SparkMllib 算法库介绍-基础编程接口-Vector 向量使用
	10:20-11:05	用户画像	品牌偏好	项目实战-品牌偏好度计	任务 75：项目实战-

	11:15-12:00			算-评语的情感分类	品牌偏好度计算-评语的情感分类
	13:30-14:15			项目实战-文本特征向量化算法-hashing 映射-TF-IDF 特征值	任务 76：项目实战-文本特征向量化算法-hashing 映射-TF-IDF 特征值
	14:25-15:10				
	15:20-16:05		评论分类	项目实战-评论语义情感分类-代码实现	任务 77：项目实战-评论语义情感分类-代码实现
	16:15-17:00				
Day11	08:30-09:15			项目需求-流失率标签计算-逻辑回归算法代码实现	任务 78：项目实战-流失率标签计算-逻辑回归算法代码实现
	09:25-10:10				
	10:20-11:05		流失预测	项目需求-行为性别预测-项目大会-需求、特征选择、算法选择-讨论分析	任务 79：项目实战-行为性别预测-用朴素贝叶斯算法实现
	11:15-12:00				
	13:30-14:15		性别预测	项目需求-行为性别预测-用朴素贝叶斯算法实现	任务 80：项目实战-风险预测-逻辑回归算法代码实现
	14:25-15:10				
	15:20-16:05				
Day12	16:15-17:00		风险预测	方案和流程都雷同“流失概率”标签的计算，只是特征选取不同	任务 81：用户画像标签更新
	08:30-09:15				
	09:25-10:10		用户画像标签更新	增量更新-全量更新-兴趣类标签的衰减问题	任务 82：部署并配置
	10:20-11:05				
	11:15-12:00		实时数据	OLAP 平台概述、平台	
	13:30-14:15				
	14:25-15:10				
Day13	15:20-16:05	数据精细			
	16:15-17:00				

		化运营分析	分析	需求示例、Presto 基本认识、集群部署、配置 Hive 数据源的连接及测试、对接 Hive 数据源的原理机制、配置 MySQL 数据源的连接及测试、对接 Kafka 数据、对接 Kafka-csv 格式数据、对接 Kafka-json 格式数据	Presto
	09:25-10:10	OLAP 平台开发			任务 83：配置 Hive 数据源的连接及测试
	10:20-11:05				任务 84：对接 Hive 数据源的原理机制
	11:15-12:00				任务 85：配置 MySQL 数据源的连接及测试
	13:30-14:15				任务 86：对接 Kafka 数据
	14:25-15:10				任务 87：对接 Kafka-csv 格式数据
	15:20-16:05				任务 88：对接 Kafka-json 格式数据
	16:15-17:00				
Day14	08:30-09:15		运营分析	技术框架选型、前后端分离开发、Web 页面框架搭建、前后端数据交互接口规范、流量概况数据查询接口、日活总数查询接口、日活维度数据查询接口、常规漏斗数据查询接口、自定义漏斗数据查询接口、流量概况数据查询服务。	任务 90：前后端数据交互接口规范
	09:25-10:10				任务 91：流量概况数据查询接口
	10:20-11:05				任务 92：日活总数查询接口
	11:15-12:00				任务 93：日活维度数据查询接口
	13:30-14:15		画像分析	画像数据查询分析、用户分群查询分析、自定义标签查询分析	任务 94：常规漏斗数据查询接口
	14:25-15:10				任务 95：自定义漏斗数据查询接口
	15:20-16:05				任务 96：流量概况数据查询服务
	16:15-17:00				任务 97：画像数据查

					询分析
Day15	08:30-09:15	项目总结	项目架构回顾	业务背景需求、项目整体技术架构、实现功能、数据分析报告撰写方法、答辩 PPT 准备、问题答疑	任务 98：用户分群查询分析
	09:25-10:10				任务 99：自定义标签查询分析
	10:20-11:05				任务 100：项目整体技术架构
	11:15-12:00				任务 101：项目实现功能
	13:30-14:15	项目答辩	分组项目答辩	PPT 分享、项目功能展示、文档展示、总结、评审	任务 102：解决问题
	14:25-15:10				任务 103：数据分析报告撰写
	15:20-16:05				任务 104：PPT 分享
	16:15-17:00				任务 105：项目功能展示

## 二、需求分析

### 1. 云计算模块

模块	功能需求	技术需求
虚拟化平台	提供资源管理和分配功能	OpenStack 搭建虚拟化平台
	支持虚拟机的创建、启动、停止和销毁	
	资源调度和弹性扩展功能	
容器化平台	实现快速部署和管理应用程序的能力	Docker 搭建容器化平台
		K8S 容器编排工具
大数据组件	提供分布式存储和计算能力	配置 Hadoop 集群
	支持大规模数据处理和分布式	配置 Spark 集群

	计算	
	提供数据仓库和 SQL 查询功能	配置 Hive
	实现高吞吐量的实时数据流处理	配置 Kafka
	收集、聚合和移动大量的日志数据	配置 Flume
	实现分布式数据库与关系型数据库之间的数据传输	配置 Sqoop
	提供高可扩展性和高性能的 NoSQL 数据库	配置 HBase
	将 MySQL 数据库的 binlog 解析为结构化 JSON 消息，并实现数据同步和实时更新	配置 Maxwell，解析 MySQL 数据库的 binlog 日志
高可用和容灾组件	在不同的数据中心部署主备节点	使用多数据中心备份
	实现服务的高可用性和容灾能力	使用 Kubernetes 的容器编排功能
	提供服务发现和负载均衡功能	使用 Kubernetes 的服务发现和负载均衡功能
数据备份和恢复组件	定期将数据备份到可靠的存储设备	使用可靠的存储设备
	实现数据备份和恢复的自动化	使用 Hadoop 的备份和恢复机制
安全性保障组件	配置身份验证和授权机制，保护数据平台免受未经授权访问	身份验证和授权机制的配置
	使用网络安全措施，如防火墙和访问控制列表，保护数据平台免受恶意攻击	防火墙和访问控制列表的配置

## 2. 大数据模块

**数据集中存储：**数据仓库应能够集中存储来自多个数据源的数据，并按照不同的层次和主题进行划分，以便于数据管理和使用。在本项目中，我们设计了 ODS 层、DWD 层、DWS 层、DIM 层和 ADS 层，每个层次又包含不同的主题，例如流量主题、用户主题、商品主题和交易主题。

**数据一致性和准确性：**数据仓库需要经过清洗、整合和转换处理，以提供高度一致和准确的数据。为实现一致性和准确性，我们采取每日增量表和全量表的策略。每日增量表记录每天新增的数据，而全量表包含完整的数据集。通过每日装载脚本对表进行自动维护，按照分层结构从底向上加载数据，确保全量数据和增量数据的同步。此外，采用覆写模式的数据库模式可以屏蔽不同步数据的影响，确保数据的一致性和准确性。

**快速查询和分析：**数据仓库应使用优化的结构和索引，以支持复杂的查询和分析操作。为提高查询性能，我们设计了高效的 SQL 查询语句，并避免不必要的表连接、子查询和复杂操作符。合理使用 WHERE 条件和 ORDER BY 子句，以及聚合函数和 GROUP BY 子句进行结果集汇总，也是提升数据仓库性能的关键。对于大规模数据，可以考虑进行数据分区或分表存储，以加快数据访问速度。

**决策支持：**数据仓库应提供全面的业务视图和洞察，以支持企业运营人员和研究学者做出明智的业务决策。为快速支持新建指标，仓库应具备敏捷决策能力。引入新指标并快速获取相关数据，可以帮助仓库管理层快速识别问题、监测绩效，并及时作出决策，以应对不断变化的业务环境。

**预测和优化：**数据仓库需要支持机器学习，以进行预测分析和业务优化。为此，数仓需要进行特征值的提前筛选和适当的特征工程，以为机器学习算法提供有意义的输入。数据清洗、转换、降噪和特征选择等操作是必要的。此外，数仓还应支持机器学习模型的训练和优化过程，包括选择合适的算法和模型架构，调整超参数，使用交叉验证等技术来提高模型性能。为支持大规模数据的训练和模型优化，数仓需要提供适当的计算资源和工具。

综上所述，数据仓库的需求分析包括数据集中存储、一致性和准确性、快速查询和分析、决策支持以及预测和优化。此外，还涉及到数据仓库的架构设计，如数据仓库汇总层、分析数据存储层、数据仓库明细层、维度层、操作数据存储层、数据安全层和数据治理层，用于支持数据的聚合、加工、存储、安全和质量管理。

---

## 3. 人工智能模块

### 3.1 用户画像分析

为了更好地理解用户，并提供更个性化、更精细化的服务，我们需要实现一个用户画像分析功能。这个功能将基于我们的数据仓库中的用户数据（包括用户的地理位置、性别、行为、购买

历史等) 进行分析, 使用 KNN 聚类算法来区分不同类型的用户, 并为每个用户创建一个个性化的用户画像, 并且计算各类用户的占比能够帮助分析用户类型, 从而采取不同的销售策略。

#### 功能需求:

- 数据处理: 此功能需要能处理大量的用户数据, 并能处理各种类型的数据, 包括数值、类别、文本等。它还需要能对数据进行清洗, 包括处理缺失值、异常值和重复值等。
- 用户分类: 此功能需要能将用户分为不同的类别, 如新用户、活跃用户、沉默用户和流失用户。此分类应基于用户的行为和购买历史, 例如, 新用户可能是那些刚注册或刚开始使用我们的服务的用户, 而活跃用户可能是那些经常在我们的平台上活动的用户。
- 用户画像创建: 此功能需要能为每个用户创建一个个性化的用户画像。用户画像应包含用户的基本信息(如地理位置、性别)和行为特征(如购买历史、活动频率等)。
- 用户画像更新: 随着用户行为的变化, 用户画像也应进行相应的更新。此功能需要能定期或者根据特定事件(如购买行为)来更新用户画像。

#### 其他需求:

- 性能需求: 由于我们可能需要处理大量的用户数据, 此功能需要具有较高的处理速度和数据吞吐量。同时, 用户画像的更新应能在短时间内完成, 以确保我们能及时根据用户画像提供服务。
- 可用性需求: 此功能的使用应尽可能简单, 以便我们的工作人员能快速理解和使用。此外, 它还应提供易于理解的用户画像报告, 以便我们的工作人员能快速理解用户的行为和需求。
- 安全性需求: 用户数据是敏感的, 此功能需要保护用户数据的安全, 防止数据泄露。所有的数据处理和存储都应符合相关的数据保护法规。

通过实现用户画像分析功能, 我们可以更深入地理解我们的用户, 为他们提供更个性化的服务, 从而增加用户满意度, 提高用户留存率, 最终提高业务收入。

## 3.2 用户性别预测

用户性别是用户画像中的重要信息之一, 对于许多业务场景和推荐系统都具有重要意义。通过准确地预测用户性别, 企业可以更好地了解用户需求, 进行精准的市场定位和个性化推荐。因此, 设计和实现一个用户性别预测功能对于提升业务效果和用户体验至关重要。

#### 功能需求:

##### 1. 数据输入和集成:

- 系统应该能够从数据仓库中获取用户的购买数据, 包括品类、品牌、购买单数和消费总金额等信息。
- 数据集成应该是实时或定期的, 以确保用户性别预测功能的数据始终是最新的。

##### 2. 特征提取和选择:

- 系统应该能够从用户购买数据中提取特征，如最近 30 天内购买最多的品类、第二多的品类、第三多的品类，最多的品牌、第二多的品牌、第三多的品牌，购买单数和消费总金额等。
- 特征选择应该基于数据分析和业务需求，选择具有较高区分度和预测性别能力的特征。

### 3. 数据预处理：

- 对于连续型特征（购买单数和消费总金额），系统应该提供数据标准化或归一化的功能，确保不同特征之间的数值范围相同。
- 对于离散型特征（品类和品牌），系统应该提供特征编码的功能，如独热编码，将其表示为适用于机器学习算法的特征向量。

### 4. 模型训练和优化：

- 系统应该提供朴素贝叶斯算法的训练功能，根据训练数据集计算先验概率和条件概率，构建预测模型。
- 在模型训练过程中，系统应该提供模型性能评估的功能，包括准确率、精确率、召回率等指标，以便进行模型的优化和改进。

### 5. 模型预测和输出：

- 系统应该能够根据新用户的购买行为提取相应特征，并使用训练好的模型进行性别预测。
- 预测结果应该是用户性别的概率分布，以便系统可以根据具体需求进行灵活的结果处理。

### 安全和隐私需求：

- 在设计用户性别预测功能时，需要考虑用户数据的安全和隐私保护。
- 系统应该采取必要的安全措施，如数据加密、访问控制等，确保用户数据的安全性和保密性。

## 3.3 用户购买意向分析

用户购买意向分析功能的需求分析旨在深入理解该功能的目标和业务需求，以便更好地设计和实现相应的解决方案。以下是对用户购买意向分析功能进行需求分析的一些建议，包括功能目标、用户需求、数据需求和业务规则等方面。

### 1. 功能目标：

- 预测用户的购买意向：该功能的主要目标是根据用户在数据仓库中的相关数据，通过机器学习算法预测用户的购买意向。预测结果可以作为营销决策和个性化推荐的依据，帮助提升用户购买转化率和用户体验。

### 2. 用户需求：

- 提供个性化推荐：用户希望平台能够根据他们的购买意向提供个性化的产品推荐，以减少



搜索和浏览时间，提高购买效率。

- 接收优惠和促销信息：用户希望能够及时了解到与其购买意向相关的优惠和促销信息，以享受更好的购物体验和价值。

- 获取个性化购买建议：用户希望能够获得基于其购买意向的个性化购买建议，例如搭配商品、相似商品等，以辅助购买决策。

### 3. 数据需求：

- 用户行为数据：需要收集和分析用户在平台上的行为数据，如浏览记录、搜索记录、购物车数据等。这些数据可反映用户的兴趣、偏好和购买行为。

- 商品数据：需要获取商品的相关信息，如类别、属性、销售数据等。这些数据可用于分析用户对不同商品类别的兴趣和购买倾向。

- 营销数据：需要获取营销活动的相关数据，如优惠券使用情况、活动参与情况等。这些数据可用于分析用户对促销活动的响应和购买行为的影响。

### 4. 业务规则：

- 定义购买意向度量指标：根据业务需求，可以定义购买意向的度量指标，如购买概率、购买金额等。这些指标可用于量化和比较用户的购买意向程度。

- 制定个性化推荐策略：基于购买意向分析的结果，制定个性化推荐策略，如基于用户相似度的协同过滤、基于内容的推荐、热门商品推荐等。这些策略可提高用户购买的相关性和满意度。

- 优化营销决策：将购买意向分析结果应用于营销决策，如优化促销活动的目标用户群体、优化优惠券的发放策略等，以提高营销活动的效果和回报率。

### 5. 需求限制：

- 数据可靠性和隐私保护：确保用户数据的可靠性和隐私保护，遵守相关法律法规和隐私政策，采取必要的数据安全措施。

- 实时性和准确性：在分析用户购买意向时，需要尽可能保持数据的实时性和准确性，以提供及时和准确的预测结果。

- 系统性能和可扩展性：随着用户数量和数据规模的增加，系统需要具备良好的性能和可扩展性，以支持大规模数据处理和快速响应。

## 3.4 售后服务满意度分析

以下是针对售后服务满意度分析功能的需求分析示例：

### 1. 功能性需求：

- 文本情感分类：系统应能够对用户对售后服务的评价文本进行情感分类，将其归类为积极、中性或消极等情感类别。

- 情感强度分析：系统应能够分析用户评价文本中的关键词，并计算出文本的情感强度得分，用于评估用户对售后服务的情感倾向。

- 综合评估：系统应能够将情感分类和情感强度分析的结果综合起来，得出对售后服务的满意度评估。

## 2. 数据需求：

- 用户评价数据：系统需要获取和存储用户对售后服务的评价数据，包括评价文本、评分、投诉数量、退单件数、退单金额、评价时间等。
- 关键词数据：系统可能需要使用情感词典或自行构建的关键词数据，用于情感强度分析的关键词匹配和情感得分计算。

## 3. 用户界面需求：

- 数据展示：系统应提供一个用户界面，用于展示售后服务满意度分析的结果，如满意度评分、情感分类结果、关键词的情感得分等。
- 可视化图表：界面应支持以图表、表格等形式呈现分析结果，以使用户能够直观地理解和分析数据。
- 数据筛选和深入分析：界面应提供交互式功能，使用户能够根据不同的维度进行数据筛选和深入分析，如按时间、产品类别等进行筛选。

## 4. 性能需求：

- 实时性：系统对售后服务评价数据的处理和分析应具有较高的实时性，以便及时了解和响应用户的反馈。
- 处理能力：系统需要能够处理大规模的用户评价数据，以保证分析的准确性和可靠性。
- 响应时间：系统的响应时间应在合理范围内，以提供良好的用户体验。

## 5. 安全性需求：

- 数据保护：系统应具备适当的安全措施，保护用户评价数据的隐私和机密性，遵循相关的数据保护法规。
- 访问控制：系统应设置访问控制机制，确保只有授权用户能够访问和使用售后服务满意度分析功能。

## 6. 可维护性需求：

- 可扩展性：系统的架构和设计应具备良好的可扩展性，以方便未来对功能和性能的扩展和改进。
- 可更新性：系统的模型和算法应能够进行定期更新和重新训练，以保持模型的准确性和适应性。

# 3.5 产品销售趋势分析

下面是产品销售趋势分析功能的需求分析：

## 1. 数据收集与预处理：

- 系统能够从数据仓库中获取每个月的产品销售额数据。
- 系统能够获取历史数据计算得出的季节性指数。
- 系统能够获取产品类别在总销售额中的占比。
- 系统能够获取最近 7、30 日的复购率。

- 系统能够获取最近 1、7、30 日的订单数、订单人数、退单数、退单人数和订单总额数据。

## 2. 数据清洗与转换：

- 系统能够对获取的数据进行清洗，处理缺失值、异常值和重复值。
- 系统能够将数据转换为适合时间序列分析的格式，例如将日期作为时间索引。

## 3. 季节性指数计算：

- 系统能够基于历史数据计算季节性指数，以反映产品销售在不同季节的波动情况。

## 4. 时间序列分析算法选择：

- 系统能够提供多种时间序列分析算法供选择，例如 ARIMA 模型和指数平滑模型。

## 5. 模型训练与参数调优：

- 系统能够根据历史数据训练选定的时间序列分析模型，并通过参数调优来提高模型的准确性和预测能力。

## 6. 趋势分析与预测：

- 系统能够分析产品销售趋势并进行预测，提供未来时间段内的销售额预测和趋势分析。

## 7. 复购率分析：

- 系统能够计算最近 7、30 日的复购率，即在特定时间段内再次购买同一产品的比例。

## 8. 订单数、订单人数、退单数、退单人数分析：

- 系统能够分析最近 1、7、30 日的订单数、订单人数、退单数和退单人数的变化趋势。

## 9. 订单总额分析：

- 系统能够分析最近 1、7、30 日的订单总额的变化趋势。

## 10. 可视化与报告输出：

- 系统能够将分析和预测结果可视化展示，例如绘制销售趋势图、复购率曲线、订单数和退单数的变化图等。
- 系统能够生成相应的报告输出，以便用户更好地理解 and 利用销售趋势分析的结果。

# 3.6 商品热度分析

1. 商品热度评估：希望通过该功能对商品的热度进行评估和分类，以了解商品的受欢迎程度和销售潜力。

2. 数据源：功能需要基于数据仓库中的数据进行分析。需要明确数据仓库的结构和数据的可用性，确保能够获取到所需的数据，包括销售增长率、商品页面访问量以及各项统计数据。

3. 特征选择：需要确定用于评估商品热度的特征。根据您提供的信息，可以选择销售增长率、商品页面访问量以及最近 1、7、30 日的复购率、订单数、订单人数、收藏数、收藏人数、退单数和退单人数等作为特征。
4. 数据预处理：需要对从数据仓库中获取的数据进行预处理，包括处理缺失值、异常值和数据格式转换等。确保数据的质量和可靠性，以提高分析结果的准确性。
5. 聚类算法选择：基于需求，选择 KNN 聚类算法用于商品热度分析。KNN 算法可以根据商品的特征将其归类为不同的热度级别，从而帮助我们了解商品的受欢迎程度和销售潜力。
6. 热度级别定义：需要定义不同热度级别的标准。例如，可以将商品划分为高热度、中热度和低热度等级别，以便更好地理解 and 比较商品之间的热度。
7. 分析结果呈现：需要将商品热度分析的结果以可视化的方式呈现出来，以使用户更直观地理解和分析数据。可以使用散点图、热力图、雷达图等方式展示不同热度级别下的商品特征。
8. 用户界面和交互：考虑到功能的实际使用场景，可以设计一个用户界面，使用户能够方便地输入参数、查看分析结果，并根据需求进行交互操作，如筛选特定商品类别或时间段的热度分析结果等。
9. 实时性：需求是否需要实时更新商品热度分析结果，或者是基于历史数据进行分析。这将决定功能的数据更新频率和实时性要求。
10. 可扩展性：需求是否要求功能能够处理大规模的数据集和高并发请求。如果需要，可能需要考虑采用分布式计算或其他技术手段来提高功能的性能和扩展性。

通过对这些需求进行详细分析，可以确保商品热度分析功能满足用户的期望，提供准确、可靠和有用的分析结果，帮助用户做出更明智的决策和制定有效的营销策略。

### 3.7 优惠券效果分析

#### 1. 功能概述：

该功能旨在通过回归分析算法对优惠券的效果进行分析，以评估优惠券对用户消费行为的影响。它将基于数据仓库中的相关数据，包括优惠券信息、优惠券使用记录、用户消费历史、用户行为记录、订单价格等，计算和生成与优惠券效果相关的指标和报告。

#### 2. 输入：

- 优惠券信息：包括优惠券的类型、折扣金额或折扣比例等。
- 优惠券使用记录：记录哪些用户使用了哪些优惠券，以及使用时间等。
- 用户消费历史：包括用户的购买记录、订单价格等。

- 用户行为记录：记录用户的点击、浏览、加入购物车等行为。
- 数据仓库连接信息：用于连接和获取相关数据的数据库信息。

### 3. 输出：

- 优惠券使用率：计算并展示优惠券的使用率，即实际使用优惠券的比例。
- 平均优惠价格：计算并展示每个订单中优惠券的平均抵扣金额。
- 优惠券转化率：计算并展示通过优惠券促成的购买行为的比例。
- 订单价格与优惠价格关联分析：分析订单价格与应用优惠券后的价格之间的关系，并生成相关报告。
- 优惠券效果报告：生成综合的优惠券效果报告，包括上述指标的统计数据、图表和可视化结果。

### 4. 数据处理和分析：

- 数据清洗：对输入的优惠券信息、使用记录、用户消费历史、用户行为记录等数据进行清洗，处理缺失值、异常值和重复值，确保数据的准确性和一致性。
- 特征工程：从用户消费历史和行为记录中提取用户特征，从优惠券信息中提取优惠券特征，构建相关特征用于回归分析。
- 回归分析：使用合适的回归算法，基于特征和目标变量进行回归分析，得出回归系数和模型结果。
- 结果解释和评估：解释回归模型结果，评估模型的拟合程度和预测效果，例如通过均方误差（MSE）和决定系数（R-squared）等指标进行评估。

### 5. 可视化和报告生成：

- 结果可视化：将计算得到的指标和分析结果通过图表、统计摘要等形式进行可视化展示，使用户能够更直观地理解优惠券的效果。
- 报告生成：根据分析结果生成综合的优惠券效果报告，包括各项指标的统计数据、趋势分析和图表等，以使用户进行业务决策和优化。

### 6. 用户界面和交互：

- 用户登录和权限管理：提供用户登录功能，并根据用户的权限控制对数据和功能的访问。
- 数据选择和筛选：允许用户根据需求选择和筛选特定的优惠券数据和时间范围。
- 报告导出和分享：提供报告导出功能，允许用户将分析结果导出为常见的文件格式（如PDF、Excel），方便分享和进一步分析。

## 3.8 用户流失预测

### 1. 功能目标：

- 提供用户流失预测：根据用户行为数据、用户个人信息和用户反馈数据，预测用户是否会流失。
- 辅助决策制定：为企业提供用户留存策略和决策支持，以减少用户流失。

## 2. 数据源:

- 用户行为数据: 包括用户的活动记录、访问频率、购买历史等。
- 用户个人信息: 包括用户的性别、年龄、地理位置等。
- 用户反馈数据: 包括用户的满意度调查结果、投诉反馈等。
- 服务使用数据: 包括用户的使用情况、服务记录等。

## 3. 功能要求:

- 数据预处理: 对原始数据进行清洗、特征选择、特征编码和数据划分等预处理步骤, 以准备数据用于建模。
- 特征工程: 提取有意义的特征, 包括用户行为特征、个人信息特征和反馈特征, 并进行特征缩放、构建和降维等处理。
- 模型训练: 选择适当的算法 (如支持向量机、随机森林等) 进行模型训练, 通过训练集优化模型参数。
- 模型评估: 使用测试集评估模型的性能, 计算准确率、召回率、F1 值等指标。
- 用户流失预测: 对新用户的数据进行预测, 判断其是否有流失的趋势。
- 结果输出: 向企业提供用户流失预测结果, 可能是二元分类 (流失/非流失) 或概率预测值。
- 可视化展示: 可视化用户流失预测结果和相关指标, 以便企业更直观地理解数据和模型效果。

## 4. 性能要求:

- 模型准确性: 需要建立准确可靠的用户流失预测模型, 以便为企业提供有价值的决策支持。
- 模型实时性: 如果需要对实时数据进行流失预测, 模型需要能够快速处理数据并给出预测结果。
- 系统稳定性: 功能应该在长时间运行和大规模数据处理下保持稳定和可靠。

## 5. 安全与隐私:

- 数据安全: 确保用户数据和敏感信息的安全存储和处理, 符合相关隐私法规和公司政策。
- 数据访问权限: 确保只有授权人员可以访问和使用用户数据, 限制敏感信息的访问权限。

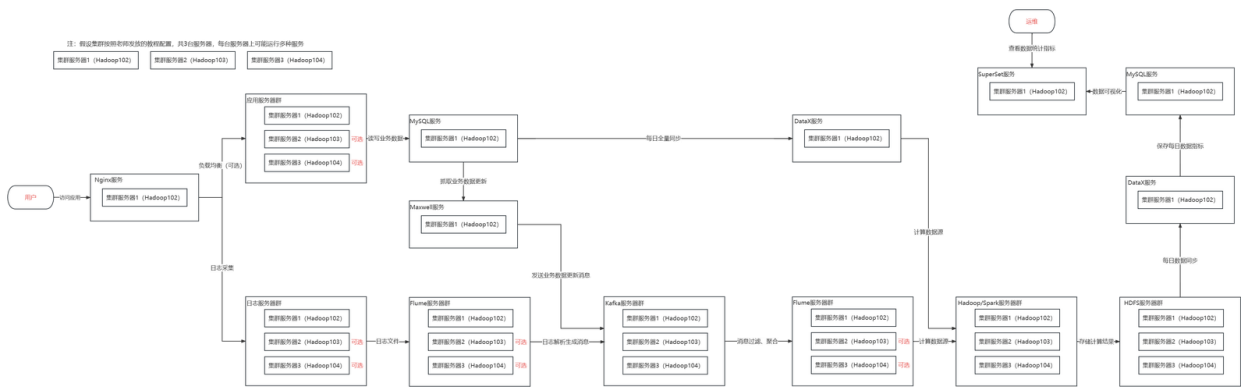
## 6. 部署与集成:

- 可扩展性: 能够处理不断增长的数据量和用户规模, 支持系统的水平扩展和性能优化。
- 数据仓库集成: 与现有数据仓库和数据源进行无缝集成, 获取所需的用户行为数据、个人信息和反馈数据。
- 可部署性: 能够将用户流失预测功能部署到生产环境, 并与企业现有系统和 workflows 集成。

### 三、系统设计

#### 1. 云计算模块

##### 1.1 云计算整体设计



#### 2. 大数据模块

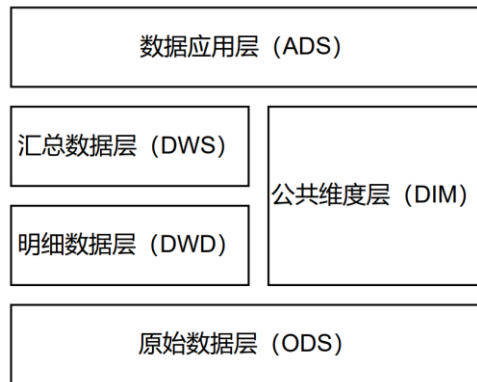
##### 2.1 数据仓库

一个优秀而可靠的数据仓库体系需要一个良好的数据分层结构，通过合理的分层可以使数据体系更加清晰，从而简化处理复杂问题的过程。合理的数据分层可以带来多重好处。首先，它可以将数据按照其特征和用途进行组织，使得数据的管理和查询更加高效。其次，通过在不同层次上进行数据的汇总和聚合，可以提供不同粒度的分析和报告，以满足不同用户的需求。此外，合理的数据分层还可以提供更好的数据安全性和可维护性，通过将敏感数据隔离在更高层次或更受保护的区域，确保数据的安全性。最后，合理的数据分层还可以支持数据生命周期管理，使得数据的存储和维护成本更加可控。

在设计数据分层结构时，需要考虑多个因素。首先，需要根据业务需求和数据特点确定合适的分层标准，例如按照主题、时间、地理位置等进行分层。其次，需要平衡数据的规模和复杂度，避免过多的层次导致管理和维护难度增加。此外，还需要确保不同层次之间的数据流畅和一致性，以便实现数据的无缝集成和交互。最后，定期审查和更新数据分层结构，以适应业务的变化和发展。

总之，一个良好的数据分层结构是构建优秀可靠的数仓体系的关键要素之一。通过合理的分层，可以提高数据管理和分析的效率，简化复杂问题的处理，同时也为数据安全和可维护性提供了保障。

以下是该项目的分层规划。



在数据仓库体系中， 以下是各个层次的详细说明：

- 原始数据层（ODS）：作为数据仓库的数据准备区，存放未经过处理的原始数据，并且保持与源系统结构的一致性。该层负责数据的抽取、清洗和初步处理，以确保数据的质量和完整性。
- 公共维度层（DIM）：基于维度建模理论进行构建，存放维度模型中的维度表，保存了一致性的维度信息。该层提供了可共享的维度数据，例如时间、地理位置、产品等，以支持不同业务领域的统一分析和报告。
- 明细数据层（DWD）：基于维度建模理论进行构建，存放维度模型中的事实表，保存各业务过程的最小粒度操作数据。该层提供了详细的事实数据，例如交易记录、用户行为等，以支持深入的分析和挖掘。
- 汇总数据层（DWS）：基于上层的指标需求，以分析的主题对象作为建模驱动，构建公共统计粒度的汇总表。该层根据业务需求进行数据汇总和聚合，提供了更高层次的统计指标结果，以支持决策和管理层面的分析。
- 数据应用层（ADS）：存放各项统计指标结果，是最接近用户需求的数据层。该层根据具体的业务需求构建各种报表、指标和分析结果，提供了对数据仓库的最终应用和展示。

通过以上的层次划分，数据仓库体系能够实现数据的整合、加工和应用。从原始数据层到公共维度层和明细数据层，再到汇总数据层和数据应用层，每个层次都有其特定的功能和用途，通过合理的分层结构，可以提高数据管理和分析的效率，同时满足不同用户的需求，简化复杂问题的处理，并为数据安全性和可维护性提供保障。

## 2.2 数据表定义

通过对数仓进行分层，我们设计每层的数据表定义如下。

### 2.2.1 ODS 层

#### 1. 活动信息表

ods\_activity\_info\_full

```
(  
    `id` STRING COMMENT '活动 id',  
    `activity_name` STRING COMMENT '活动名称',  
    `activity_type` STRING COMMENT '活动类型',  
    `activity_desc` STRING COMMENT '活动描述',  
    `start_time` STRING COMMENT '开始时间',  
    `end_time` STRING COMMENT '结束时间',  
    `create_time` STRING COMMENT '创建时间'
```



```
) COMMENT '活动信息表'
```

## 2. 活动规则表

```
ods_activity_rule_full
```

```
(  
  `id` STRING COMMENT '编号',  
  `activity_id` STRING COMMENT '类型',  
  `activity_type` STRING COMMENT '活动类型',  
  `condition_amount` DECIMAL(16, 2) COMMENT '满减金额',  
  `condition_num` BIGINT COMMENT '满减件数',  
  `benefit_amount` DECIMAL(16, 2) COMMENT '优惠金额',  
  `benefit_discount` DECIMAL(16, 2) COMMENT '优惠折扣',  
  `benefit_level` STRING COMMENT '优惠级别'  
) COMMENT '活动规则表'
```

## 3. 省份表

```
ods_base_province_full
```

```
(  
  `id` STRING COMMENT '编号',  
  `name` STRING COMMENT '省份名称',  
  `region_id` STRING COMMENT '地区 ID',  
  `area_code` STRING COMMENT '地区编码',  
  `iso_code` STRING COMMENT '旧版 ISO-3166-2 编码, 供可视化使用',  
  `iso_3166_2` STRING COMMENT '新版 IOS-3166-2 编码, 供可视化使用'  
) COMMENT '省份表'
```

## 4. 地区表

```
ods_base_region_full
```

```
(  
  `id` STRING COMMENT '编号',  
  `region_name` STRING COMMENT '地区名称'  
) COMMENT '地区表'
```

## 2.2.2 DIM 层

### 1. 商品维度表

```
dim_sku_full
```

```
(  
  `id` STRING COMMENT 'sku_id',  
  `price` DECIMAL(16, 2) COMMENT '商品价格',  
  `sku_name` STRING COMMENT '商品名称',  
  `sku_desc` STRING COMMENT '商品描述',  
  `weight` DECIMAL(16, 2) COMMENT '重量',  
  `is_sale` BOOLEAN COMMENT '是否在售',  
  `spu_id` STRING COMMENT 'spu 编号',  
  `spu_name` STRING COMMENT 'spu 名称',  
)
```

```

`category3_id`          STRING COMMENT '三级分类 id',
`category3_name`        STRING COMMENT '三级分类名称',
`category2_id`          STRING COMMENT '二级分类 id',
`category2_name`        STRING COMMENT '二级分类名称',
`category1_id`          STRING COMMENT '一级分类 id',
`category1_name`        STRING COMMENT '一级分类名称',
`tm_id`                 STRING COMMENT '品牌 id',
`tm_name`               STRING COMMENT '品牌名称',
`sku_attr_values`       ARRAY<STRUCT<attr_id :STRING,value_id :STRING,attr_name :STRING,value_name:STRING>> COMMENT '平台属性',
`sku_sale_attr_values`  ARRAY<STRUCT<sale_attr_id :STRING,sale_attr_value_id :STRING,sale_attr_name :STRING,sale_attr_value_name:STRING>> COMMENT '销售属性',
`create_time`           STRING COMMENT '创建时间'
) COMMENT '商品维度表'

```

## 2. 优惠券维度表

**dim\_coupon\_full**

```

(
`id`          STRING COMMENT '购物券编号',
`coupon_name` STRING COMMENT '购物券名称',
`coupon_type_code` STRING COMMENT '购物券类型编码',
`coupon_type_name` STRING COMMENT '购物券类型名称',
`condition_amount` DECIMAL(16, 2) COMMENT '满额数',
`condition_num` BIGINT COMMENT '满件数',
`activity_id` STRING COMMENT '活动编号',
`benefit_amount` DECIMAL(16, 2) COMMENT '减金额',
`benefit_discount` DECIMAL(16, 2) COMMENT '折扣',
`benefit_rule` STRING COMMENT '优惠规则:满元*减*元, 满*件打*折',
`create_time` STRING COMMENT '创建时间',
`range_type_code` STRING COMMENT '优惠范围类型编码',
`range_type_name` STRING COMMENT '优惠范围类型名称',
`limit_num` BIGINT COMMENT '最多领取次数',
`taken_count` BIGINT COMMENT '已领取次数',
`start_time` STRING COMMENT '可以领取的开始日期',
`end_time` STRING COMMENT '可以领取的结束日期',
`operate_time` STRING COMMENT '修改时间',
`expire_time` STRING COMMENT '过期时间'
) COMMENT '优惠券维度表'

```

## 3. 用户维度表

**dim\_user\_zip**

```

(
`id`          STRING COMMENT '用户 id',
`login_name`  STRING COMMENT '用户名称',
`nick_name`   STRING COMMENT '用户昵称',
`name`        STRING COMMENT '用户姓名',
`phone_num`   STRING COMMENT '手机号码',

```

```

`email`          STRING COMMENT '邮箱',
`user_level`     STRING COMMENT '用户等级',
`birthday`       STRING COMMENT '生日',
`gender`         STRING COMMENT '性别',
`create_time`    STRING COMMENT '创建时间',
`operate_time`   STRING COMMENT '操作时间',
`start_date`     STRING COMMENT '开始日期',
`end_date`       STRING COMMENT '结束日期'
) COMMENT '用户表'

```

## 2.2.3 DWD 层

### 1. 交易域加购事务事实表

```

dwd_trade_cart_add_inc
(
    `id`          STRING COMMENT '编号',
    `user_id`     STRING COMMENT '用户 id',
    `sku_id`      STRING COMMENT '商品 id',
    `date_id`     STRING COMMENT '时间 id',
    `create_time` STRING COMMENT '加购时间',
    `source_id`   STRING COMMENT '来源类型 ID',
    `source_type_code` STRING COMMENT '来源类型编码',
    `source_type_name` STRING COMMENT '来源类型名称',
    `sku_num`     BIGINT COMMENT '加购物车件数'
) COMMENT '交易域加购物车事务事实表'

```

### 2. 交易域下单事务事实表

```

dwd_trade_order_detail_inc
(
    `id`          STRING COMMENT '编号',
    `order_id`    STRING COMMENT '订单 id',
    `user_id`     STRING COMMENT '用户 id',
    `sku_id`      STRING COMMENT '商品 id',
    `province_id` STRING COMMENT '省份 id',
    `activity_id` STRING COMMENT '参与活动规则 id',
    `activity_rule_id` STRING COMMENT '参与活动规则 id',
    `coupon_id`   STRING COMMENT '使用优惠券 id',
    `date_id`     STRING COMMENT '下单日期 id',
    `create_time` STRING COMMENT '下单时间',
    `source_id`   STRING COMMENT '来源编号',
    `source_type_code` STRING COMMENT '来源类型编码',
    `source_type_name` STRING COMMENT '来源类型名称',
    `sku_num`     BIGINT COMMENT '商品数量',
    `split_original_amount` DECIMAL(16, 2) COMMENT '原始价格',
    `split_activity_amount` DECIMAL(16, 2) COMMENT '活动优惠分摊',
    `split_coupon_amount` DECIMAL(16, 2) COMMENT '优惠券优惠分摊',
    `split_total_amount` DECIMAL(16, 2) COMMENT '最终价格分摊'
) COMMENT '交易域下单明细事务事实表'

```

### 3. 交易域取消订单事务事实表

dwd\_trade\_cancel\_detail\_inc

```
(
    `id` STRING COMMENT '编号',
    `order_id` STRING COMMENT '订单 id',
    `user_id` STRING COMMENT '用户 id',
    `sku_id` STRING COMMENT '商品 id',
    `province_id` STRING COMMENT '省份 id',
    `activity_id` STRING COMMENT '参与活动规则 id',
    `activity_rule_id` STRING COMMENT '参与活动规则 id',
    `coupon_id` STRING COMMENT '使用优惠券 id',
    `date_id` STRING COMMENT '取消订单日期 id',
    `cancel_time` STRING COMMENT '取消订单时间',
    `source_id` STRING COMMENT '来源编号',
    `source_type_code` STRING COMMENT '来源类型编码',
    `source_type_name` STRING COMMENT '来源类型名称',
    `sku_num` BIGINT COMMENT '商品数量',
    `split_original_amount` DECIMAL(16, 2) COMMENT '原始价格',
    `split_activity_amount` DECIMAL(16, 2) COMMENT '活动优惠分摊',
    `split_coupon_amount` DECIMAL(16, 2) COMMENT '优惠券优惠分摊',
    `split_total_amount` DECIMAL(16, 2) COMMENT '最终价格分摊'
) COMMENT '交易域取消订单明细事务事实表'
```

## 2.2.4 DWS 层

### 1. 最近 1 日汇总表

dws\_trade\_user\_sku\_order\_1d

```
(
    `user_id` STRING COMMENT '用户 id',
    `sku_id` STRING COMMENT 'sku_id',
    `sku_name` STRING COMMENT 'sku 名称',
    `category1_id` STRING COMMENT '一级分类 id',
    `category1_name` STRING COMMENT '一级分类名称',
    `category2_id` STRING COMMENT '一级分类 id',
    `category2_name` STRING COMMENT '一级分类名称',
    `category3_id` STRING COMMENT '一级分类 id',
    `category3_name` STRING COMMENT '一级分类名称',
    `tm_id` STRING COMMENT '品牌 id',
    `tm_name` STRING COMMENT '品牌名称',
    `order_count_1d` BIGINT COMMENT '最近 1 日下单次数',
    `order_num_1d` BIGINT COMMENT '最近 1 日下单件数',
    `order_original_amount_1d` DECIMAL(16, 2) COMMENT '最近 1 日下单原始金额',
    `activity_reduce_amount_1d` DECIMAL(16, 2) COMMENT '最近 1 日活动优惠金额',
    `coupon_reduce_amount_1d` DECIMAL(16, 2) COMMENT '最近 1 日优惠券优惠金额',
    `order_total_amount_1d` DECIMAL(16, 2) COMMENT '最近 1 日下单最终金额'
)
```

```
'
) COMMENT '交易域用户商品粒度订单最近 1 日汇总事实表'
```

dws\_trade\_user\_sku\_order\_refund\_1d

```
(
  `user_id`          STRING COMMENT '用户 id',
  `sku_id`           STRING COMMENT 'sku_id',
  `sku_name`         STRING COMMENT 'sku 名称',
  `category1_id`     STRING COMMENT '一级分类 id',
  `category1_name`   STRING COMMENT '一级分类名称',
  `category2_id`     STRING COMMENT '一级分类 id',
  `category2_name`   STRING COMMENT '一级分类名称',
  `category3_id`     STRING COMMENT '一级分类 id',
  `category3_name`   STRING COMMENT '一级分类名称',
  `tm_id`            STRING COMMENT '品牌 id',
  `tm_name`          STRING COMMENT '品牌名称',
  `order_refund_count_1d` BIGINT COMMENT '最近 1 日退单次数',
  `order_refund_num_1d`  BIGINT COMMENT '最近 1 日退单件数',
  `order_refund_amount_1d` DECIMAL(16, 2) COMMENT '最近 1 日退单金额'
) COMMENT '交易域用户商品粒度退单最近 1 日汇总事实表'
```

## 2. 最近 n 日汇总表

dws\_trade\_user\_sku\_order\_nd

```
(
  `user_id`          STRING COMMENT '用户 id',
  `sku_id`           STRING COMMENT 'sku_id',
  `sku_name`         STRING COMMENT 'sku 名称',
  `category1_id`     STRING COMMENT '一级分类 id',
  `category1_name`   STRING COMMENT '一级分类名称',
  `category2_id`     STRING COMMENT '一级分类 id',
  `category2_name`   STRING COMMENT '一级分类名称',
  `category3_id`     STRING COMMENT '一级分类 id',
  `category3_name`   STRING COMMENT '一级分类名称',
  `tm_id`            STRING COMMENT '品牌 id',
  `tm_name`          STRING COMMENT '品牌名称',
  `order_count_7d`   STRING COMMENT '最近 7 日下单次数',
  `order_num_7d`     BIGINT COMMENT '最近 7 日下单件数',
  `order_original_amount_7d` DECIMAL(16, 2) COMMENT '最近 7 日下单原始金
额',
  `activity_reduce_amount_7d` DECIMAL(16, 2) COMMENT '最近 7 日活动优惠金
额',
  `coupon_reduce_amount_7d`  DECIMAL(16, 2) COMMENT '最近 7 日优惠券优惠
金额',
  `order_total_amount_7d`    DECIMAL(16, 2) COMMENT '最近 7 日下单最终金
额',
  `order_count_30d`         BIGINT COMMENT '最近 30 日下单次数',
  `order_num_30d`          BIGINT COMMENT '最近 30 日下单件数',
  `order_original_amount_30d` DECIMAL(16, 2) COMMENT '最近 30 日下单原始
金额',
```

```

`activity_reduce_amount_30d` DECIMAL(16, 2) COMMENT '最近 30 日活动优惠金额',
`coupon_reduce_amount_30d` DECIMAL(16, 2) COMMENT '最近 30 日优惠券优惠金额',
`order_total_amount_30d` DECIMAL(16, 2) COMMENT '最近 30 日下单最终金额'
) COMMENT '交易域用户商品粒度订单最近 n 日汇总事实表'

```

dws\_trade\_user\_sku\_order\_refund\_nd

```

(
  `user_id` STRING COMMENT '用户 id',
  `sku_id` STRING COMMENT 'sku_id',
  `sku_name` STRING COMMENT 'sku 名称',
  `category1_id` STRING COMMENT '一级分类 id',
  `category1_name` STRING COMMENT '一级分类名称',
  `category2_id` STRING COMMENT '一级分类 id',
  `category2_name` STRING COMMENT '一级分类名称',
  `category3_id` STRING COMMENT '一级分类 id',
  `category3_name` STRING COMMENT '一级分类名称',
  `tm_id` STRING COMMENT '品牌 id',
  `tm_name` STRING COMMENT '品牌名称',
  `order_refund_count_7d` BIGINT COMMENT '最近 7 日退单次数',
  `order_refund_num_7d` BIGINT COMMENT '最近 7 日退单件数',
  `order_refund_amount_7d` DECIMAL(16, 2) COMMENT '最近 7 日退单金额',
  `order_refund_count_30d` BIGINT COMMENT '最近 30 日退单次数',
  `order_refund_num_30d` BIGINT COMMENT '最近 30 日退单件数',
  `order_refund_amount_30d` DECIMAL(16, 2) COMMENT '最近 30 日退单金额'
) COMMENT '交易域用户商品粒度退单最近 n 日汇总事实表'

```

### 3. 历史至今汇总表

dws\_trade\_user\_order\_td

```

(
  `user_id` STRING COMMENT '用户 id',
  `order_date_first` STRING COMMENT '首次下单日期',
  `order_date_last` STRING COMMENT '末次下单日期',
  `order_count_td` BIGINT COMMENT '下单次数',
  `order_num_td` BIGINT COMMENT '购买商品件数',
  `original_amount_td` DECIMAL(16, 2) COMMENT '原始金额',
  `activity_reduce_amount_td` DECIMAL(16, 2) COMMENT '活动优惠金额',
  `coupon_reduce_amount_td` DECIMAL(16, 2) COMMENT '优惠券优惠金额',
  `total_amount_td` DECIMAL(16, 2) COMMENT '最终金额'
) COMMENT '交易域用户粒度订单历史至今汇总事实表'

```

dws\_trade\_user\_payment\_td

```

(
  `user_id` STRING COMMENT '用户 id',
  `payment_date_first` STRING COMMENT '首次支付日期',

```

```

    `payment_date_last` STRING COMMENT '末次支付日期',
    `payment_count_td` BIGINT COMMENT '最近 7 日支付次数',
    `payment_num_td` BIGINT COMMENT '最近 7 日支付商品件数',
    `payment_amount_td` DECIMAL(16, 2) COMMENT '最近 7 日支付金额'
) COMMENT '交易域用户粒度支付历史至今汇总事实表'

```

## 2.2.5 ADS 层

### 1. 流量主题

```

ads_traffic_stats_by_channel
(
    `dt` STRING COMMENT '统计日期',
    `recent_days` BIGINT COMMENT '最近天数,1:最近 1 天,7:最近 7 天,30:最近 30 天',
    `channel` STRING COMMENT '渠道',
    `uv_count` BIGINT COMMENT '访客人数',
    `avg_duration_sec` BIGINT COMMENT '会话平均停留时长,单位为秒',
    `avg_page_count` BIGINT COMMENT '会话平均浏览页面数',
    `sv_count` BIGINT COMMENT '会话数',
    `bounce_rate` DECIMAL(16, 2) COMMENT '跳出率'
) COMMENT '各渠道流量统计'

```

```

ads_page_path
(
    `dt` STRING COMMENT '统计日期',
    `recent_days` BIGINT COMMENT '最近天数,1:最近 1 天,7:最近 7 天,30:最近 30 天',
    `source` STRING COMMENT '跳转起始页面 ID',
    `target` STRING COMMENT '跳转终到页面 ID',
    `path_count` BIGINT COMMENT '跳转次数'
) COMMENT '页面浏览路径分析'

```

### 2. 用户主题

```

ads_user_change
(
    `dt` STRING COMMENT '统计日期',
    `user_churn_count` BIGINT COMMENT '流失用户数',
    `user_back_count` BIGINT COMMENT '回流用户数'
) COMMENT '用户变动统计'

```

```

ads_user_retention
(
    `dt` STRING COMMENT '统计日期',
    `create_date` STRING COMMENT '用户新增日期',
    `retention_day` INT COMMENT '截至当前日期留存天数',
    `retention_count` BIGINT COMMENT '留存用户数量',
    `new_user_count` BIGINT COMMENT '新增用户数量',
    `retention_rate` DECIMAL(16, 2) COMMENT '留存率'
)

```

```
) COMMENT '用户留存率'
```

**ads\_user\_stats**

```
(  
    `dt` STRING COMMENT '统计日期',  
    `recent_days` BIGINT COMMENT '最近 n 日,1:最近 1 日,7:最近 7 日,30:最近 30 日',  
    `new_user_count` BIGINT COMMENT '新增用户数',  
    `active_user_count` BIGINT COMMENT '活跃用户数'  
) COMMENT '用户新增活跃统计'
```

**ads\_user\_action**

```
(  
    `dt` STRING COMMENT '统计日期',  
    `recent_days` BIGINT COMMENT '最近天数,1:最近 1 天,7:最近 7 天,30:最近 30 天',  
    `home_count` BIGINT COMMENT '浏览首页人数',  
    `good_detail_count` BIGINT COMMENT '浏览商品详情页人数',  
    `cart_count` BIGINT COMMENT '加入购物车人数',  
    `order_count` BIGINT COMMENT '下单人数',  
    `payment_count` BIGINT COMMENT '支付人数'  
) COMMENT '漏斗分析'
```

**ads\_new\_buyer\_stats**

```
(  
    `dt` STRING COMMENT '统计日期',  
    `recent_days` BIGINT COMMENT '最近天数,1:最近 1 天,7:最近 7 天,30:最近 30 天',  
    `new_order_user_count` BIGINT COMMENT '新增下单人数',  
    `new_payment_user_count` BIGINT COMMENT '新增支付人数'  
) COMMENT '新增交易用户统计'
```

### 3. 商品主题

**ads\_repeat\_purchase\_by\_tm**

```
(  
    `dt` STRING COMMENT '统计日期',  
    `recent_days` BIGINT COMMENT '最近天数,7:最近 7 天,30:最近 30 天',  
    `tm_id` STRING COMMENT '品牌 ID',  
    `tm_name` STRING COMMENT '品牌名称',  
    `order_repeat_rate` DECIMAL(16, 2) COMMENT '复购率'  
) COMMENT '各品牌复购率统计'
```

**ads\_trade\_stats\_by\_tm**

```
(  
    `dt` STRING COMMENT '统计日期',  
    `recent_days` BIGINT COMMENT '最近天数,1:最近 1 天,7:最近 7 天,30:最近 30 天',
```



```

    `tm_id`                STRING COMMENT '品牌 ID',
    `tm_name`              STRING COMMENT '品牌名称',
    `order_count`          BIGINT COMMENT '订单数',
    `order_user_count`     BIGINT COMMENT '订单人数',
    `order_refund_count`   BIGINT COMMENT '退单数',
    `order_refund_user_count` BIGINT COMMENT '退单人数'
) COMMENT '各品牌商品交易统计'

```

#### ads\_trade\_stats\_by\_cate

```

(
    `dt`                STRING COMMENT '统计日期',
    `recent_days`       BIGINT COMMENT '最近天数,1:最近 1 天,7:最近 7
天,30:最近 30 天',
    `category1_id`      STRING COMMENT '一级分类 id',
    `category1_name`    STRING COMMENT '一级分类名称',
    `category2_id`      STRING COMMENT '二级分类 id',
    `category2_name`    STRING COMMENT '二级分类名称',
    `category3_id`      STRING COMMENT '三级分类 id',
    `category3_name`    STRING COMMENT '三级分类名称',
    `order_count`       BIGINT COMMENT '订单数',
    `order_user_count`  BIGINT COMMENT '订单人数',
    `order_refund_count` BIGINT COMMENT '退单数',
    `order_refund_user_count` BIGINT COMMENT '退单人数'
) COMMENT '各分类商品交易统计'

```

#### ads\_sku\_cart\_num\_top3\_by\_cate

```

(
    `dt`                STRING COMMENT '统计日期',
    `category1_id`      STRING COMMENT '一级分类 ID',
    `category1_name`    STRING COMMENT '一级分类名称',
    `category2_id`      STRING COMMENT '二级分类 ID',
    `category2_name`    STRING COMMENT '二级分类名称',
    `category3_id`      STRING COMMENT '三级分类 ID',
    `category3_name`    STRING COMMENT '三级分类名称',
    `sku_id`            STRING COMMENT '商品 id',
    `sku_name`          STRING COMMENT '商品名称',
    `cart_num`          BIGINT COMMENT '购物车中商品数量',
    `rk`                BIGINT COMMENT '排名'
) COMMENT '各分类商品购物车存量 Top10'

```

## 4. 交易主题

#### ads\_trade\_stats

```

(
    `dt`                STRING COMMENT '统计日期',
    `recent_days`       BIGINT COMMENT '最近天数,1:最近 1 日,7:最近 7
天,30:最近 30 天',
    `order_total_amount` DECIMAL(16, 2) COMMENT '订单总额,GMV',
    `order_count`       BIGINT COMMENT '订单数',

```

```

    `order_user_count`          BIGINT COMMENT '下单人数',
    `order_refund_count`        BIGINT COMMENT '退单数',
    `order_refund_user_count`   BIGINT COMMENT '退单人数'
) COMMENT '交易统计'

```

#### ads\_order\_by\_province

```

(
    `dt`          STRING COMMENT '统计日期',
    `recent_days` BIGINT COMMENT '最近天数,1:最近 1 天,7:最近 7 天,30:最近 30 天',
    `province_id` STRING COMMENT '省份 ID',
    `province_name` STRING COMMENT '省份名称',
    `area_code`    STRING COMMENT '地区编码',
    `iso_code`     STRING COMMENT '国际标准地区编码',
    `iso_code_3166_2` STRING COMMENT '国际标准地区编码',
    `order_count`  BIGINT COMMENT '订单数',
    `order_total_amount` DECIMAL(16, 2) COMMENT '订单金额'
) COMMENT '各地区订单统计'

```

#### 5. 优惠券主题

##### ads\_coupon\_stats

```

(
    `dt`          STRING COMMENT '统计日期',
    `coupon_id`   STRING COMMENT '优惠券 ID',
    `coupon_name` STRING COMMENT '优惠券名称',
    `start_date`  STRING COMMENT '发布日期',
    `rule_name`   STRING COMMENT '优惠规则,例如满 100 元减 10 元',
    `reduce_rate` DECIMAL(16, 2) COMMENT '补贴率'
) COMMENT '优惠券统计'

```

#### 6. 活动主题

##### ads\_activity\_stats

```

(
    `dt`          STRING COMMENT '统计日期',
    `activity_id`  STRING COMMENT '活动 ID',
    `activity_name` STRING COMMENT '活动名称',
    `start_date`   STRING COMMENT '活动开始日期',
    `reduce_rate`  DECIMAL(16, 2) COMMENT '补贴率'
) COMMENT '活动统计'

```

## 2.3 数仓数据同步策略

业务数据是数据仓库的重要数据来源，我们需要每日定时从业务数据库中抽取数据，传输到数据仓库中，之后再对数据进行分析统计。

为保证统计结果的正确性，需要保证数据仓库中的数据与业务数据库是同步的，离线数仓的计

算周期通常为天，所以数据同步周期也通常为天，即每天同步一次即可。

数据的同步策略有全量同步和增量同步。

#### 全量同步：

每天固定时间点，触发全量同步任务。

在同步任务开始前，首先对数据仓库中的对应表进行清空或标记为过期状态，以确保新的全量数据可以正确加载。

从业务数据库中抽取全部数据，并按照预定义的数据转换规则进行清洗、转换和规范化。

将转换后的数据加载到数据仓库的对应表中，以替换旧的数据。

完成数据加载后，进行必要的索引重建和数据校验，确保数据仓库中的数据与业务数据库一致性。

在同步过程中，记录同步日志和错误信息，以便后续排查和监控。

#### 增量同步：

首先进行一次全量同步，将业务数据库中的全部数据加载到数据仓库。

每天固定时间点，触发增量同步任务。

通过比较业务数据库中的数据更新时间戳或增量标识，筛选出当天发生变化的数据。

从业务数据库中抽取新增和变化数据，并按照转换规则进行清洗、转换和规范化。

将转换后的增量数据加载到数据仓库的对应表中，以追加方式更新数据。

在同步过程中，记录同步日志和错误信息，以便后续排查和监控。

为确保数据同步的可靠性和稳定性，我们还可以考虑了以下细节：

异常处理：处理抽取、转换和加载过程中的异常情况，例如网络中断、数据格式错误等。可以设置预警机制和错误处理流程。

并发处理：处理并发情况下的数据冲突，例如多个任务同时修改同一条记录。可以采用锁机制或乐观锁定策略来处理并发更新。

数据校验：在数据加载后，进行数据校验和一致性检查，确保数据的准确性和完整性。

监控和报告：建立监控系统，对数据同步任务进行实时监控，并生成报告，记录同步状态、延迟和错误信息，以便及时发现和解决问题。

通过合理设计和实施数据同步策略，可以确保数据仓库中的数据与业务数据库保持同步，为后续的数据分析和统计提供准确可靠的基础。

## 3. 人工智能模块

### 3.1 用户画像分析

#### 用户画像分析

基于数据仓库中的数据	用户省份 id 用户性别 用户行为 7、30 天内买得最多的品类 7、30 天内买得第二多的品类 7、30 天内买得第三多的品类 7、30 天内买得最多的品牌 7、30 天内买得第二多的品牌 7、30 天内买得第三多的品牌 7、30 天内的购买单数 7、30 天内的消费总金额
使用的算法	KNN 聚类算法
完成的需求分析	<p>将用户进行分类是一种常见的市场细分策略，有助于更好地理解用户群体的需求和行为，从而提供更精准的服务和产品。以下是一种基于用户行为和消费习惯的分类方法：</p> <p>新用户（New Users）：这个群体包括那些刚注册或者刚开始使用服务的用户。他们可能对平台的功能和操作还不熟悉，所以需要更多的引导和教程。也可能需要更多的优惠和激励措施来促使他们进行第一次购买。</p> <p>活跃用户（Active Users）：这个群体的用户在平台上有频繁的活动，比如经常浏览、搜索和购买商品。他们可能是平台的忠实用户，对平台的运作和商品非常熟悉。</p> <p>沉默用户（Silent Users）：这类用户虽然已经注册，但是很少在平台上有活动。他们可能只是偶尔浏览，或者在需要时才会使用平台。</p> <p>流失用户（Churned Users）：这类用户在一段时间内完全没有活动，可能是因为他们找到了其他的服务，或者对当前的服务不满意。</p>

## 3.2 用户性别预测

用户性别预测	
基于数据仓库中的数据	30 天内买得最多的品类 30 天内买得第二多的品类 30 天内买得第三多的品类 30 天内买得最多的品牌 30 天内买得第二多的品牌 30 天内买得第三多的品牌 30 天内的购买单数 30 天内的消费总金额
使用的算法	朴素贝叶斯
完成的需求分析	用户行为性别预测分析

## 3.3 用户购买意向分析

用户购买意向分析	
基于数据仓库中的数据	用户每次浏览的平均时长 用户浏览过的不同商品种类的数量，反映用户的兴趣广泛程度 用户对平台热门商品的浏览次数占总浏览次数的比例 用户每次搜索的平均关键词数量 购物车中商品结算数量与加入数量的比例

	最近 30 日下单次数 最近 30 日下单商品件数 最近 30 日下单原始金额 最近 30 日活动优惠金额 最近 30 日优惠券优惠金额 最近 30 日加购次数 最近 30 日加购商品件数 最近 30 日支付次数 最近 30 日支付商品件数
使用的算法	决策树、随机森林
完成的需求分析	用户购买意向分析

3.4 售后服务满意度分析

售后服务满意度分析	
基于数据仓库中的数据	用户对售后服务的平均评分 用户对售后服务提出的投诉数量 根据用户的评价和反馈中的关键词 退单件数 退单金额 评价时间 评价名称
使用的算法	文本情感分类、情感强度分析
完成的需求分析	售后服务满意度分析

3.5 产品销售趋势分析

产品销售趋势分析	
基于数据仓库中的数据	每个月的产品销售额 基于历史数据计算得出的季节性指数 产品类别在总销售额中的占比 最近 7、30 日复购率 最近 1、7、30 日订单数 最近 1、7、30 日订单人数 最近 1、7、30 日退单数 最近 1、7、30 日退单人数 最近 1、7、30 日订单总额
使用的算法	时间序列分析算法（ARIMA、指数平滑）
完成的需求分析	产品销售趋势分析

3.6 商品热度分析

商品热度分析
--------

基于数据仓库中的数据	销售增长率 商品页面访问量 最近 7、30 日复购率 最近 1、7、30 日订单数 最近 1、7、30 日订单人数 最近 1、7、30 日收藏数 最近 1、7、30 日收藏人数 最近 1、7、30 日退单数 最近 1、7、30 日退单人数 最近 1、7、30 日订单总额
使用的算法	KNN 聚类算法
完成的需求分析	商品热度分析，分析最近一段时间内的热卖产品 热销商品（Hot-selling Products）：这类商品的销售量非常高，或者销售增长速度非常快。它们可能是一些流行的或者具有独特优势的商品。 潜力商品（Potential Products）：这类商品的销售量可能不高，但是浏览次数多，用户评价好，或者销售增长速度快。这可能意味着这些商品具有很大的增长潜力。 稳定商品（Stable Products）：这类商品的销售量和用户评价都稳定在一个较高的水平。它们可能是一些经典的或者必需的商品。 长尾商品（Long-tail Products）：这类商品的销售量可能不高，但是它们能够满足一些特定的用户需求。长尾商品的存在，可以提高商品的多样性，吸引更多的用户。

### 3.7 优惠券效果分析

优惠券效果分析	
基于数据仓库中的数据	优惠券信息 优惠券使用记录 用户消费历史 用户行为记录 订单价格 优惠价格 优惠后订单价格 优惠券使用数量
使用的算法	回归分析
完成的需求分析	优惠券效果分析，分析优惠券是否达到了发放的目的

### 3.8 用户流失预测

用户流失预测	
基于数据仓库中的数据	用户行为数据 用户个人信息 用户反馈数据 服务使用数据
使用的算法	支持向量机、随机森林
完成的需求分析	用户流失预测

## 4.软件开发模块

从用户跟踪、流量分析、商品分析和活动优惠这四个主题方面入手，采取以下思路分别设计一个可视化的页面：

### 4.1 用户跟踪

用户增长：使用折线图或区域图展示每天、每周或每月的用户增长趋势，帮助您了解用户的  
增长情况。用户分布：使用地图或热力图显示用户的地理分布，帮助您了解用户所在地区的  
分布情况。用户行为：使用漏斗图或路径图展示用户在网站或应用中的行为流程，帮助您分析  
用户在不同阶段的转化率和用户流失点。

### 4.2 流量分析

流量来源：使用饼图或柱状图显示不同来源（如搜索引擎、社交媒体、直接访问等）带来  
的流量比例，帮助您了解流量的来源分布。流量趋势：使用折线图展示每天、每周或每月的总  
流量趋势，帮助您了解流量的变化和高峰时段。受访页面：使用词云或热图展示最受访页面的  
关键词或点击热度，帮助您了解用户感兴趣的内容和热门页面。

### 4.3 商品分析

销售趋势：使用折线图或柱状图展示每天、每周或每月的销售额或销售量趋势，帮助您了  
解销售的季节性和趋势。产品对比：使用雷达图或柱状图展示不同产品的销售比较，帮助您了  
解产品的表现和市场占有率。价格分布：使用箱线图或直方图显示产品价格的分布情况，帮助  
您了解产品价格的范围和分布。

### 4.4 活动优惠

优惠效果：使用漏斗图或堆叠柱状图展示不同活动的参与人数和转化率，帮助您了解不同  
活动的效果和影响力。优惠渠道：使用漏斗图或饼图展示不同渠道（如电子邮件、社交媒体、  
广告等）带来的优惠参与量，帮助您了解不同渠道的效果和成本效益。优惠类型：使用堆叠柱  
状图或雷达图展示不同类型的优惠活动的参与人数和转化率，帮助您了解哪些优惠类型最受用  
户欢迎。

## 四、 开发过程

### 1. 云计算模块

#### 1.1. 配置 Hadoop 集群

##### 1.1.1. 配置 Hadoop 集群

- a. 编辑 Hadoop 配置文件。主要的配置文件是 core-site.xml, hdfs-site.xml, mapred-site.xml 和 yarn-site.xml。这些文件位于 Hadoop 的 etc/hadoop/目录下。
- b. 配置 core-site.xml, 指定 Hadoop 集群的名称和 HDFS 的默认文件系统 URI。
- c. 配置 hdfs-site.xml, 指定 Hadoop 分布式文件系统的副本数量和数据存储路径。
- d. 配置 mapred-site.xml, 指定 MapReduce 的框架和资源管理器。
- e. 配置 yarn-site.xml, 指定 YARN 资源管理器的配置参数。
- f. 配置其他相关的参数, 如日志路径、权限等。

##### 1.1.2. 配置集群的主节点和从节点

- a. 在其中一个服务器上配置 Hadoop 集群的主节点 (NameNode)。
- b. 在其他服务器上配置 Hadoop 集群的从节点 (DataNode)。
- c. 将从节点的主机名或 IP 地址添加到主节点的 etc/hadoop/slaves 文件中, 每行一个从节点。

##### 1.1.3. 分发 Hadoop 文件和配置

- g. 使用 SSH 将 Hadoop 文件和配置从主节点复制到所有从节点。
- h. 确保每个节点上的 Hadoop 文件和配置相同。

##### 1.1.4. 启动 Hadoop 集群

- a. 在主节点上执行启动 Hadoop 集群的命令: start-dfs.sh 和 start-yarn.sh。
- b. 检查集群的运行状态, 确保各个组件都正常启动。

#### 1.2. 配置 ZooKeeper 集群

##### 1.2.1. 配置 ZooKeeper 集群

- a. 创建一个用于存储 ZooKeeper 数据和日志的目录, 例如/var/lib/zookeeper.
- b. 在每个服务器上创建一个独立的配置文件。可以从 ZooKeeper 的 conf/目录中复制 zoo\_sample.cfg 为 zoo.cfg。



c. 编辑每个服务器上的 zoo.cfg 文件，指定以下参数：

1. **tickTime**: ZooKeeper 的基本时间单位（以毫秒为单位）。
2. **dataDir**: 数据目录的路径，用于存储 ZooKeeper 的事务日志和快照。
3. **clientPort**: 客户端连接 ZooKeeper 的端口。
4. **initLimit** 和 **syncLimit**: 用于配置 ZooKeeper 集群中的 Leader 选举和数据同步的参数。
5. **server.X**: 定义每个服务器在集群中的角色和通信地址，其中 **X** 是服务器的唯一标识符。

**zoo.cfg:**

```
tickTime=2000
dataDir=/var/lib/zookeeper
clientPort=2181
initLimit=10
syncLimit=5
server.1=server1.example.com:2888:3888
server.2=server2.example.com:2888:3888
server.3=server3.example.com:2888:3888
```

d. 在每个服务器上创建一个名为 myid 的文件，其中包含服务器的唯一标识符。该标识符与 zoo.cfg 中的 server.X 相对应。

### 1.2.2. 分发 ZooKeeper 文件和配置

- a. 使用 SSH 将 ZooKeeper 文件和配置从一个服务器复制到其他服务器。可以使用 scp 命令或其他文件传输工具。
- b. 确保每个服务器上的 ZooKeeper 文件和配置相同。

### 1.2.3. 启动 ZooKeeper 集群

在每个服务器上执行启动 ZooKeeper 的命令：zkServer.sh start，并检查集群的运行状态，确保各个节点都已启动。

## 1.3. 配置 Kafka 集群

### 1.3.1. 配置 Kafka 集群

- a. 创建一个目录用于存储 Kafka 的数据和日志，例如 /var/lib/kafka。
- b. 在每个服务器上创建一个独立的配置文件。可以从 Kafka 的 config/ 目录中复制 server.properties 为 serverX.properties，其中 X 是服务器的唯一标识符。
- c. 编辑每个服务器上的 serverX.properties 文件，指定以下参数：

1. **broker.id**: 每个服务器在集群中的唯一标识符。
2. **listeners**: Kafka 监听器的主机和端口，用于客户端和其他 Kafka 节点的通信。
3. **log.dirs**: 用于存储 Kafka 消息日志的目录路径。
4. **zookeeper.connect**: ZooKeeper 集群的连接字符串。

#### **serverX.properties:**

```
broker.id=1
```

```
listeners=PLAINTEXT://server1.example.com:9092
```

```
log.dirs=/var/lib/kafka
```

```
zookeeper.connect=server1.example.com:2181,server2.example.com:2181,server3.example.com:2181
```

### **1.3.2. 分发 Kafka 文件和配置**

使用 SSH 将 Kafka 文件和配置从一个服务器复制到其他服务器。可以使用 scp 命令或其他文件传输工具，确保每个服务器上的 Kafka 文件和配置相同。

### **1.3.3. 启动 Kafka 集群**

在每个服务器上执行启动 Kafka 的命令：bin/kafka-server-start.sh -daemon config/serverX.properties，并检查集群的运行状态，确保各个节点都已启动。

## **1.4. 配置 Flume 集群**

### **1.4.1. 下载和解压 Flume**

a. 在每台服务器上下载相同版本的 Apache Flume。可以从 Apache Flume 官方网站下载稳定版本。

b. 解压下载的 Flume 二进制文件到一个适当的目录，例如/opt/flume。

### **1.4.2. 创建 Flume 配置文件**

- a. 在一个服务器上创建一个 Flume Agent 的配置文件，例如 flume-agent1.conf。
- b. 在另一个服务器上创建另一个 Flume Agent 的配置文件，例如 flume-agent2.conf。
- c. 编辑每个配置文件，指定以下参数：
  - i. **agent.sources**: 指定 Agent 的数据源类型和配置。
  - ii. **agent.sinks**: 指定 Agent 的数据接收器类型和配置。
  - iii. **agent.channels**: 指定 Agent 的通道类型和配置。
  - iv. **agent.sources.<source\_name>.channels**: 指定数据源连接的通道。
  - v. **agent.sinks.<sink\_name>.channel**: 指定数据接收器连接的通道。
  - vi. **agent.channels.<channel\_name>.type**: 指定通道的类型。

#### **flume-agent1.conf:**

```
agent.sources = source1
```

```
agent.sinks = sink1
```

```
agent.channels = channel1
```

```
agent.sources.source1.type = <source_type>
```

```
agent.sources.source1.channels = channel1
```

```
agent.sources.source1.<source_specific_properties> = ...
```

```
agent.sinks.sink1.type = <sink_type>
```

```
agent.sinks.sink1.channel = channel1
```

```
agent.sinks.sink1.<sink_specific_properties> = ...
```

```
agent.channels.channel1.type = <channel_type>
```

```
agent.channels.channel1.<channel_specific_properties> = ...
```

### **1.4.3. 分发 Flume 文件和配置**

使用 SSH 将 Flume 文件和配置从一个服务器复制到其他服务器。可以使用 **scp** 命令或其他文件传输工具。确保每个服务器上的 Flume 文件和配置相同。

#### **1.4.4. 启动 Flume Agents**

在每个服务器上分别启动 Flume Agent，指定不同的配置文件：bin/flume-ng agent --conf conf --conf-file /path/to/flume-agent1.conf --name agent1 -Dflume.root.logger=INFO,console 和 bin/flume-ng agent --conf conf --conf-file /path/to/flume-agent2.conf --name agent2 -Dflume.root.logger=INFO,console。检查各个 Agent 的运行状态，确保数据能够流动。

## **1.5. 编写 Flume 拦截器**

### **1.5.1. 创建拦截器项目**

- a. 创建一个新的 Java 项目，或在现有 Java 项目中添加 Flume 拦截器模块。
- b. 导入 Flume 的依赖库，包括 flume-ng-core 和其他必要的库。

### **1.5.2. 创建拦截器类**

- a. 创建一个类，实现 Flume 的 Interceptor 接口。
- b. 实现 initialize 方法，用于初始化拦截器。

- c. 实现 intercept 方法，用于对事件进行处理和转换。
- d. 实现 close 方法，用于拦截器的清理工作。

#### **UserLogInterceptor.java:**

```
import org.apache.flume.Context;
import org.apache.flume.Event;
import org.apache.flume.interceptor.Interceptor;

import java.nio.charset.StandardCharsets;
import java.util.List;
import java.util.Map;

public class UserLogInterceptor implements Interceptor {

    @Override
    public void initialize() {
        // 初始化拦截器
    }

    @Override
    public Event intercept(Event event) {
        // 对事件进行处理和转换
        byte[] body = event.getBody();
        String logData = new String(body, StandardCharsets.UTF_8);

        // 解析 JSON 日志数据
        try {
            Map<String, Object> parsedData = JSONParser.parse(logData);

            // 处理特定字段
            String userId = (String) parsedData.get("userId");
            String action = (String) parsedData.get("action");

            // 添加额外的字段
            parsedData.put("processedField", "some value");

            // 更新日志数据为处理后的数据
            String processedLogData = JSONParser.toJson(parsedData);
```

```
event.setBody(processedLogData.getBytes(StandardCharsets.UTF_8));
    } catch (Exception e) {
        // 处理日志解析异常
        e.printStackTrace();
    }

    return event;
}

@Override
public List<Event> intercept(List<Event> events) {
    for (Event event : events) {
        intercept(event);
    }
    return events;
}

@Override
public void close() {
    // 清理拦截器
}

public static class Builder implements Interceptor.Builder {

    @Override
    public void configure(Context context) {
        // 配置拦截器
    }

    @Override
    public Interceptor build() {
        return new UserLogInterceptor();
    }
}
}
```

### 1.5.3. 配置 Flume Agent

在 Flume 的配置文件中，定义使用拦截器的 Agent。在 Agent 的配置中，指定拦截器的类型和属性。

#### flume\_agent.conf

```
agent.sources = source1
```

```
agent.sinks = sink1
```

```
agent.channels = channel1
```

```
agent.sources.source1.type = <source_type>
```

```
agent.sources.source1.interceptors = interceptor1
```

```
agent.sources.source1.interceptors.interceptor1.type =
```

```
com.example.UpperCaseInterceptor$Builder
```

```
agent.sinks.sink1.type = <sink_type>
```

```
agent.sinks.sink1.channel = channel1
```

```
agent.channels.channel1.type = <channel_type>
```

### 1.5.4. 编译和部署

- 将拦截器类编译为一个可执行的 JAR 文件。
- 将拦截器的 JAR 文件复制到 Flume 的 lib 目录中。
- 启动 Flume Agent，让拦截器生效。

完成以上步骤后，Flume 会在数据传输过程中应用拦截器，并根据拦截器的逻辑对事件进行处理和转换。

## 1.6. 生成模拟数据

### 1.6.1. 启动 Kafka 集群

- 确保你已经搭建和配置了一个 Kafka 集群，包括至少一个 Broker 和一个 Topic。
- 启动 Kafka 集群，确保各个组件都正常运行。

### 1.6.2. 编写数据生产者

- 创建一个简单的数据生产者程序，使用 Kafka 的 Java 客户端库。

- b. 在程序中指定 Kafka 集群的地址和端口，以及要发送数据的 Topic 名称。
- c. 生成模拟数据并将其发送到指定的 Topic 中。

### DataProducer.java

```
import org.apache.kafka.clients.producer.*;

import java.util.Properties;

public class DataProducer {

    private static final String BOOTSTRAP_SERVERS =
"kafka1.example.com:9092,kafka2.example.com:9092,kafka3.example.com:9092";

    private static final String TOPIC_NAME = "user_logs";

    public static void main(String[] args) {

        // 配置 Kafka 生产者
        Properties props = new Properties();
        props.put("bootstrap.servers", BOOTSTRAP_SERVERS);
        props.put("key.serializer",
"org.apache.kafka.common.serialization.StringSerializer");
        props.put("value.serializer",
"org.apache.kafka.common.serialization.StringSerializer");

        // 创建 Kafka 生产者
        Producer<String, String> producer = new KafkaProducer<>(props);

        try {

            // 模拟生成并发送数据
            for (int i = 1; i <= 10; i++) {
                String userId = "user" + i;
                String action = "action" + i;
                String logData = "{\\"userId\\":\\" + userId + "\\",
\\"action\\":\\" + action + "\\}";

                // 创建 Kafka 记录
                ProducerRecord<String, String> record = new
ProducerRecord<>(TOPIC_NAME, logData);
```

```

        // 发送记录
        producer.send(record, new Callback() {
            @Override
            public void onCompletion(RecordMetadata metadata,
Exception exception) {
                if (exception != null) {
                    exception.printStackTrace();
                } else {
                    System.out.println("Sent record: " +
logData);
                }
            }
        });
    } finally {
        producer.close();
    }
}
}

```

在上述示例中，Kafka 生产者会生成 10 条模拟数据，并将其发送到名为 **user\_logs** 的 Topic 中。

### 1.6.3. 启动 Kafka 消费者

- 创建一个简单的 Kafka 消费者程序，使用 Kafka 的 Java 客户端库。
- 在程序中指定 Kafka 集群的地址和端口，以及要消费的 Topic 名称。
- 订阅指定的 Topic，并处理接收到的数据。

以下是一个简单的 Java 代码示例，用于启动 Kafka 消费者并消费来自 **user\_logs** Topic 的数据：

#### DataConsumer.java

```

import org.apache.kafka.clients.consumer.*;
import org.apache.kafka.common.TopicPartition;
import java.util.*;

public class DataConsumer {

```



```

    private static final String BOOTSTRAP_SERVERS =
"kafka1.example.com:9092,kafka2.example.com:9092,kafka3.example.com:9092";

    private static final String TOPIC_NAME = "user_logs";


    public static void main(String[] args) {
        // 配置 Kafka 消费者
        Properties props = new Properties();
        props.put("bootstrap.servers", BOOTSTRAP_SERVERS);
        props.put("group.id", "user_logs_consumer");
        props.put("key.deserializer",
"org.apache.kafka.common.serialization.StringDeserializer");
        props.put("value.deserializer",
"org.apache.kafka.common.serialization.StringDeserializer");


        // 创建 Kafka 消费者
        KafkaConsumer<String, String> consumer = new
KafkaConsumer<>(props);


        // 订阅 Topic
        consumer.subscribe(Collections.singletonList(TOPIC_NAME));


        // 消费数据
        try {
            while (true) {
                ConsumerRecords<String, String> records =
consumer.poll(Duration.ofMillis(100));
                for (ConsumerRecord<String, String> record : records) {
                    String logData = record.value();
                    System.out.println("Received record: " + logData);
                }
            }
        } finally {
            consumer.close();
        }
    }
}

```

在上述示例中，Kafka 消费者会订阅名为 **user\_logs** 的 Topic，并持续消费从该 Topic 中接收到的数据。

#### 1.6.4. 运行程序并观察数据消费情况

- a. 编译并运行数据生产者程序，确保模拟数据成功发送到 Kafka 集群。
- b. 编译并运行 Kafka 消费者程序，观察控制台输出，确认消费者是否能够接收和消费到发送的数据。

观察到的数据类似于以下内容：

```
Received record: {"userId":"user1", "action":"action1"}
```

```
Received record: {"userId":"user2", "action":"action2"}
```

```
...
```

```
Received record: {"userId":"user10", "action":"action10"}
```

### 1.7. 搭建部署 Maxwell

#### 1.7.1. 下载和安装 Maxwell

- a. 在服务器上下载 Maxwell 的二进制文件。
- b. 解压下载的 Maxwell 二进制文件到一个适当的目录，例如/opt/maxwell.
- c. 在 Maxwell 的安装目录下创建一个配置文件，例如 config.properties。
- d. 编辑配置文件，指定以下参数：
  - i. **host**: MySQL 数据库的主机名或 IP 地址。
  - ii. **port**: MySQL 数据库的端口号。
  - iii. **user**: 连接 MySQL 数据库的用户名。
  - iv. **password**: 连接 MySQL 数据库的密码。
  - v. **database**: 要监听变更事件的数据库名称。
  - vi. **producer**: 指定消息队列的配置，如 Kafka 的连接信息。
  - vii. **include\_table**: 指定要捕获变更事件的表名。

**config.properties:**

```
host=mysql.example.com
```

```
port=3306
```

```
user=maxwell
```

```
password=maxwell_password
```

```
database=my_database
```

```
producer=kafka
```

```
kafka.bootstrap.servers=kafka1.example.com:9092,kafka2.example.com:9092,kafka3.example.com:9092  
include_table=my_table1,my_table2
```

### 1.7.2. 启动 Maxwell

- 在服务器上执行启动 Maxwell 的命令：`./bin/maxwell --config /opt/maxwell/config.properties`。
- 检查 Maxwell 的日志输出，确保 Maxwell 能够连接到 MySQL 数据库并捕获变更事件。
- 确保变更事件已经转发到指定的消息队列（如 Kafka）中。

这就搭建并配置了一个基本的 Maxwell 实例，并将捕获的 MySQL 变更事件转发到消息队列中。可以根据具体需求扩展 Maxwell 的功能，例如添加过滤条件、自定义字段映射等。

## 1.8. 搭建 DataX

### 1.8.1. 下载和安装 DataX

- 在服务器上下载 DataX 的二进制文件。可以从 DataX 的官方 GitHub 仓库下载最新稳定版本。
- 解压下载的 DataX 二进制文件到一个适当的目录，例如 `/opt/datax`。
- 在 DataX 的安装目录下创建一个同步任务配置文件，例如 `mysql_to_hdfs.json`。
- 编辑配置文件，指定以下参数：
  - job.content**: 定义同步任务的具体配置，包括数据源、目标表、字段映射等。
  - job.setting**: 定义同步任务的全局设置，如并发数、错误限制等。
  - job.reader.parameter**: 配置 MySQL 数据源的连接信息、查询语句等。
  - job.writer.parameter**: 配置 HDFS 目标表的存储路径、格式等。

`mysql_to_hdfs.json`:

```
{  
  "job": {  
    "setting": {  
      "speed": {  
        "channel": 3  
      }  
    },  
    "content": [  
      {  
        "reader": {  
          "name": "mysqlreader",
```

```

        "parameter": {
            "connection": [
                {
                    "jdbcUrl":
"jdbc:mysql://mysql.example.com:3306/my_database",
                    "username": "datax_user",
                    "password": "datax_password"
                }
            ],
            "splitPk": "id",
            "column": ["id", "name", "age"],
            "where": "1=1"
        }
    },
    "writer": {
        "name": "hdfswriter",
        "parameter": {
            "path": "/data/output",
            "fileName": "output.txt",
            "defaultFS": "hdfs://hadoop.example.com:9000",
            "fileType": "text",
            "writeMode": "append"
        }
    }
}
]
}
}

```

### 1.8.2. 执行同步任务

- a. 在服务器上执行 DataX 的启动命令：/opt/datax/bin/datax.py /opt/datax/job/mysql\_to\_hdfs.json。
- b. 观察命令行输出，确保同步任务成功执行，并将数据从 MySQL 数据库全量同步到 HDFS 中。

以上步骤搭建并配置了一个基本的 DataX 实例，并成功将 MySQL 数据库中的数据全量同步到了 HDFS 中。可以根据具体需求扩展 DataX 的功能，如增量同步、数据过滤等。

## 1.9. 搭建部署 Hive 集群

### 1.9.1. 下载和安装 Hive

- a. 在每台服务器上下载相同版本的 Apache Hive。可以从 Apache Hive 官方网站下载稳定版本。
- b. 解压下载的 Hive 二进制文件到一个适当的目录，例如/opt/hive。
- c. 在 Hive 的安装目录下的 conf 目录中创建一个 Hive 配置文件，例如 hive-site.xml。
- d. 编辑配置文件，指定以下参数：
  - i. **javax.jdo.option.ConnectionURL**: 指定 Hive 元数据存储的数据库连接 URL。
  - ii. **javax.jdo.option.ConnectionUserName**: 连接元数据存储的数据库的用户名。
  - iii. **javax.jdo.option.ConnectionPassword**: 连接元数据存储的数据库的密码。
  - iv. **hive.metastore.warehouse.dir**: 指定 Hive 表数据的存储路径。
  - v. **hive.exec.scratchdir**: 指定 Hive 的临时文件目录。
  - vi. **hive.querylog.location**: 指定 Hive 查询日志文件目录。
  - vii. 其他可选配置参数，如数据压缩、执行引擎等。

**hive-site.xml:**

```
<?xml version="1.0"?>
<?xml-stylesheet type="text/xsl" href="configuration.xsl"?>
<configuration>
  <property>
    <name>javax.jdo.option.ConnectionURL</name>
    <value>jdbc:mysql://mysql.example.com:3306/hive_metastore</value>
  </property>
  <property>
    <name>javax.jdo.option.ConnectionUserName</name>
    <value>hive_user</value>
  </property>
  <property>
    <name>javax.jdo.option.ConnectionPassword</name>
    <value>hive_password</value>
  </property>
  <property>
    <name>hive.metastore.warehouse.dir</name>
    <value>/user/hive/warehouse</value>
  </property>
  <property>
    <name>hive.exec.scratchdir</name>
    <value>/tmp/hive</value>
```

```

</property>
<property>
  <name>hive.querylog.location</name>
  <value>/var/log/hive</value>
</property>
<!-- 其他配置参数 -->
</configuration>

```

## 1.9.2. 配置 Hive 集群

- 在每台服务器上编辑 Hive 的环境变量配置文件，例如 hive-env.sh，设置 Hadoop 和 Java 的路径。
- 在一个服务器上创建一个 Hive 集群配置文件，例如 hive-cluster.xml，指定 Hive 集群的元数据仓库和其他配置。
- 将 Hive 配置文件和 Hadoop 的配置文件分发到其他服务器上。

## 1.9.3. 启动 Hive Metastore

在其中一台服务器上启动 Hive Metastore 服务：\$HIVE\_HOME/bin/hive --service metastore.

## 1.9.4. 启动 HiveServer2

在其中一台服务器上启动 HiveServer2 服务：\$HIVE\_HOME/bin/hiveserver2.

# 1.10 备份与容灾

## 1.10.1. 备份

```

#!/bin/bash
# 定义备份目录和日期
backup_dir="/path/to/backup"
date=$(date +"%Y%m%d")
# 备份 HDFS 数据
hdfs_backup() {
  # 备份 HDFS 数据到本地目录
  hdfs dfs -get /path/to/hdfs/data $backup_dir/hdfs_backup_$date
}
# 备份 YARN 日志
yarn_backup() {
  # 备份 YARN 日志到本地目录
  cp -R /path/to/yarn/logs $backup_dir/yarn_backup_$date
}

```

```
# 备份 Hive 元数据
hive_backup() {
    # 备份 Hive 元数据到本地目录
    cp -R /path/to/hive/metastore $backup_dir/hive_backup_$date
}

# 备份操作
backup() {
    # 创建备份目录
    mkdir -p $backup_dir
    # 备份 HDFS 数据
    hdfs_backup
    # 备份 YARN 日志
    yarn_backup
    # 备份 Hive 元数据
    hive_backup
    # 压缩备份文件
    tar -czvf $backup_dir/backup_$date.tar.gz $backup_dir
    # 清理临时备份文件
    rm -rf $backup_dir
}

# 恢复操作
restore() {
    # 解压备份文件
    tar -xzvf $backup_dir/backup_$date.tar.gz -C $backup_dir

    # 恢复 HDFS 数据
    hdfs dfs -put $backup_dir/hdfs_backup_$date /path/to/hdfs/data
    # 恢复 YARN 日志
    cp -R $backup_dir/yarn_backup_$date /path/to/yarn/logs
    # 恢复 Hive 元数据
    cp -R $backup_dir/hive_backup_$date /path/to/hive/metastore
    # 清理临时文件
    rm -rf $backup_dir
}

# 主程序
main() {
```

```

if [ "$1" == "backup" ]; then
    backup
elif [ "$1" == "restore" ]; then
    restore
else
    echo "Usage: ./backup_script.sh [backup|restore]"
    exit 1
fi
}
# 执行主程序
main "$@"

```

### 1.10.2.容灾

```

#!/bin/bash
# 定义主节点和备节点信息
master="hostname1"
backup="hostname2"
# 检查主节点状态
check_master_status() {
    # 使用 ssh 连接到主节点，检查主节点的状态
    ssh $master "jps | grep -q 'NameNode'"

    if [ $? -eq 0 ]; then
        echo "主节点正常运行中"
        return 0
    else
        echo "主节点异常，启动备节点..."
        return 1
    fi
}
# 启动备节点
start_backup_node() {
    # 使用 ssh 连接到备节点，启动备节点的相关服务
    ssh $backup "start-dfs.sh"
    ssh $backup "start-yarn.sh"
}
# 切换为备节点

```



```

switch_to_backup_node() {
    # 使用 ssh 连接到主节点，停止主节点的相关服务
    ssh $master "stop-dfs.sh"
    ssh $master "stop-yarn.sh"

    # 使用 ssh 连接到备节点，启动备节点的相关服务
    start_backup_node

    echo "已切换为备节点"
}
# 主程序
main() {
    # 检查主节点状态
    check_master_status

    if [ $? -eq 0 ]; then
        echo "主节点正常运行中，无需切换"
    else
        echo "切换为备节点..."
        switch_to_backup_node
    fi
}
# 执行主程序
main

```

### 1.10.3 负载均衡和高可用性

```

#!/bin/bash
# 定义主节点和备节点信息
master="hostname1"
backup="hostname2"
# 检查主节点状态
check_master_status() {
    # 使用 ssh 连接到主节点，检查主节点的状态
    ssh $master "jps | grep -q 'NameNode'"

    if [ $? -eq 0 ]; then
        echo "主节点正常运行中"
    else
        echo "主节点未运行，尝试切换到备节点"
        switch_to_backup_node
    fi
}
# 执行主程序
main

```

```

    return 0
else
    echo "主节点异常，启动备节点..."
    return 1
fi
}

# 启动备节点
start_backup_node() {
    # 使用 ssh 连接到备节点，启动备节点的相关服务
    ssh $backup "start-dfs.sh"
    ssh $backup "start-yarn.sh"
}

# 切换为备节点
switch_to_backup_node() {
    # 使用 ssh 连接到主节点，停止主节点的相关服务
    ssh $master "stop-dfs.sh"
    ssh $master "stop-yarn.sh"

    # 使用 ssh 连接到备节点，启动备节点的相关服务
    start_backup_node

    echo "已切换为备节点"
}

# 监控集群状态
monitor_cluster() {
    while true; do
        # 检查主节点状态
        check_master_status

        if [ $? -ne 0 ]; then
            # 发送警报或通知，例如邮件或短信
            echo "警告：主节点异常！"
            switch_to_backup_node
        fi
        sleep 60 # 每分钟检查一次
    done
}

```

```
# 主程序
main() {
    # 启动监控集群状态
    monitor_cluster
}
# 执行主程序
main
```

## 2. 大数据模块

### 2.1. Hive 环境搭建

#### 2.1.1. Hive on Spark 配置

1. 官网下载 Hive3.1.2 源码，修改 pom 文件中引用的 Spark 版本为 3.0.0，如果编译通过，直接打包获取 jar 包。
2. 在 Hive 所在节点部署 Spark，上传并解压解压 spark-3.0.0-bin-hadoop3.2.tgz。
3. 配置 SPARK\_HOME 环境变量。
4. 在 hive 中创建 spark 配置文件，并在 HDFS 创建如下路径，用于存储历史日志。
5. 向 HDFS 上传 Spark 纯净版 jar 包
6. 修改 hive-site.xml 文件。

#### 2.1.2. Hive on Spark 测试

1. 启动 hive 客户端并创建测试表。
2. 通过 insert 测试效果。

### 2.2. Yarn 环境配置

1. 在 hadoop101 的/opt/module/hadoop-3.1.3/etc/hadoop/capacity-scheduler.xml 文件中修改参数。
2. 分发 capacity-scheduler.xml 配置文件。
3. 关闭正在运行的任务，重新启动 yarn 集群。

### 2.3. 数据仓库开发环境配置

1. 启动 HiveServer2。

2. 配置 DataGrip 连接。
3. 测试使用，创建并查看数据库。
4. 修改连接，指明连接数据库。

## 2.4. 模拟数据准备

我们组设定的数仓上线日期为 2020-01-01，为模拟真实场景，需准备数据。

### 2.4.1. 用户行为日志

1. 启动日志采集通道，包括 Flume、Kafak 等
2. 修改两个日志服务器（hadoop102、hadoop103）中的 /opt/module/applog/application.yml 配置文件，将 mock.date 参数改为 2020-01-01。
3. 执行日志生成脚本 lg.sh。
4. 观察 HDFS 是否出现相应文件。

### 2.4.2. 业务数据

业务数据准备的是 2020-01-01 至 2020-04-09 的数据。

1. 生成模拟数据，修改 hadoop102 节点上的 /opt/module/db\_log/application.properties 文件，将 mock.date、mock.clear，mock.clear.user 三个参数调整为如图所示的值。
2. 执行模拟生成业务数据的命令，生成第一天 2020-01-01 的历史数据。
3. 修改 /opt/module/db\_log/application.properties 文件，调整 mock.date、mock.clear，mock.clear.user 三个参数的值。
4. 执行模拟生成业务数据的命令，生成第二天 2020-01-02 的历史数据。
5. 之后只修改 /opt/module/db\_log/application.properties 文件中的 mock.date 参数，依次改为生成对应日期的数据。

### 2.4.3. 全量表同步

1. 执行全量表同步脚本，观察 HDFS 上是否出现全量表数据。
2. 增量表首日全量同步。
3. 清除 Maxwell 断点记录。由于 Maxwell 支持断点续传，而上述重新生成业务数据的过程，会产生大量的 binlog 操作日志，这些日志我们并不需要。故此处需清除 Maxwell 的断点记录，另其从 binlog 最新的位置开始采集。
4. 关闭 Maxwell，清空 Maxwell 数据库，相当于初始化 Maxwell。
5. 修改 Maxwell 配置文件中的 mock\_date 参数。

6. 启动增量表数据通道，包括 Maxwell、Kafka、Flume。
7. 执行增量表首日全量同步脚本，观察 HDFS 上是否出现全量表数据。

## 3. 人工智能模块

### 3.1 用户画像分析

#### 3.1.1 设计思路

我们将所有的用户分为四类：

- 新用户 (New Users)：这个群体包括那些刚注册或者刚开始使用服务的用户。他们可能对平台的功能和操作还不熟悉，因此需要更多的引导和教程。同时，他们可能还未形成固定的购买习惯，所以需要更多的优惠和激励措施来促使他们进行第一次购买。对于新用户，我们可以进行新用户引导教程，提供新用户专属优惠券，推荐热销商品，通过这些方式引导他们熟悉平台并进行初次购买。
- 活跃用户 (Active Users)：这个群体的用户在平台上有频繁的活动，比如经常浏览、搜索和购买商品。他们可能是平台的忠实用户，对平台的运作和商品非常熟悉。对于这些用户，我们需要保持高水平的服务品质，并提供他们喜欢的商品。对于活跃用户，我们可以进行个性化推荐，比如推荐他们经常购买的品类或者品牌，为他们定制优惠活动，如购买多件折扣，积分回馈等，通过这些方式保持他们的活跃度和满意度。
- 沉默用户 (Silent Users)：这类用户虽然已经注册，但是很少在平台上有活动。他们可能只是偶尔浏览，或者在需要时才会使用平台。这可能是因为他们没有找到自己喜欢的商品，或者是因为他们更喜欢在其他平台购买。对于这类用户，我们需要重新吸引他们的注意力，唤醒他们的购买欲望。比如，我们可以通过电子邮件或短信通知他们最新的优惠活动，或者推荐他们可能感兴趣的物品。
- 流失用户 (Churned Users)：这类用户在一段时间内完全没有活动，可能是因为他们找到了其他的服 务，或者对当前的服务不满意。这可能是因为我们 的服务质量不足，或者是因为他们对我们的商品不满意。对于这类用户，我们需要找出他们离开的原因，并尝试挽回他们。我们可以通过问卷调查，了解他们的不满意点，并根据他们的反馈改进我们的服务和商品。同时，我们可以通过一些挽回活动，如提供折扣券或礼品，来吸引他们回到我们的平台。

总的来说，每一类用户都有自己的特性和需求，我们需要根据这些特性和需求，采取相应的措施，以提供更好的服务，满足他们的需求。

在实现用户画像分析功能时，我们将使用 KNN (K-Nearest Neighbors) 聚类算法，这是一种非参数的、基于实例的学习方法。它通过计算新样本与已有样本之间的距离，找出 K 个最近的邻居，并根据这些邻居的类别进行预测。

在此场景中，我们首先需要理解数据中包含的信息。我们拥有用户的省份 ID、性别、行为，以及在过去 7 天和 30 天内的购买信息（如购买品类、品牌、单数和消费总金额等）。这些数据都是有助于形成用户画像的重要因素。

设计思路如下：

1. 数据预处理：首先，我们需要对数据进行清洗，包括处理缺失值、异常值、重复值等。接着，进行特征工程，例如：将购买品类和品牌进行独热编码，将连续变量（如购买单数、消费总金额）进行归一化处理，以减少不同特征尺度间的影响。
2. 用户分类：然后，利用 KNN 聚类算法进行用户分类。我们可以根据上述定义的用户类型，如新用户、活跃用户、沉默用户和流失用户，进行分类。这可能需要我们事先设定一些阈值，如用户在一定时间内的行为频率、购买次数等。
3. 模型训练：接着，我们需要将处理后的数据输入 KNN 模型进行训练。在此过程中，可能需要调整一些参数，如选择最佳的邻居数量  $K$ ，以实现最佳的分类效果。
4. 用户画像创建：最后，根据聚类结果，我们可以为每个用户生成一个个性化的用户画像。用户画像应包含用户的基本信息（如省份、性别）和行为特征（如购买品类、品牌、购买频率、消费金额等）。

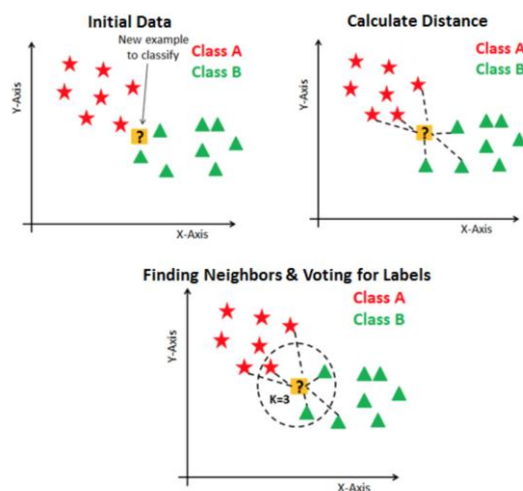
这样的设计不仅可以帮助我们更好地理解 and 分类用户，还可以根据用户的不同特征，提供更精细化、个性化的服务。例如，对于新用户，我们可以提供更多的引导和教程，以及优惠券或者促销活动；对于活跃用户，我们可以推荐他们经常购买的品类和品牌；对于沉默用户，我们可以通过发送推广信息或者优惠活动，尝试唤醒他们的购买欲望；对于流失用户，我们可以通过调查或者反馈，了解他们离开的原因，从而改善服务和产品。

### 3.1.2 算法原理

首先，值得注意的是，KNN (K-Nearest Neighbors) 算法其实是一种用于分类和回归的监督学习算法，并非常规意义上的聚类算法。但是，KNN 算法的思想可以被用于聚类，常见的实现方法如基于密度的聚类方法 DBSCAN。为了回答你的问题，我将先介绍 KNN 的基础概念和工作原理，然后再解释如何将其用于聚类。

KNN 算法的基本概念：

KNN 算法，即 K-最近邻算法，是一种基于实例的学习，或者称为懒散学习的方法。算法中的  $K$  代表在预测时要参考的邻近样本的数量。这个算法的主要思想是假定相似的样本存在于相似的类别中。如果我们有一个未标记的样本，我们会看它的“邻居”（在特征空间中距离最近的样本）来确定它的类别。



K-近邻（KNN）算法是一种基于实例的学习，或者称为懒惰学习的方法，常用于分类和回归任务。以下是 KNN 算法的基本步骤：

1. 选择邻居数 K: K 是用户指定的正整数，常通过交叉验证等方式选择。
2. 计算未知样本与训练集中每个样本的距离: 通常使用欧氏距离（Euclidean distance）或曼哈顿距离（Manhattan distance）来计算距离。但实际上，可以根据数据的特性选择适当的距离计算方法。
3. 排序: 将计算得到的所有距离进行排序，得到一个按距离大小排序的序列。
4. 确定最近的 K 个样本: 根据排序的结果，选取距离最近的 K 个样本。这 K 个样本就是我们所说的 K 个“最近邻”。
5. 投票决定未知样本的类别: 对于分类问题，一般采取多数投票法，即未知样本的类别由最近邻的 K 个样本中出现次数最多的类别确定；对于回归问题，未知样本的值由最近邻的 K 个样本的平均值（或中位数）确定。

## 3.2 用户性别预测

### 3.2.1 设计思路

#### 1. 数据收集和准备：

- 收集用户的购买数据，包括品类、品牌、购买单数和消费总金额等信息。这些数据可以从数据仓库中获取。
- 将用户购买数据按照时间窗口进行划分，例如选择最近 30 天的购买数据作为训练数据集。
- 根据已知的用户性别数据，将用户购买数据和对应的性别标签进行关联，构建带有性别标签的训练数据集。

#### 2. 特征提取：

- 从用户购买数据中提取特征。在这个设计中，我们可以选择最近 30 天内买得最多的品类、第二多的品类、第三多的品类、最多的品牌、第二多的品牌、第三多的品牌、购买单数和消费总金额作为特征。

- 对于品类和品牌，可以将其转化为离散的特征表示，例如使用独热编码将每个品类和品牌表示为二进制向量。

### 3. 数据预处理：

- 对于连续型的特征（购买单数和消费总金额），可以进行数据标准化或归一化处理，确保不同特征之间的数值范围相同。
- 对于离散型的特征（品类和品牌），可以进行特征编码，例如独热编码。

### 4. 模型训练：

- 使用准备好的训练数据集进行朴素贝叶斯模型的训练。
- 计算每个类别的先验概率，即在训练数据集中每个性别的出现频率。
- 计算每个特征在每个类别下的条件概率，即在每个性别下，每个特征取某个值的频率。

### 5. 模型预测：

- 对于一个新的用户，根据其购买行为提取相应的特征。
- 使用训练好的朴素贝叶斯模型，根据先验概率和条件概率计算用户属于每个性别的概率。
- 根据计算得到的概率，选择概率最高的性别作为预测结果。

### 6. 模型评估和优化：

- 使用一部分已知性别的测试数据集对模型进行评估，计算准确率、精确率、召回率等指标，了解模型的性能。
- 如果模型性能不够理想，可以考虑调整特征的选择、数据预处理方法，或者使用其他分类算法进行比较。

用户性别预测功能基于朴素贝叶斯算法，通过用户的购买行为作为特征，利用训练数据集计算先验概率和条件概率，实现对用户性别的预测。

## 3.2.2 算法原理

朴素贝叶斯算法是一种基于概率统计的分类算法，它基于贝叶斯定理和特征条件独立性的假设，通过计算先验概率和条件概率来进行分类。下面详细介绍朴素贝叶斯算法的原理：

### 1. 贝叶斯定理：

贝叶斯定理是概率论中的重要定理，它描述了在给定先验知识的情况下，如何根据新的观测数据来更新对事件发生概率的估计。对于两个事件 A 和 B，贝叶斯定理可以表达为：

$$P(A|B) = (P(B|A) * P(A)) / P(B)$$

其中， $P(A|B)$  表示在事件 B 发生的条件下事件 A 发生的概率， $P(B|A)$  表示在事件 A 发生的条件下事件 B 发生的概率， $P(A)$  和  $P(B)$  分别表示事件 A 和 B 单独发生的概率。

### 2. 朴素贝叶斯假设：

朴素贝叶斯算法假设所有的特征之间是相互独立的，即每个特征对于分类的贡献是独立的。这个假设使得朴素贝叶斯算法计算简单，但在现实问题中可能不符合实际情况。



### 3. 模型训练：

朴素贝叶斯算法的模型训练过程主要包括先验概率的计算和条件概率的计算。

- 先验概率：

先验概率是指在没有任何特征信息的情况下，每个类别发生的概率。在训练阶段，通过统计训练集中每个类别出现的频率来计算先验概率。

- 条件概率：

条件概率是指在已知特征的条件下，某个类别发生的概率。在训练阶段，对于每个特征，需要计算在每个类别下，该特征取某个值的概率。如果特征是离散的，可以直接计算各个类别下特征值的频率；如果特征是连续的，可以假设特征值服从某种概率分布，如高斯分布，然后通过训练数据来估计分布的参数。

### 4. 模型预测：

在预测阶段，给定一个新的样本，朴素贝叶斯算法根据贝叶斯定理计算每个类别的后验概率，然后选择具有最高后验概率的类别作为预测结果。

- 特征提取：

对于新样本，需要从样本中提取特征，这些特征应该与训练阶段使用的特征一致。

- 后验概率计算：

对于每个类别，根据训练阶段计算得到的先验概率和条件概率，使用贝叶斯定理计算后验概率。对于多个特征的情况下，可以使用条件独立性假设，将条件概率的计算拆分为单个特征的条件概率的乘积。

- 预测结果：

根据计算得到的后验概率，选择具有最高后验概率的类别作为预测结果。

### 5. 模型评估和优化：

在朴素贝叶斯算法中，常用的模型评估指标包括准确率、精确率、召回率等。可以使用交叉验证等方法来评估模型的性能，并进行模型的优化，例如调整特征的选择、数据的预处理方法或使用更复杂的贝叶斯模型。

朴素贝叶斯算法是一种基于概率统计的分类算法，通过贝叶斯定理和特征条件独立性假设，计算先验概率和条件概率来进行分类。它具有简单、高效的特点，在许多实际应用中取得了良好的效果。然而，朴素贝叶斯算法的特征条件独立性假设可能与实际问题不符，因此在使用时需要根据具体情况进行适当调整和改进。

## 3.3 用户购买意向分析

### 3.3.1 设计思路

用户购买意向分析功能旨在根据用户在数据仓库中的相关数据进行分析，以预测用户的购买意向。为了实现这一功能，可以考虑使用决策树和随机森林两种机器学习算法。

设计思路如下：

#### 1. 数据收集和预处理：

- 从数据仓库中获取所需的用户行为数据，包括用户每次浏览的平均时长、用户浏览过的不同商品种类的数量、用户对平台热门商品的浏览次数占总浏览次数的比例、用户每次搜索的平均关键词数量、购物车中商品结算数量与加入数量的比例等信息。
- 对数据进行清洗和预处理，包括处理缺失值、异常值和重复值，将数据转换为可用于机器学习算法的数值型特征。

#### 2. 特征选择：

- 根据领域知识和经验，选择对购买意向有潜在影响的特征。例如，用户浏览过的不同商品种类的数量可能反映用户的兴趣广泛程度，用户对平台热门商品的浏览次数占总浏览次数的比例可能反映用户对热门商品的偏好程度等。
- 还可以使用特征选择算法，如信息增益、卡方检验等，从原始特征中选择最相关的特征。

#### 3. 数据划分：

- 将数据集划分为训练集和测试集。训练集用于模型的训练和参数调整，测试集用于评估模型的性能。

#### 4. 决策树模型：

- 使用训练集数据构建决策树模型。决策树是一种树形结构，每个内部节点表示一个特征，每个叶节点表示一个类别或决策结果。
- 在构建决策树时，可以使用不同的划分准则，如基尼系数或信息增益，选择最佳的特征和划分点。
- 针对决策树的过拟合问题，可以通过剪枝等方法进行优化，以避免模型过于复杂而导致泛化性能下降。

#### 5. 随机森林模型：

- 随机森林是一种集成学习方法，通过同时构建多个决策树，并通过投票或平均的方式进行预测。
- 在随机森林中，每个决策树的构建使用不同的训练数据子集和特征子集，以增加模型的多样性和泛化能力。
- 随机森林模型可以通过调整决策树的数量、最大树深度等超参数进行优化。

#### 6. 模型训练和评估：

- 使用训练集对决策树和随机森林模型进行训练。
- 使用测试集对模型进行评估，计算预测准确率、召回率、F1 值等指标，评估模型的性能和泛化能力。

#### 7. 模型应用：

- 完成模型的训练和评估后，可以将模型应用于新的数据，预测用户的购买意向。
- 根据模型的预测结果，可以制定相应的营销策略，例如向购买意向较高的用户推送个性化的推荐商品、优惠券等。

需要注意的是，在实际应用中，还需考虑以下因素：

- 数据质量和可靠性：确保数据的准确性和完整性，避免错误的对模型训练和预测结果的影响。
- 模型的更新和迭代：随着时间推移和新数据的积累，需要定期更新和优化模型，以保持其准确性和可靠性。
- 隐私和安全性：在使用用户数据进行分析时，需要遵守相关的隐私政策和法规，并采取必要的安全措施保护用户的个人信息。

通过这些步骤，可以利用机器学习算法对用户行为数据进行分析，预测用户的购买意向，并为营销决策提供有价值的参考。

### 3.3.2 算法原理

用户购买意向分析功能可以使用决策树和随机森林两种机器学习算法进行实现。下面将详细介绍这两种算法的原理。

#### 1. 决策树算法原理：

决策树是一种基于树形结构的机器学习算法，用于分类和回归问题。它通过构建一系列的决策规则，将数据集划分为不同的类别或取值，从而实现预测和分类的功能。

##### - 决策树的构建过程：

决策树的构建过程从根节点开始，根据某个特征对数据集进行划分，将数据分配到相应的子节点。该划分过程基于一个划分准则，如基尼系数或信息增益，以选择最佳的特征和划分点。然后，对每个子节点递归地重复这个划分过程，直到达到终止条件，如达到叶节点或达到最大树深度等。

##### - 决策树的预测过程：

对于新的数据样本，通过沿着决策树的分支路径，根据样本的特征值来到达叶节点。叶节点对应着一个类别或取值，即预测的结果。

##### - 决策树的优势和注意事项：

决策树算法具有易于理解、解释和可视化的特点。它可以处理数值型和类别型的特征，并且在一定程度上具有抗噪性。然而，决策树容易产生过拟合问题，即模型过于复杂而在训练集上表现很好，但在新数据上的泛化能力较差。为了克服过拟合，可以采用剪枝等技术，或者结合其他方法，如随机森林。

#### 2. 随机森林算法原理：

随机森林是一种集成学习方法，通过构建多个决策树并进行集成来进行预测和分类。随机森

林综合了决策树的优势，并通过随机化和投票的方式提高模型的鲁棒性和泛化能力。

- 随机森林的构建过程：

随机森林的构建过程包括两个主要的随机性来源：样本随机性和特征随机性。对于每棵决策树的构建，从原始数据集中通过有放回的采样方式随机选择一个样本子集。同时，对于每个节点的特征选择，在特征集合中随机选择一部分特征作为候选特征，再从中选择最佳的特征进行划分。

- 随机森林的预测过程：

对于新的数据样本，通过对随机森林中的每棵决策树进行预测，根据投票或平均的方式来确定最终的预测结果。分类问题中，采用投票的方式选择类别；回归问题中，采用平均的方式得到预测值。

- 随机森林的优势和注意事项：

随机森林在决策树的基础上引入了随机性，减少了模型的方差和过拟合风险，提高了模型的泛化能力。随机森林能够处理大规模的数据集，并且对于缺失值和异常值具有较好的鲁棒性。然而，随机森林的模型解释性较差，而且构建和预测的计算开销相对较大。

针对购买意向分析功能，可以使用决策树和随机森林这两种算法进行建模和预测。决策树算法通过构建一系列的决策规则来实现预测和分类，而随机森林算法通过构建多个决策树并进行集成来提高模型的鲁棒性和泛化能力。

## 3.4 售后服务满意度分析

### 3.4.1 设计思路

使用循环神经网络（RNN）实现售后服务满意度分析功能是一种常见的方法。RNN 是一种适用于序列数据建模的神经网络，能够捕捉文本中的上下文信息和长期依赖关系，因此在处理自然语言处理任务中广泛应用。下面将详细介绍使用 RNN 实现售后服务满意度分析功能的算法原理。

#### 1. 数据准备：

首先，需要准备一个包含用户评价和标签（满意、中性、不满意）的训练数据集。数据集应包括经过预处理的用户评价文本，以及对应的标签。

#### 2. 文本表示：

在使用 RNN 进行训练之前，需要将文本数据转换为适合 RNN 处理的表示形式。常见的文本表示方法包括词嵌入（Word Embedding）和整数编码（Integer Encoding）：

- 词嵌入：通过将每个词语映射到一个高维向量空间中的实数向量，捕捉词语之间的语义关系。可以使用预训练的词嵌入模型，如 Word2Vec、GloVe 等，也可以在训练过程中学习得到。

- 整数编码：将每个词语编码为一个整数，构成一个整数序列。可以使用词频或其他统计信息来为词语赋予整数编码。

### 3. 序列填充：

RNN 需要输入具有相同长度的序列数据，但用户评价文本往往长度不一致。因此，需要对文本序列进行填充（Padding）或截断（Truncation），使它们具有相同的长度。常见的做法是根据最长的文本长度，将其他文本进行填充或截断为相同的长度。

### 4. 构建 RNN 模型：

在构建 RNN 模型之前，需要选择 RNN 的类型，如简单循环神经网络（Simple RNN）、长短期记忆（LSTM）或门控循环单元（GRU）等。

### 5. 模型训练：

在完成 RNN 模型的构建后，需要将准备好的训练数据集输入模型进行训练。训练的过程包括以下步骤：

- 前向传播：将训练数据输入到 RNN 模型中，通过前向传播计算模型的预测输出。
- 计算损失：将模型的预测输出与真实标签进行比较，计算损失（损失函数可以选择交叉熵损失等）。
- 反向传播：根据损失函数的梯度，通过反向传播更新模型的参数，以最小化损失函数。
- 参数优化：使用优化算法（如梯度下降）对模型参数进行更新，逐步优化模型的性能。
- 迭代训练：重复进行前向传播、损失计算、反向传播和参数优化的步骤，直到达到指定的训练轮次或收敛条件。

### 6. 模型评估：

在模型训练完成后，需要对模型进行评估以确定其在新数据上的性能。可以使用验证集或交叉验证等方法，计算模型在新数据上的准确率、精确率、召回率、F1 分数等指标。

### 7. 预测与应用：

在模型经过评估后，可以使用该模型对新的用户评价文本进行情感分类预测。将新的文本输入训练好的 RNN 模型中，通过前向传播得到模型的预测输出。根据预测输出的概率或阈值，将文本归类为满意、中性或不满意等情感类别。

通过以上步骤，使用 RNN 可以对售后服务的满意度进行分析和预测。RNN 模型能够捕捉文本序列中的上下文信息和长期依赖关系，有助于更准确地理解用户评价的情感倾向。但在实际应用中，还需要根据具体场景和数据特点进行调优和改进，以提高算法的性能和实用性。

## 3.4.2 算法思路

LSTM（Long Short-Term Memory，长短期记忆）是一种常用于处理序列数据的循环神经网络（RNN）的变种。相比于传统的 RNN 结构，LSTM 通过引入门控机制，有效地解决了传统 RNN 中的梯度消失和梯度爆炸问题，能够更好地捕捉和处理序列数据中的长期依赖关系。下面将详细介绍 LSTM 的结构和工作原理。

### 1. LSTM 结构:

LSTM 由一个或多个 LSTM 单元组成, 每个 LSTM 单元内部包含多个关键组件:

- 输入门 (Input Gate): 决定哪些信息将被加入到 LSTM 单元的内部状态中。
- 遗忘门 (Forget Gate): 决定哪些信息将被从 LSTM 单元的内部状态中遗忘。
- 输出门 (Output Gate): 决定哪些信息将从 LSTM 单元的内部状态输出。
- 细胞状态 (Cell State): 负责存储和传递长期依赖关系的信息。
- 隐藏状态 (Hidden State): 根据输入和前一个时间步的隐藏状态计算得到, 用于传递 LSTM 单元的输出。

### 2. 工作原理:

在每个时间步, LSTM 接收当前时间步的输入和前一个时间步的隐藏状态作为输入, 并通过一系列的门控机制来控制信息的流动和处理。

- 输入门: 通过输入门决定哪些信息将被添加到细胞状态中。输入门的计算公式为:
  - 输入门 =  $\text{sigmoid}(\text{输入} \times \text{权重} + \text{偏置})$
- 遗忘门: 通过遗忘门决定哪些信息将从细胞状态中遗忘。遗忘门的计算公式为:
  - 遗忘门 =  $\text{sigmoid}(\text{输入} \times \text{权重} + \text{偏置})$
- 细胞状态更新: 通过输入门和遗忘门的控制, 更新细胞状态的值。更新公式为:
  - 细胞状态 = 遗忘门  $\odot$  细胞状态 + 输入门  $\odot \tanh(\text{输入} \times \text{权重} + \text{偏置})$
- 输出门: 通过输出门决定哪些信息将从细胞状态输出。输出门的计算公式为:
  - 输出门 =  $\text{sigmoid}(\text{输入} \times \text{权重} + \text{偏置})$
- 隐藏状态更新: 根据输出门和更新后的细胞状态, 计算当前时间步的隐藏状态。更新公式为:
  - 隐藏状态 = 输出门  $\odot \tanh(\text{细胞状态})$

### 3. 长期依赖关系处理:

LSTM 通过细胞状态的传递和控制, 能够有效地处理序列数据中的长期依赖关系。细胞状态充当了一个内部存储器, 能够保留和传递关键的信息, 避免了传统 RNN 中梯度消失和梯度爆炸的问题。遗忘门和输入门的控制使得 LSTM 能够选择性地遗忘或添加信息到细胞状态中, 而输出门的控制则决定哪些信息将从细胞状态输出, 以用于当前时间步的预测或下一时间步的计算。

### 4. 多层 LSTM:

LSTM 可以通过堆叠多个 LSTM 层来增加模型的表达能力和复杂度。在多层 LSTM 中, 每一层的隐藏状态作为下一层的输入。通过多层 LSTM 的堆叠, 模型能够更好地学习和捕捉数据中的复杂关系和模式。

## 3.5 产品销售趋势分析

### 3.5.1 设计思路

#### 1. 数据收集与预处理：

首先，需要从数据仓库中获取每个月的产品销售额数据以及其他相关数据，包括历史数据计算得出的季节性指数、产品类别在总销售额中的占比以及最近 7、30 日的订单、退单等数据。

#### 2. 数据清洗与转换：

对获取的数据进行清洗和转换，确保数据的完整性和准确性。这包括处理缺失值、异常值和重复值等。同时，将数据转换为适合时间序列分析的格式，例如将日期作为时间索引。

#### 3. 季节性指数计算：

使用历史数据计算季节性指数。季节性指数反映了产品销售在不同季节的波动情况，可以帮助我们理解销售趋势中的季节性因素。可以使用加法模型或乘法模型计算季节性指数。

#### 4. 时间序列分析算法选择：

选择合适的时间序列分析算法，如 ARIMA 模型或指数平滑模型，用于分析产品销售趋势。ARIMA 模型适用于具有明显趋势和季节性的数据，指数平滑模型则适用于较平滑的数据。

#### 5. 模型训练与参数调优：

根据历史数据，训练选定的时间序列分析模型，并通过参数调优来提高模型的准确性和预测能力。这可以使用交叉验证等技术来评估模型的性能，并进行参数的优化选择。

#### 6. 趋势分析与预测：

使用训练好的时间序列分析模型，对产品销售趋势进行分析和预测。可以通过模型预测未来时间段内的销售额，并进一步分析销售趋势的走向和变化。

#### 7. 复购率分析：

计算最近 7、30 日的复购率，即在特定时间段内再次购买同一产品的比例。可以根据订单数据来计算复购率，考虑用户 ID 和购买日期等信息。

#### 8. 订单数、订单人数、退单数、退单人数分析：

分析最近 1、7、30 日的订单数、订单人数、退单数和退单人数的变化趋势，可以使用时间窗口来获取特定时间段内的数据，并进行统计和计算。

#### 9. 订单总额分析：

分析最近 1、7、30 日的订单总额的变化趋势，可以根据订单数据计算订单总额，并使用时间窗口来获取特定时间段内的数据。

#### 10. 结果可视化与报告输出：

将分析和预测的结果可视化展示，例如绘制销售趋势图、复购率曲线、订单数和退单数的变化图等。同时，生成相应的报告输出，以便用户更好地理解 and 利用销售趋势分析的结果。

### 3.5.2 算法思路

ARIMA（自回归滑动平均模型，Autoregressive Integrated Moving Average）模型是一种常用的时间序列分析方法，用于预测和建模具有趋势和季节性的数据。ARIMA 模型基于时间序列的自相关性（AR 部分）和移动平均性（MA 部分），并结合差分运算（I 部分）来处理非平稳性。

ARIMA 模型的原理如下：

#### 1. 自回归（AR）部分：

自回归是指时间序列在当前时刻的值与之前时刻的值之间存在关联性。AR 部分使用过去时刻的观测值作为自变量来预测当前时刻的观测值。AR(p)模型中，p 表示使用的过去时刻的观测值个数，模型公式为： $Y(t) = c + \varphi(1)Y(t-1) + \varphi(2)Y(t-2) + \dots + \varphi(p)Y(t-p) + \varepsilon(t)$ ，其中 c 为常数， $\varphi$  为模型参数， $\varepsilon$  为误差项。

#### 2. 移动平均（MA）部分：

移动平均是指时间序列在当前时刻的值与过去时刻的误差项之间存在关联性。MA 部分使用过去时刻的误差项作为自变量来预测当前时刻的观测值。MA(q)模型中，q 表示使用的过去时刻的误差项个数，模型公式为： $Y(t) = c + \theta(1)\varepsilon(t-1) + \theta(2)\varepsilon(t-2) + \dots + \theta(q)\varepsilon(t-q) + \varepsilon(t)$ ，其中 c 为常数， $\theta$  为模型参数， $\varepsilon$  为误差项。

#### 3. 差分运算（I 部分）：

差分运算是为了处理非平稳性的时间序列。非平稳性指时间序列的均值和方差随时间变化，无法在时间上稳定。通过对时间序列进行差分，可以将非平稳时间序列转化为平稳时间序列。I 部分使用差分后的时间序列作为自变量，模型公式为： $Y'(t) = Y(t) - Y(t-1)$ ，其中  $Y'(t)$  表示差分后的时间序列。

#### 4. ARIMA 模型的组合：

ARIMA 模型将 AR、MA 和 I 部分结合起来，形成一个综合的时间序列模型。ARIMA(p, d, q)模型中，p 表示 AR 部分的阶数，d 表示差分的阶数，q 表示 MA 部分的阶数。ARIMA 模型的公式为： $Y'(t) = c + \varphi(1)Y'(t-1) + \varphi(2)Y'(t-2) + \dots + \varphi(p)Y'(t-p) + \theta(1)\varepsilon(t-1) + \theta(2)\varepsilon(t-2) + \dots + \theta(q)\varepsilon(t-q) + \varepsilon(t)$ ，其中 c 为常数， $\varphi$  和  $\theta$  为模型参数， $\varepsilon$  为误差项。

#### 5. 参数估计和模型选择：

在实际应用中，需要对 ARIMA 模型的参数进行估计和选择。可以使用最大似然估计（Maximum Likelihood Estimation）或其他优化算法来估计模型参数。同时，可以使用信息准则（如 AIC、BIC）来选择最优的 ARIMA 模型，以获得更好的预测性能。



## 3.6 商品热度分析

### 3.6.1 设计思路

实现商品热度分析功能可以使用 KNN (K-Nearest Neighbors) 聚类算法。KNN 算法是一种常用的无监督学习算法，用于将数据点分为不同的簇群。在商品热度分析中，我们可以使用 KNN 算法将商品分为不同的热度级别，从而帮助我们了解不同商品的受欢迎程度和销售潜力。以下是实现商品热度分析功能的设计思路：

#### 1. 数据预处理：

首先，我们需要从数据仓库中获取所需的数据，包括销售增长率、商品页面访问量以及各项统计数据（最近 7、30 日复购率、最近 1、7、30 日订单数等）。对于缺失值和异常值，需要进行处理，可以选择删除或填充缺失值，或者使用合适的统计方法进行修正。

#### 2. 特征选择：

在设计热度分析模型时，需要选择合适的特征。根据我们的目标，可以选取销售增长率、商品页面访问量和各项统计数据作为特征。这些特征能够反映商品的销售表现、用户行为和用户反馈等方面。

#### 3. 特征标准化：

由于选取的特征可能具有不同的量纲和取值范围，为了确保它们对聚类结果的影响具有平等的权重，需要进行特征标准化。一种常用的方法是使用标准化 (Standardization) 或归一化 (Normalization) 技术，将所有特征值映射到相似的数值范围内。

#### 4. 聚类算法选择：

在本案例中，我们选择 KNN 算法作为聚类算法。KNN 算法是一种简单有效的聚类算法，它通过计算数据点之间的距离来确定最近的邻居，并将相似的数据点归为一类。KNN 算法可以根据商品的销售表现和用户行为等特征将商品进行分组。

#### 5. 确定簇数：

在 KNN 算法中，需要事先确定聚类的簇数。可以使用 Elbow 方法或 Silhouette 系数等评估指标来选择合适的簇数。Elbow 方法通过绘制簇数和聚类误差之间的关系图，找到拐点处的簇数。Silhouette 系数衡量了簇内的紧密度和簇间的分离度，选择 Silhouette 系数最大的簇数。

#### 6. 模型训练和预测：

在确定了簇数后，使用 KNN 算法对数据进行训练。训练过程中，算法将根据数据点之间的距离计算确定最近的邻居，并将相似的数据点归为同一簇。训练完成后，可以使用该模型对新的商品数据进行预测，将其归为相应的热度级别簇群。

#### 7. 热度级别分析：

根据聚类结果，可以将商品划分为不同的热度级别，如高热度、中热度和低热度。可以根据销售增长率、商品页面访问量等指标对不同簇群的商品进行详细的分析，了解其销售情况和用

户行为，制定相应的营销策略。

#### 8. 结果可视化：

最后，可以通过可视化工具将热度分析的结果呈现出来，以直观地展示不同簇群的商品分布情况。可以使用散点图、热力图或雷达图等方式，将商品在不同热度级别下的特征进行可视化展示，帮助决策者更好地理解和分析数据。

通过使用 KNN 聚类算法实现商品热度分析功能，我们可以根据商品的销售表现、用户行为和用户反馈等特征将商品进行分组，并将其划分为不同的热度级别。这可以帮助我们更好地了解商品的受欢迎程度和销售潜力，为制定相应的营销策略提供参考。同时，合适的预处理、特征选择、特征标准化和簇数确定等步骤，以及结果的可视化展示，都对于实现有效的商品热度分析功能至关重要。

### 3.6.2 算法思路

见 3.1.2

## 3.7 优惠券效果分析

### 3.7.1 设计思路

#### 1. 确定目标和指标：

首先，我们需要明确优惠券效果分析的目标。通常，优惠券的效果分析旨在评估优惠券对用户消费行为的影响，例如优惠券的使用率、促销效果以及与订单价格之间的关联。在此基础上，我们可以选择以下指标来衡量优惠券的效果：

- 优惠券使用率：衡量用户实际使用优惠券的比例。
- 平均优惠价格：计算每个订单中优惠券的平均抵扣金额。
- 优惠券转化率：衡量通过优惠券促成的购买行为的比例。
- 订单价格与优惠价格之间的关联：分析订单价格和应用优惠券后的价格之间的关系。

#### 2. 数据准备：

为了实现优惠券效果分析，我们需要收集和整理以下数据：

- 优惠券信息：包括优惠券的类型、折扣金额或折扣比例等。
- 优惠券使用记录：记录哪些用户使用了哪些优惠券，以及使用时间等。
- 用户消费历史：包括用户的购买记录、订单价格等。
- 用户行为记录：记录用户的点击、浏览、加入购物车等行为，用于分析用户的购买意愿和行为偏好。
- 订单价格：记录订单的原始价格。
- 优惠价格：记录优惠券应用后的价格。
- 优惠后订单价格：记录应用优惠券后的订单价格。
- 优惠券使用数量：记录每个优惠券的使用数量。

### 3. 数据清洗和整合：

在数据准备完成后，我们需要对数据进行清洗和整合，以便后续分析使用。这包括处理缺失值、异常值、重复值，对数据进行格式化和标准化等操作。同时，将不同数据源的数据进行整合，建立关联关系，以便进行后续的分析。

### 4. 特征工程：

在进行回归分析之前，我们需要对数据进行特征工程处理，以提取和构建适合回归模型的特征。根据前面确定的指标，可能需要进行以下操作：

- 用户特征提取：从用户消费历史和行为记录中提取有关用户的特征，如用户的购买频率、购买金额、点击次数等。
- 优惠券特征提取：从优惠券信息中提取相关特征，如优惠券类型、折扣金额等。
- 订单特征提取：从订单数据中提取特征，如订单价格。
- 目标特征构建：根据前面确定的指标，构建相应的目标特征，如优惠券使用率、优惠价格等。

### 5. 回归模型建立和训练：

在特征工程完成后，我们可以选择适合的回归分析算法建立模型。常见的回归算法包括线性回归、决策树回归、随机森林回归等。根据实际情况选择合适的算法，并使用训练集对模型进行训练和调优，以获得较好的预测效果。

### 6. 模型评估和解释：

完成模型的训练后，需要对模型进行评估和解释。可以使用评估指标如均方误差（MSE）、决定系数（R-squared）等来评估模型的拟合程度和预测效果。此外，还可以通过模型的系数（如线性回归的系数）来解释各个特征对于目标的影响程度。

### 7. 结果可视化和报告生成：

最后，根据分析结果，可以进行结果的可视化展示和报告生成，以便更好地理解 and 传达优惠券的效果分析结果。可以使用图表、统计摘要等方式呈现分析结果，帮助业务决策和优化。

## 3.7.2 算法思路

线性回归是一种广泛应用于预测和建模的统计方法，其原理基于线性关系的假设。线性回归的目标是通过拟合一个线性方程来描述自变量（输入特征）和因变量（输出）之间的关系。下面是线性回归的原理：

#### 1. 假设函数：

线性回归假设因变量  $Y$  和自变量  $X$  之间存在线性关系，可以用以下的线性方程来表示：

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n + \varepsilon$$

其中， $Y$  是因变量， $X_1, X_2, \dots, X_n$  是自变量， $\beta_0, \beta_1, \beta_2, \dots, \beta_n$  是待估计的系数（也称为回归系数）， $\varepsilon$  是误差项，表示模型无法完全解释的随机噪声。

#### 2. 最小二乘法：

线性回归的目标是找到最优的回归系数，使得观测数据与预测值之间的误差平方和最小化。这称为最小二乘法。具体来说，我们通过最小化残差平方和（Sum of Squared Residuals, SSR）来求解回归系数的最优值。

### 3. 模型参数估计：

为了求解最优的回归系数，可以使用最小二乘法来估计模型参数。最小二乘法的目标是最小化观测数据与模型预测值之间的差距，即最小化误差的平方和。

### 4. 参数估计方法：

通过最小化误差平方和，我们可以得到回归系数的估计值。一种常用的方法是使用正规方程（Normal Equation），通过对误差平方和的偏导数等于零来求解回归系数的闭式解。

### 5. 模型评估：

在获得回归系数的估计值后，我们需要评估模型的拟合程度和预测能力。常见的评估指标包括均方误差（Mean Squared Error, MSE）、决定系数（Coefficient of Determination, R-squared）、残差分析等。

### 6. 预测和推断：

通过获得的回归模型和估计的回归系数，可以进行预测和推断。给定新的自变量值，可以利用回归模型预测相应的因变量值，并进行推断和分析。

## 3.8 用户流失预测

### 3.8.1 设计思路

用户流失预测是一项重要的数据分析任务，可以帮助企业识别可能会流失的用户并采取相应的措施来留住他们。基于数据仓库中的用户行为数据、用户个人信息和用户反馈数据，以及服务使用数据，可以利用支持向量机（Support Vector Machine, SVM）和随机森林（Random Forest）算法来完成用户流失预测功能。下面是该功能的设计思路。

#### 1. 数据预处理：

首先，需要从数据仓库中提取用户行为数据、用户个人信息、用户反馈数据和服务使用数据。这些数据可能包括用户的活动记录、访问频率、购买历史、个人信息（如性别、年龄、地理位置等）、用户满意度调查结果等。

在数据预处理阶段，需要进行以下步骤：

- 数据清洗：检查数据中是否存在缺失值、异常值或错误值，并进行相应的处理，例如删除或填补缺失值。
- 特征选择：从所有可用的特征中选择与用户流失相关的重要特征。可以使用统计方法、相关性分析、特征重要性评估等技术来选择特征。
- 特征编码：将非数值型的特征转换为数值型特征，以便于算法的处理。例如，可以使用独热编码（One-Hot Encoding）将分类特征转换为二进制向量表示。

- 数据划分：将数据集划分为训练集和测试集。训练集用于模型的训练和参数调整，测试集用于评估模型的性能。

## 2. 特征工程：

在数据预处理后，进行特征工程是非常重要的步骤。特征工程的目的是提取更有代表性的特征，以提高模型的性能。以下是一些常用的特征工程技术：

- 特征缩放：对数值型特征进行缩放，使其具有相似的尺度。常用的缩放方法包括标准化（Standardization）和归一化（Normalization）。
- 特征构建：根据领域知识或统计分析，构建新的特征来捕捉用户行为和个人信息之间的关联。例如，可以计算用户的活跃度指标、购买频率指标等。
- 特征降维：对于维度较高的数据，可以使用降维技术（如主成分分析）将其转换为较低维度的表示，以减少计算复杂度和避免维度灾难。

## 3. 模型选择与训练：

在用户流失预测中，可以选择支持向量机（SVM）和随机森林（Random Forest）算法进行建模。

- 支持向量机（SVM）：SVM 是一种二分类模型，可以通过在高维特征空间中构造最优超平面来实现分类。在用户流失预测中，SVM 可以根据用户的行为和个人信息特征，学习一个分类模型，将用户划分为流失和非流失两类。

- 随机森林（Random Forest）：随机森林是一种集成学习方法，通过构建多个决策树并进行集成来实现分类。随机森林可以对特征进行自动选择，并且具有较好的鲁棒性和泛化能力。在用户流失预测中，随机森林可以利用用户行为和个人信息特征，构建一个多棵决策树的集成模型，用于判断用户是否会流失。

对于 SVM 和随机森林，可以使用交叉验证等技术选择最优的模型参数，并使用训练集进行模型训练。训练过程中，可以使用常见的性能评估指标（如准确率、召回率、F1 值等）来评估模型的效果，并根据需要进行模型调整和优化。

## 4. 模型评估与应用：

完成模型训练后，可以使用测试集来评估模型的性能。通过对测试集中的样本进行预测，并与真实的用户流失情况进行比较，可以计算出模型的准确率、召回率、F1 值等指标。

模型评估的结果可以为企业提供有关用户流失的洞察，并支持制定相应的用户留存策略。例如，可以通过对流失概率较高的用户进行个性化推荐、优惠券发放、定制化服务等方式来留住用户。

此外，模型也可以应用于新用户的流失预测。当有新的用户行为和个人信息数据进入数据仓库时，可以将其输入模型，预测其是否有流失的趋势，并及时采取相应的措施。

综上所述，实现用户流失预测功能的设计思路包括数据预处理、特征工程、模型选择与训练以

及模型评估与应用。通过这些步骤，可以构建一个有效的用户流失预测模型，为企业提供有针对性的用户留存策略和决策支持。

### 3.8.2 算法思路

#### 1. 支持向量机 (SVM):

支持向量机是一种二分类模型，其主要原理是在特征空间中构建一个最优超平面，以最大化不同类别之间的间隔，从而实现分类。

核心原理:

- 将数据映射到高维特征空间: SVM 的基本思想是通过将低维数据映射到高维特征空间，使得数据在高维空间中更容易分割。这种映射可以通过核函数来实现，常用的核函数包括线性核函数、多项式核函数和径向基函数 (RBF) 核函数。
- 构建最优超平面: 在高维特征空间中，SVM 通过寻找一个最优超平面来将不同类别的数据分开。最优超平面被定义为距离两个类别的最近样本（支持向量）最远的超平面。这个超平面可以被用于将新样本分类到相应的类别。

优点:

- 可处理高维数据: SVM 适用于高维特征空间的数据，可以处理大量特征的情况。
- 泛化能力强: SVM 在处理少量样本的情况下，仍然具有较好的泛化能力，可以有效地解决小样本问题。
- 鲁棒性好: 由于 SVM 的目标是最大化间隔，对于噪声数据和异常值具有较好的鲁棒性。

#### 2. 随机森林 (Random Forest):

随机森林是一种集成学习方法，通过构建多个决策树并进行集成，来实现分类和回归任务。

核心原理:

- 随机抽样: 随机森林采用自助采样法 (Bootstrap Sampling) 从原始数据集中随机抽取一部分样本，形成多个不同的训练子集。同时，对于每个决策树的节点划分，也会从所有特征中随机选择一部分特征进行考虑。
- 决策树的构建: 对于每个训练子集，随机森林构建一棵决策树。决策树通过递归地将训练集分割成更小的子集，使得每个子集中的样本属于同一类别或拥有相似的属性。
- 集成投票: 在预测阶段，随机森林中的所有决策树会对新样本进行预测，并将每棵树的预测结果进行投票或平均，得到最终的预测结果。

优点:

- 高度并行化: 由于每棵决策树的构建和预测是相互独立的，随机森林可以有效地进行并行计算，加快模型的训练和预测速度。
- 可处理高维数据: 随机森林对于高维数据的处理能力较强，可以处理大量特征的情况。
- 鲁棒性好: 随机森林对于噪声数据和缺失值有较好的鲁棒性，能够有效地应对数据质量问题。

## 4. 软件开发模块

### 4.1 开发工具介绍

Superset 是一个由 Airbnb 开源的现代化数据探索和可视化平台。它提供了一个直观易用的界面，用于探索和可视化各种数据源中的数据。Superset 具有以下主要特点：

- (1) 数据连接和集成：Superset 支持多种数据源的连接，包括常见的关系型数据库（如 MySQL、PostgreSQL）和大数据技术（如 Apache Hive、Apache Druid）。可以轻松地配置和管理这些数据源，并在一个统一的界面中访问和查询数据。
- (2) 丰富的可视化选项：Superset 提供了广泛的可视化选项，包括折线图、柱状图、散点图、热图、地图等。可以根据需求选择适当的图表类型，并通过定制化设置调整其外观和交互效果。
- (3) 交互式仪表盘：Superset 支持创建交互式仪表盘，可以将多个可视化组件组合在一起，以便用户可以通过选择过滤器、交互式控件等来探索和分析数据。这使得用户能够根据自己的需求自定义仪表盘，并与数据进行实时交互。
- (4) SQL 查询和数据切片：Superset 提供了一个内置的查询编辑器，使用户能够通过编写 SQL 查询来对数据进行更高级的探索和处理。此外，Superset 还支持数据切片（Slicing），允许用户根据不同的维度和筛选条件对数据进行切片，以便更深入地分析和比较数据。
- (5) 权限和安全性：Superset 具有灵活的权限和安全性设置，可以根据用户角色和组织结构来控制对数据和功能的访问。可以定义不同的角色和权限级别，并限制用户的数据访问和操作范围，以确保数据的安全和保密性。
- (6) 扩展和定制：Superset 是一个开源项目，具有活跃的社区支持和贡献。可以根据自己的需求进行定制和扩展，添加新的数据源连接、可视化插件或功能扩展，以满足特定的数据探索和可视化需求。

### 4.2 开发步骤

- (1) 数据准备：确保数据集合适合进行可视化分析。整理和清洗数据，确保数据质量和一致性。
- (2) 安装和配置 Superset：按照 Superset 的官方文档或指南进行安装和配置。这将涉及安装必要的依赖项、设置数据库连接和用户认证等步骤。
- (3) 创建数据源：通过 Superset 的界面，连接到数据源并配置数据源的参数。这可以是各种数据库、文件或 API 等。
- (4) 创建数据可视化：使用 Superset 的界面，创建交互式报表和可视化。可以选择不同的图

表类型，如折线图、柱状图、地图等，根据您的数据和展示需求进行配置。

(5) 配置和定制界面：根据展示需求，调整和定制 Superset 的界面。可以更改颜色、布局和样式等，以使其与品牌或主题一致。

(6) 导出和嵌入可视化：将数据可视化导出为交互式报表或图表。Superset 提供了多种导出选项，例如将可视化嵌入到网页中，作为 iframe 或图像等的形式。

(7) 创建网页前端：使用网页前端开发技术（如 HTML、CSS 和 JavaScript），创建一个网页来承载和展示 Superset 生成的可视化。您可以在网页中添加其他内容，如标题、说明或其他图表。

(8) 网页布局和样式：使用 CSS 和 JavaScript 等技术，调整网页的布局和样式，使其适应您的展示需求和品牌风格。

(9) 部署和测试：将网页部署到服务器或托管平台上，并进行测试，确保所有可视化和交互功能正常运行。

## 五、测试与总结

### 1. 云计算模块

#### 1.1 HDFS 测试

- a. 使用 **hdfs dfs -mkdir** 命令创建一个新的目录。
- b. 使用 **hdfs dfs -put** 命令将本地文件上传到 HDFS 中。
- c. 使用 **hdfs dfs -ls** 命令检查上传的文件是否存在。
- d. 使用 **hdfs dfs -cat** 命令查看文件内容。

---

**测试结果：**确认目录和文件在 HDFS 中创建成功，并且能够正确读取和显示文件内容。

---

#### 1.2 YARN 测试

- a. 使用 **yarn jar** 命令提交一个简单的 MapReduce 作业。
  - b. 使用 YARN Web 界面或命令行工具（如 **yarn application -list**）检查作业的状态和进展。
  - c. 确保作业成功完成并生成预期的输出结果。
-



**测试结果：**确认作业能够在集群上正确执行，并且输出结果与预期一致。

---

## 1.3 节点管理测试

- e. 使用 Hadoop Web 界面或命令行工具（如 **hadoop dfsadmin -report** 和 **yarn node -list**）查看节点的健康状况和可用性。
- f. 确保所有节点都处于正常状态，并且集群资源分配均衡。

**测试结果：**确认集群中的节点正常工作，并且资源管理和分配功能正常。

---

## 1.4 故障恢复测试

- a. 模拟主节点故障，例如停止主节点上的 NameNode 服务。
- b. 使用命令行或脚本自动化切换为备用节点（如果有）。
- c. 确保备用节点正常接管主节点的功能，并且集群仍然可用。

**测试结果：**确认在主节点故障的情况下，备用节点能够成功切换并保障集群的高可用性。

---

# 2. 大数据模块

## 2.1. MySQL 建库建表

### 1) 各活动补贴率

```
DROP TABLE IF EXISTS `ads_activity_stats`;
CREATE TABLE `ads_activity_stats` (
  `dt` date NOT NULL COMMENT '统计日期',
  `activity_id` varchar(16) CHARACTER SET utf8 COLLATE utf8_general_ci
NOT NULL COMMENT '活动 ID',
  `activity_name` varchar(64) CHARACTER SET utf8 COLLATE utf8_general_ci
NULL DEFAULT NULL COMMENT '活动名称',
  `start_date` varchar(16) CHARACTER SET utf8 COLLATE utf8_general_ci NULL
DEFAULT NULL COMMENT '活动开始日期',
  `reduce_rate` decimal(16, 2) NULL DEFAULT NULL COMMENT '补贴率',
  PRIMARY KEY (`dt`, `activity_id`) USING BTREE
) ENGINE = InnoDB CHARACTER SET = utf8 COLLATE = utf8_general_ci COMMENT
= '活动统计' ROW_FORMAT = Dynamic;
```

### 2) 各优惠券补贴率

```
DROP TABLE IF EXISTS `ads_coupon_stats`;
CREATE TABLE `ads_coupon_stats` (
  `dt` date NOT NULL COMMENT '统计日期',
  `coupon_id` varchar(16) CHARACTER SET utf8 COLLATE utf8_general_ci NOT
NULL COMMENT '优惠券 ID',
  `coupon_name` varchar(64) CHARACTER SET utf8 COLLATE utf8_general_ci
```

```

NULL DEFAULT NULL COMMENT '优惠券名称',
`start_date` varchar(16) CHARACTER SET utf8 COLLATE utf8_general_ci NULL
DEFAULT NULL COMMENT '发布日期',
`rule_name` varchar(64) CHARACTER SET utf8 COLLATE utf8_general_ci NULL
DEFAULT NULL COMMENT '优惠规则, 例如满 100 元减 10 元',
`reduce_rate` decimal(16, 2) NULL DEFAULT NULL COMMENT '补贴率',
PRIMARY KEY(`dt`,`coupon_id`) USING BTREE
) ENGINE = InnoDB CHARACTER SET = utf8 COLLATE = utf8_general_ci COMMENT
= '优惠券统计' ROW_FORMAT = Dynamic;

```

### 3) 新增交易用户统计

```

DROP TABLE IF EXISTS `ads_new_buyer_stats`;
CREATE TABLE `ads_new_buyer_stats` (
`dt` date NOT NULL COMMENT '统计日期',
`recent_days` bigint(20) NOT NULL COMMENT '最近天数,1:最近 1 天,7:最近 7
天,30:最近 30 天',
`new_order_user_count` bigint(20) NULL DEFAULT NULL COMMENT '新增下单人
数',
`new_payment_user_count` bigint(20) NULL DEFAULT NULL COMMENT '新增支付
人数',
PRIMARY KEY(`dt`,`recent_days`) USING BTREE
) ENGINE = InnoDB CHARACTER SET = utf8 COLLATE = utf8_general_ci COMMENT
= '新增交易用户统计' ROW_FORMAT = Dynamic;

```

### 4) 各省份订单统计

```

DROP TABLE IF EXISTS `ads_order_by_province`;
CREATE TABLE `ads_order_by_province` (
`dt` date NOT NULL COMMENT '统计日期',
`recent_days` bigint(20) NOT NULL COMMENT '最近天数,1:最近 1 天,7:最近 7
天,30:最近 30 天',
`province_id` varchar(16) CHARACTER SET utf8 COLLATE utf8_general_ci
NOT NULL COMMENT '省份 ID',
`province_name` varchar(16) CHARACTER SET utf8 COLLATE utf8_general_ci
NULL DEFAULT NULL COMMENT '省份名称',
`area_code` varchar(16) CHARACTER SET utf8 COLLATE utf8_general_ci NULL
DEFAULT NULL COMMENT '地区编码',
`iso_code` varchar(16) CHARACTER SET utf8 COLLATE utf8_general_ci NULL
DEFAULT NULL COMMENT '国际标准地区编码',
`iso_code_3166_2` varchar(16) CHARACTER SET utf8 COLLATE
utf8_general_ci NULL DEFAULT NULL COMMENT '国际标准地区编码',
`order_count` bigint(20) NULL DEFAULT NULL COMMENT '订单数',
`order_total_amount` decimal(16, 2) NULL DEFAULT NULL COMMENT '订单金额',
PRIMARY KEY(`dt`,`recent_days`,`province_id`) USING BTREE
) ENGINE = InnoDB CHARACTER SET = utf8 COLLATE = utf8_general_ci COMMENT
= '各地区订单统计' ROW_FORMAT = Dynamic;

```

### 5) 用户路径分析

```

DROP TABLE IF EXISTS `ads_page_path`;
CREATE TABLE `ads_page_path` (
`dt` date NOT NULL COMMENT '统计日期',
`recent_days` bigint(20) NOT NULL COMMENT '最近天数,1:最近 1 天,7:最近 7
天,30:最近 30 天',
`source` varchar(64) CHARACTER SET utf8 COLLATE utf8_general_ci NOT NULL
COMMENT '跳转起始页面 ID',

```

```

`target` varchar(64) CHARACTER SET utf8 COLLATE utf8_general_ci NOT NULL
COMMENT '跳转终到页面 ID',
`path count` bigint(20) NULL DEFAULT NULL COMMENT '跳转次数',
PRIMARY KEY(`dt`,`recent_days`,`source`,`target`) USING BTREE
) ENGINE = InnoDB CHARACTER SET = utf8 COLLATE = utf8_general_ci COMMENT
= '页面浏览路径分析' ROW_FORMAT = Dynamic;

```

## 6) 各品牌复购率

```

DROP TABLE IF EXISTS `ads_repeat_purchase_by_tm`;
CREATE TABLE `ads_repeat_purchase_by_tm` (
  `dt` date NOT NULL COMMENT '统计日期',
  `recent_days` bigint(20) NOT NULL COMMENT '最近天数,7:最近 7 天,30:最近 30
天',
  `tm_id` varchar(16) CHARACTER SET utf8 COLLATE utf8_general_ci NOT NULL
COMMENT '品牌 ID',
  `tm_name` varchar(32) CHARACTER SET utf8 COLLATE utf8_general_ci NULL
DEFAULT NULL COMMENT '品牌名称',
  `order_repeat_rate` decimal(16, 2) NULL DEFAULT NULL COMMENT '复购率',
PRIMARY KEY(`dt`,`recent_days`,`tm_id`) USING BTREE
) ENGINE = InnoDB CHARACTER SET = utf8 COLLATE = utf8_general_ci COMMENT
= '各品牌复购率统计' ROW_FORMAT = Dynamic;

```

## 7) 各品类商品购物车存量 topN

```

DROP TABLE IF EXISTS `ads_sku_cart_num_top3_by_cate`;
CREATE TABLE `ads_sku_cart_num_top3_by_cate` (
  `dt` date NOT NULL COMMENT '统计日期',
  `category1_id` varchar(16) CHARACTER SET utf8 COLLATE utf8_general_ci
NOT NULL COMMENT '一级分类 ID',
  `category1_name` varchar(64) CHARACTER SET utf8 COLLATE utf8_general_ci
NULL DEFAULT NULL COMMENT '一级分类名称',
  `category2_id` varchar(16) CHARACTER SET utf8 COLLATE utf8_general_ci
NOT NULL COMMENT '二级分类 ID',
  `category2_name` varchar(64) CHARACTER SET utf8 COLLATE utf8_general_ci
NULL DEFAULT NULL COMMENT '二级分类名称',
  `category3_id` varchar(16) CHARACTER SET utf8 COLLATE utf8_general_ci
NOT NULL COMMENT '三级分类 ID',
  `category3_name` varchar(64) CHARACTER SET utf8 COLLATE utf8_general_ci
NULL DEFAULT NULL COMMENT '三级分类名称',
  `sku_id` varchar(16) CHARACTER SET utf8 COLLATE utf8_general_ci NOT NULL
COMMENT '商品 id',
  `sku_name` varchar(128) CHARACTER SET utf8 COLLATE utf8_general_ci NULL
DEFAULT NULL COMMENT '商品名称',
  `cart_num` bigint(20) NULL DEFAULT NULL COMMENT '购物车中商品数量',
  `rk` bigint(20) NULL DEFAULT NULL COMMENT '排名',
PRIMARY KEY(`dt`,`sku_id`,`category1_id`,`category2_id`,`category3_id`) USING BTREE
) ENGINE = InnoDB CHARACTER SET = utf8 COLLATE = utf8_general_ci COMMENT
= '各分类商品购物车存量 Top10' ROW_FORMAT = Dynamic;

```

## 8) 交易综合统计

```

DROP TABLE IF EXISTS `ads_trade_stats`;
CREATE TABLE `ads_trade_stats` (
  `dt` date NOT NULL COMMENT '统计日期',
  `recent_days` bigint(255) NOT NULL COMMENT '最近天数,1:最近 1 日,7:最近 7
天,30:最近 30 天',
  `order_total_amount` decimal(16, 2) NULL DEFAULT NULL COMMENT '订单总

```

```

额,GMV',
`order_count` bigint(20) NULL DEFAULT NULL COMMENT '订单数',
`order_user_count` bigint(20) NULL DEFAULT NULL COMMENT '下单人数',
`order_refund_count` bigint(20) NULL DEFAULT NULL COMMENT '退单数',
`order_refund_user_count` bigint(20) NULL DEFAULT NULL COMMENT '退单人数',
PRIMARY KEY (`dt`,`recent_days`) USING BTREE
) ENGINE = InnoDB CHARACTER SET = utf8 COLLATE = utf8_general_ci COMMENT
= '交易统计' ROW_FORMAT = Dynamic;

```

## 9) 各品类商品交易统计

```

DROP TABLE IF EXISTS `ads_trade_stats_by_cate`;
CREATE TABLE `ads_trade_stats_by_cate` (
  `dt` date NOT NULL COMMENT '统计日期',
  `recent_days` bigint(20) NOT NULL COMMENT '最近天数,1:最近 1 天,7:最近 7
天,30:最近 30 天',
  `category1_id` varchar(16) CHARACTER SET utf8 COLLATE utf8_general_ci
NOT NULL COMMENT '一级分类 id',
  `category1_name` varchar(64) CHARACTER SET utf8 COLLATE utf8_general_ci
NULL DEFAULT NULL COMMENT '一级分类名称',
  `category2_id` varchar(16) CHARACTER SET utf8 COLLATE utf8_general_ci
NOT NULL COMMENT '二级分类 id',
  `category2_name` varchar(64) CHARACTER SET utf8 COLLATE utf8_general_ci
NULL DEFAULT NULL COMMENT '二级分类名称',
  `category3_id` varchar(16) CHARACTER SET utf8 COLLATE utf8_general_ci
NOT NULL COMMENT '三级分类 id',
  `category3_name` varchar(64) CHARACTER SET utf8 COLLATE utf8_general_ci
NULL DEFAULT NULL COMMENT '三级分类名称',
  `order_count` bigint(20) NULL DEFAULT NULL COMMENT '订单数',
  `order_user_count` bigint(20) NULL DEFAULT NULL COMMENT '订单人数',
  `order_refund_count` bigint(20) NULL DEFAULT NULL COMMENT '退单数',
  `order_refund_user_count` bigint(20) NULL DEFAULT NULL COMMENT '退单人数',
PRIMARY KEY (`dt`,`recent_days`,`category1_id`,`category2_id`,`category3_id`) USING BTREE
) ENGINE = InnoDB CHARACTER SET = utf8 COLLATE = utf8_general_ci COMMENT
= '各分类商品交易统计' ROW_FORMAT = Dynamic;

```

## 10) 各品牌商品交易统计

```

DROP TABLE IF EXISTS `ads_trade_stats_by_tm`;
CREATE TABLE `ads_trade_stats_by_tm` (
  `dt` date NOT NULL COMMENT '统计日期',
  `recent_days` bigint(20) NOT NULL COMMENT '最近天数,1:最近 1 天,7:最近 7
天,30:最近 30 天',
  `tm_id` varchar(16) CHARACTER SET utf8 COLLATE utf8_general_ci NOT NULL
COMMENT '品牌 ID',
  `tm_name` varchar(32) CHARACTER SET utf8 COLLATE utf8_general_ci NULL
DEFAULT NULL COMMENT '品牌名称',
  `order_count` bigint(20) NULL DEFAULT NULL COMMENT '订单数',
  `order_user_count` bigint(20) NULL DEFAULT NULL COMMENT '订单人数',
  `order_refund_count` bigint(20) NULL DEFAULT NULL COMMENT '退单数',
  `order_refund_user_count` bigint(20) NULL DEFAULT NULL COMMENT '退单人数',
PRIMARY KEY (`dt`,`recent_days`,`tm_id`) USING BTREE
) ENGINE = InnoDB CHARACTER SET = utf8 COLLATE = utf8_general_ci COMMENT
= '各品牌商品交易统计' ROW_FORMAT = Dynamic;

```

## 11) 各渠道流量统计

```
DROP TABLE IF EXISTS `ads_traffic_stats_by_channel`;
CREATE TABLE `ads_traffic_stats_by_channel` (
  `dt` date NOT NULL COMMENT '统计日期',
  `recent_days` bigint(20) NOT NULL COMMENT '最近天数,1:最近 1 天,7:最近 7 天,30:最近 30 天',
  `channel` varchar(16) CHARACTER SET utf8 COLLATE utf8_general_ci NOT NULL COMMENT '渠道',
  `uv_count` bigint(20) NULL DEFAULT NULL COMMENT '访客人数',
  `avg_duration_sec` bigint(20) NULL DEFAULT NULL COMMENT '会话平均停留时长, 单位为秒',
  `avg_page_count` bigint(20) NULL DEFAULT NULL COMMENT '会话平均浏览页面数',
  `sv_count` bigint(20) NULL DEFAULT NULL COMMENT '会话数',
  `bounce_rate` decimal(16, 2) NULL DEFAULT NULL COMMENT '跳出率',
  PRIMARY KEY (`dt`,`recent_days`,`channel`) USING BTREE
) ENGINE = InnoDB CHARACTER SET = utf8 COLLATE = utf8_general_ci COMMENT = '各渠道流量统计' ROW_FORMAT = Dynamic;
```

## 12) 用户行为漏斗分析

```
DROP TABLE IF EXISTS `ads_user_action`;
CREATE TABLE `ads_user_action` (
  `dt` date NOT NULL COMMENT '统计日期',
  `recent_days` bigint(20) NOT NULL COMMENT '最近天数,1:最近 1 天,7:最近 7 天,30:最近 30 天',
  `home_count` bigint(20) NULL DEFAULT NULL COMMENT '浏览首页人数',
  `good_detail_count` bigint(20) NULL DEFAULT NULL COMMENT '浏览商品详情页人数',
  `cart_count` bigint(20) NULL DEFAULT NULL COMMENT '加入购物车人数',
  `order_count` bigint(20) NULL DEFAULT NULL COMMENT '下单人数',
  `payment_count` bigint(20) NULL DEFAULT NULL COMMENT '支付人数',
  PRIMARY KEY (`dt`,`recent_days`) USING BTREE
) ENGINE = InnoDB CHARACTER SET = utf8 COLLATE = utf8_general_ci COMMENT = '漏斗分析' ROW_FORMAT = Dynamic;
```

## 13) 用户变动统计

```
DROP TABLE IF EXISTS `ads_user_change`;
CREATE TABLE `ads_user_change` (
  `dt` varchar(16) CHARACTER SET utf8 COLLATE utf8_general_ci NOT NULL COMMENT '统计日期',
  `user_churn_count` varchar(16) CHARACTER SET utf8 COLLATE utf8_general_ci NULL DEFAULT NULL COMMENT '流失用户数',
  `user_back_count` varchar(16) CHARACTER SET utf8 COLLATE utf8_general_ci NULL DEFAULT NULL COMMENT '回流用户数',
  PRIMARY KEY (`dt`) USING BTREE
) ENGINE = InnoDB CHARACTER SET = utf8 COLLATE = utf8_general_ci COMMENT = '用户变动统计' ROW_FORMAT = Dynamic;
```

## 14) 用户留存率

```
DROP TABLE IF EXISTS `ads_user_retention`;
CREATE TABLE `ads_user_retention` (
  `dt` date NOT NULL COMMENT '统计日期',
  `create_date` varchar(16) CHARACTER SET utf8 COLLATE utf8_general_ci NOT NULL COMMENT '用户新增日期',
```

```

`retention_day` int(20) NOT NULL COMMENT '截至当前日期留存天数',
`retention_count` bigint(20) NULL DEFAULT NULL COMMENT '留存用户数量',
`new_user_count` bigint(20) NULL DEFAULT NULL COMMENT '新增用户数量',
`retention_rate` decimal(16, 2) NULL DEFAULT NULL COMMENT '留存率',
PRIMARY KEY(`dt`,`create_date`,`retention_day`) USING BTREE
) ENGINE = InnoDB CHARACTER SET = utf8 COLLATE = utf8_general_ci COMMENT
= '留存率' ROW_FORMAT = Dynamic;

```

## 15) 用户新增活跃统计

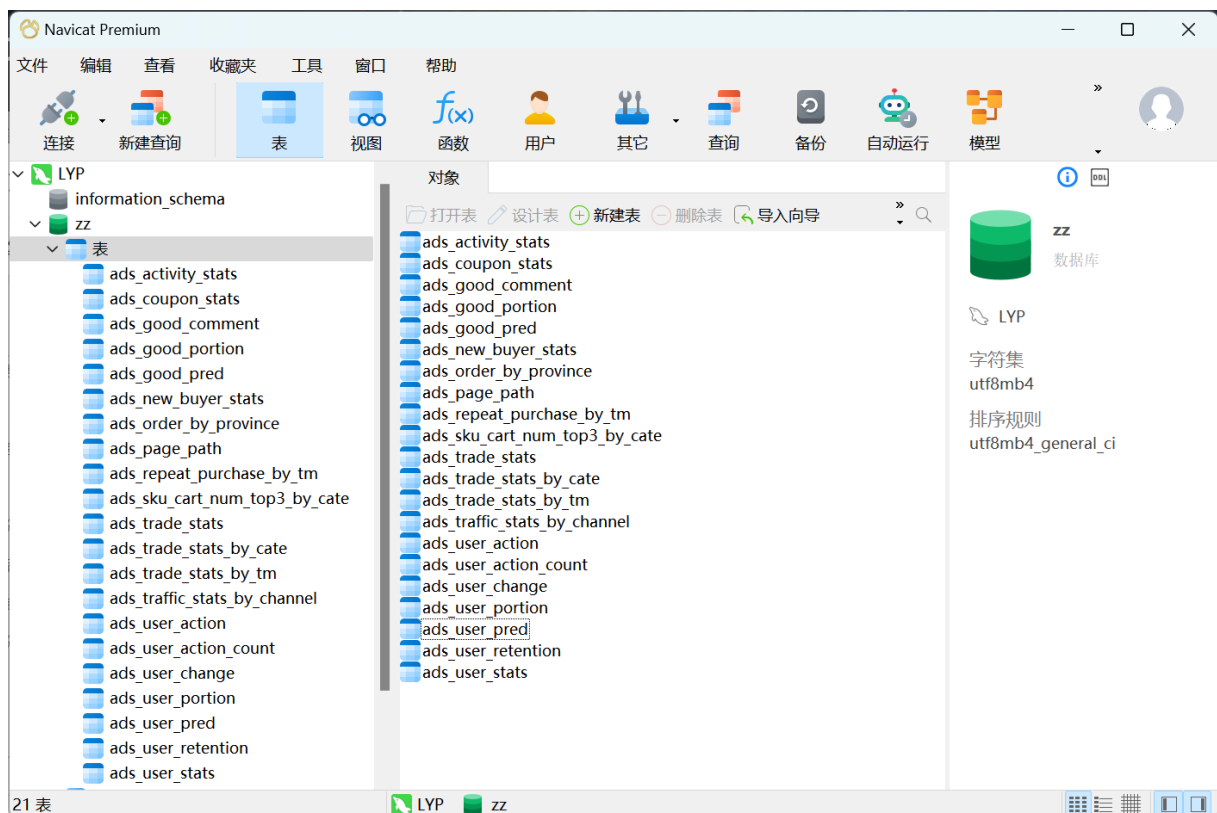
```

DROP TABLE IF EXISTS `ads_user_stats`;
CREATE TABLE `ads_user_stats` (
  `dt` date NOT NULL COMMENT '统计日期',
  `recent_days` bigint(20) NOT NULL COMMENT '最近 n 日,1:最近 1 日,7:最近 7
日,30:最近 30 日',
  `new_user_count` bigint(20) NULL DEFAULT NULL COMMENT '新增用户数',
  `active user count` bigint(20) NULL DEFAULT NULL COMMENT '活跃用户数',
  PRIMARY KEY(`dt`,`recent_days`) USING BTREE
) ENGINE = InnoDB CHARACTER SET = utf8 COLLATE = utf8_general_ci COMMENT
= '用户新增活跃统计' ROW_FORMAT = Dynamic;

```

## 2.2. 访问数据表

为了方便前后端访问数据，我们将生成的模拟数据移到云端：



数据表中是我们生成的模拟数据，以 ads\_activity\_stats 为例：

ads\_activity\_stats @zz (LYP) - 表 - Navicat Premium

文件 编辑 查看 表 收藏夹 工具 窗口 帮助

连接 新建查询 表 视图 函数 用户 其它 查询 备份 自动运行 模型

LYP

information\_schema

zz

表

ads\_activity\_stats

ads\_coupon\_stats

ads\_good\_comme

ads\_good\_portion

ads\_good\_pred

ads\_new\_buyer\_st

ads\_order\_by\_pro

ads\_page\_path

ads\_repeat\_purch

ads\_sku\_cart\_num

ads\_trade\_stats

ads\_trade\_stats\_b

ads\_trade\_stats\_b

ads\_traffic\_stats

ads\_user\_action

ads\_user\_action\_c

ads\_user\_change

ads\_user\_portion

ads\_user\_pred

ads\_user\_retention

对象 ads\_activity\_stats @zz (LYP) - 表

开始事务 文本 筛选 排序 导入 导出 数据生成

dt	activity_id	activity_name	start_date	reduce_rate
2020-06-14	1	联想专场	2020-10-21	0.04
2020-06-14	2	Apple品牌日	2020-06-10	0.03
2020-01-01	1	华为花粉日	2019-12-29	0.03
2020-01-01	2	Apple品牌日	2019-12-29	0.18
2020-01-01	3	Apple品牌日	2019-12-29	0.12
2020-01-02	4	联想专场	2019-12-30	0.10
2020-01-02	5	Apple品牌日	2019-12-30	0.12
2020-01-02	6	联想专场	2019-12-30	0.08
2020-01-03	7	华为花粉日	2019-12-31	0.18
2020-01-04	8	联想专场	2020-01-01	0.10
2020-01-04	9	联想专场	2020-01-01	0.15
2020-01-04	10	联想专场	2020-01-01	0.17
2020-01-05	11	Apple品牌日	2020-01-02	0.13
2020-01-05	12	联想专场	2020-01-02	0.12
2020-01-05	13	华为花粉日	2020-01-02	0.06
2020-01-05	14	联想专场	2020-01-02	0.03
2020-01-06	15	SUMSUNG会场	2020-01-03	0.18
2020-01-06	16	华为花粉日	2020-01-03	0.12
2020-01-06	17	华为花粉日	2020-01-03	0.18

dt

类型 date

不是 null 否

默认值 (NULL)

注释 --

SELECT \* FROM `zz`.`ads\_activity\_stats` LIMIT 0,1000

第 1 条记录 (共 246 条) 于第 1 页

3. 人工智能模块

3.1 用户画像分析

用户画像分析	
输入	用户省份 id 用户性别 用户行为 7、30 天内买得最多的品类 7、30 天内买得第二多的品类 7、30 天内买得第三多的品类 7、30 天内买得最多的品牌 7、30 天内买得第二多的品牌 7、30 天内买得第三多的品牌 7、30 天内的购买单数 7、30 天内的消费总金额
使用的算法	KNN 聚类算法
输出	新用户占比 活跃用户占比 沉默用户占比 流动用户占比

3.2 售后服务满意度分析

售后服务满意度分析	
输入	用户对售后服务的平均评分 用户对售后服务提出的投诉数量 根据用户的评价和反馈中的关键词 退单件数 退单金额 评价时间 评价名称
使用的算法	文本情感分类、情感强度分析
输出	好评率 差评率 星级

3.3 产品销售趋势分析

产品销售趋势分析	
输入	历史销售数据, 前 100 天
使用的算法	时间序列分析算法 (ARIMA、指数平滑)
输出	未来预测数据, 后 100 天

3.4 商品热度分析

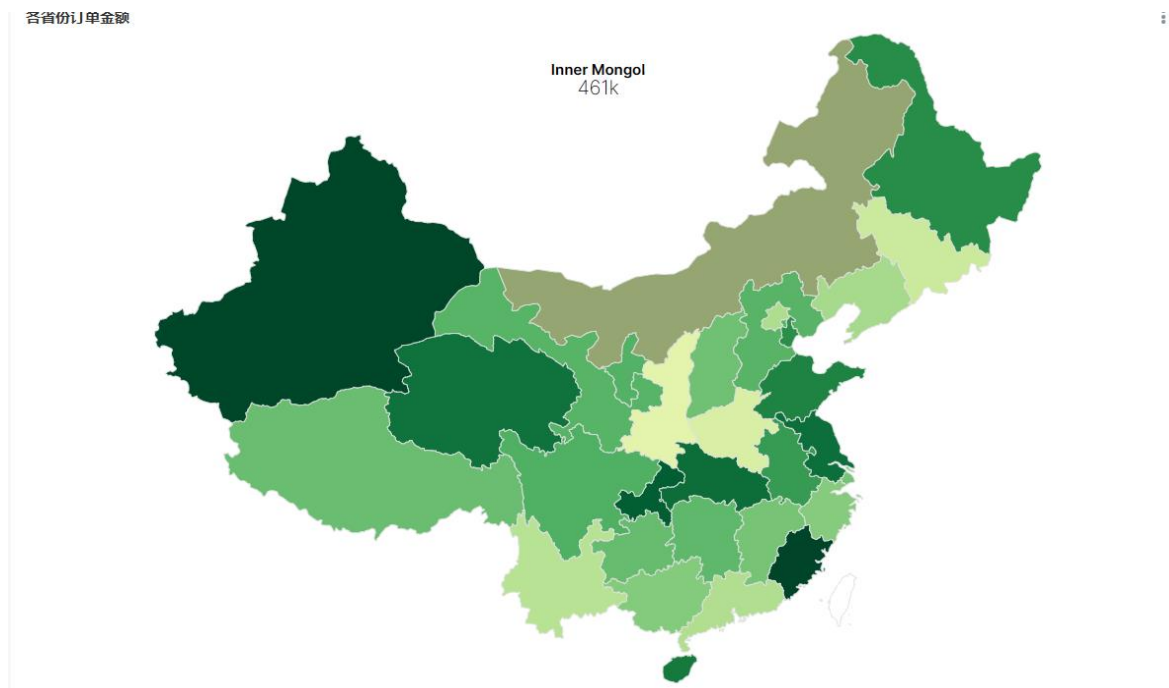
商品热度分析	
输入	销售增长率 商品页面访问量 最近 7、30 日复购率 最近 1、7、30 日订单数 最近 1、7、30 日订单人数 最近 1、7、30 日收藏数 最近 1、7、30 日收藏人数 最近 1、7、30 日退单数 最近 1、7、30 日退单人数 最近 1、7、30 日订单总额
使用的算法	KNN 聚类算法
输出	热销商品占比 潜力商品占比 稳定商品占比 长尾商品占比



## 4.软件开发模块

### 4.1 交易主题

#### 4.1.1 各省份订单金额地图



在前端展示时，鼠标移动到指定位置则会显示具体金额和省份

#### 4.1.2 各省份订单数量

##### 各省份订单数

Show All entries

Search 34 records...

province_name	SUM(order_count)	%SUM(order_count)
海南	57	4.368%
湖北	52	3.985%
山东	50	3.831%
宁夏	49	3.755%
青海	48	3.678%
新疆	47	3.602%
贵州	44	3.372%
河北	42	3.218%
浙江	42	3.218%
黑龙江	42	3.218%
广西	41	3.142%
广东	41	3.142%
江苏	39	2.989%
澳门	39	2.989%
天津	38	2.912%
上海	37	2.835%
吉林	36	2.759%
内蒙古	36	2.759%
辽宁	36	2.759%
安徽	36	2.759%
山西	35	2.682%
甘肃	31	2.375%

前端展示可以提供搜索不同省份订单数功能

4.2 优惠券主题

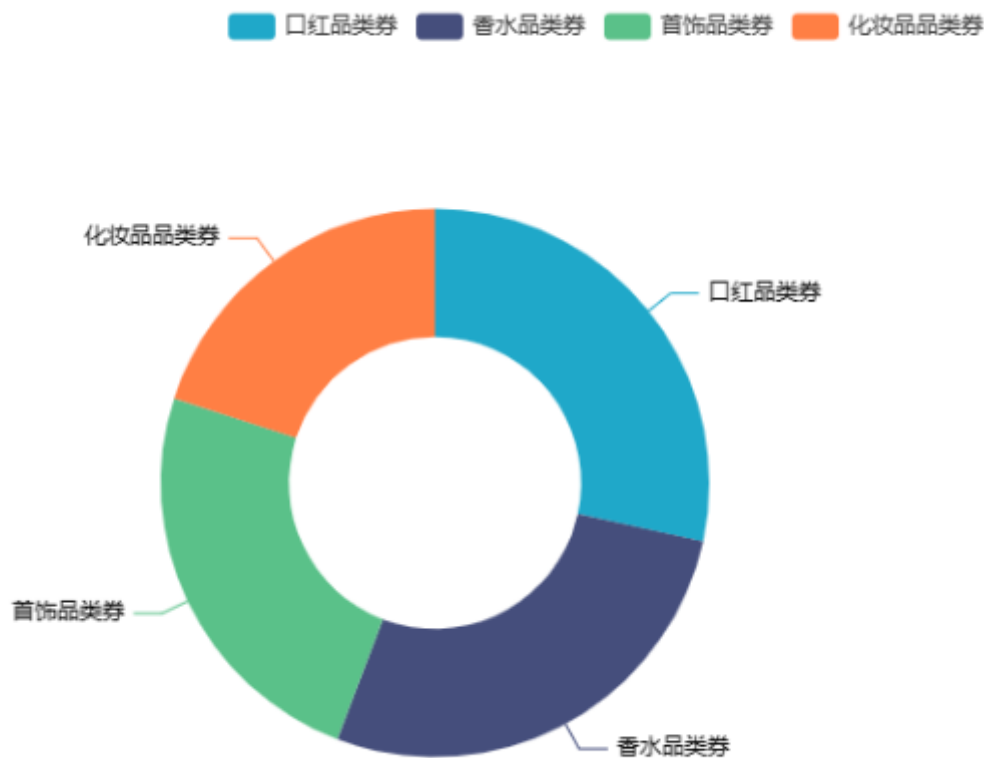


显示当日优惠券的补贴率，并提供查询历史数据，以及与前一天的比较



显示当日活动的补贴率，并提供查询历史数据，以及与前一天的比较

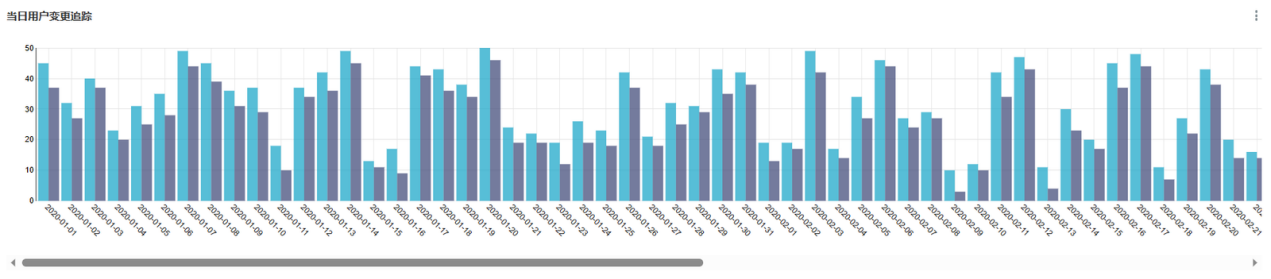
优惠券使用次数



显示优惠券的使用次数，并提供查询各类券的比较

4.3 用户主题

4.3.1 用户变动统计



当日用户变更追踪，可以提供历史数据查询

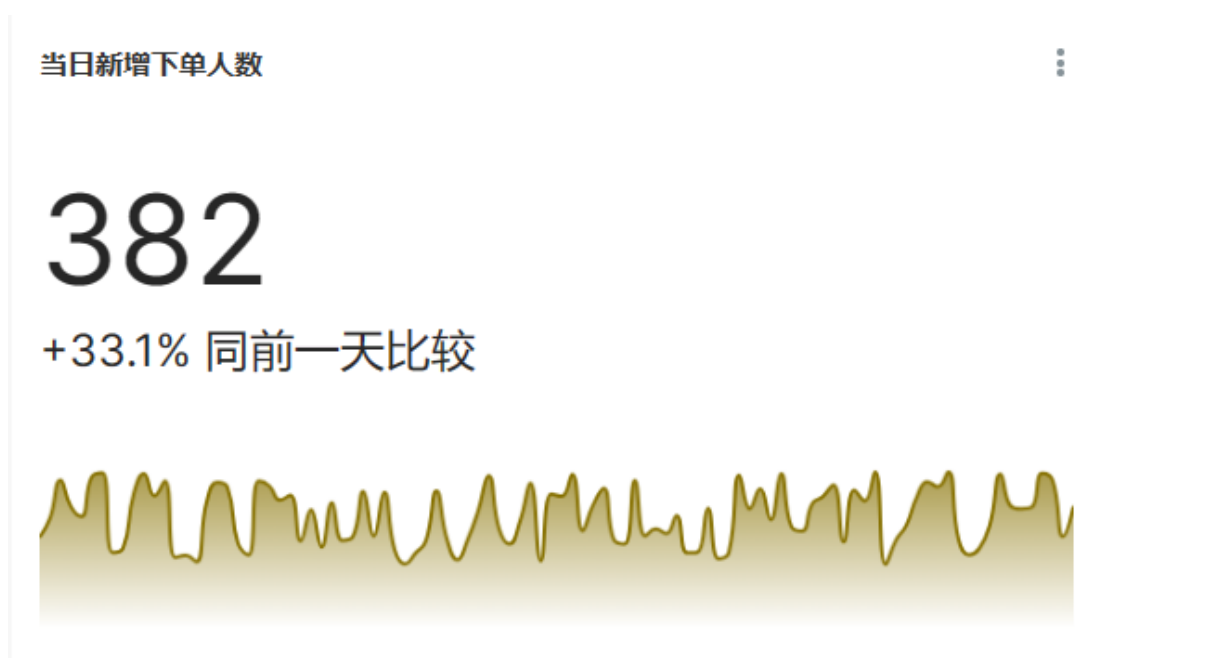


显示每日用户留存率，同时可以提供历史数据查询



记录当日新增人数，同时可以提供历史数据查询，以及与前一天的数据比较

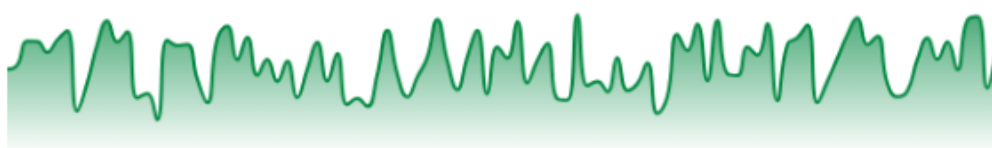
#### 4.3.2 用户行为漏斗分析



## 当日新增支付成功人数

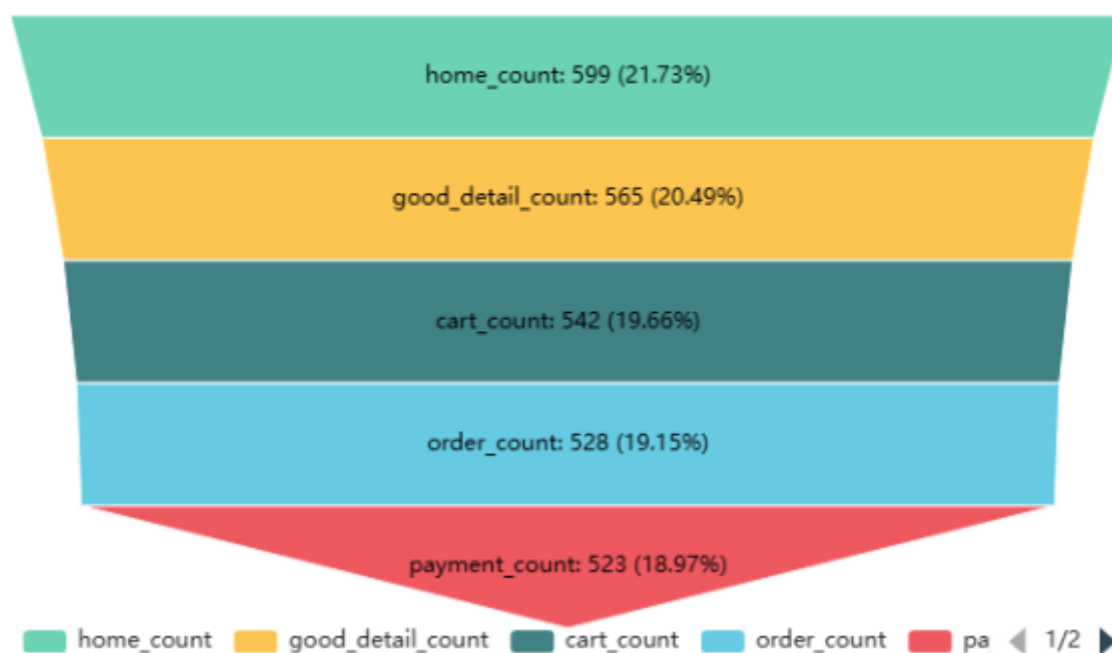
377

+71.4% 同前一天比较



显示当日新增下单人数和新增支付成功人数，并提供历史数据查询，以及与前一天的比较

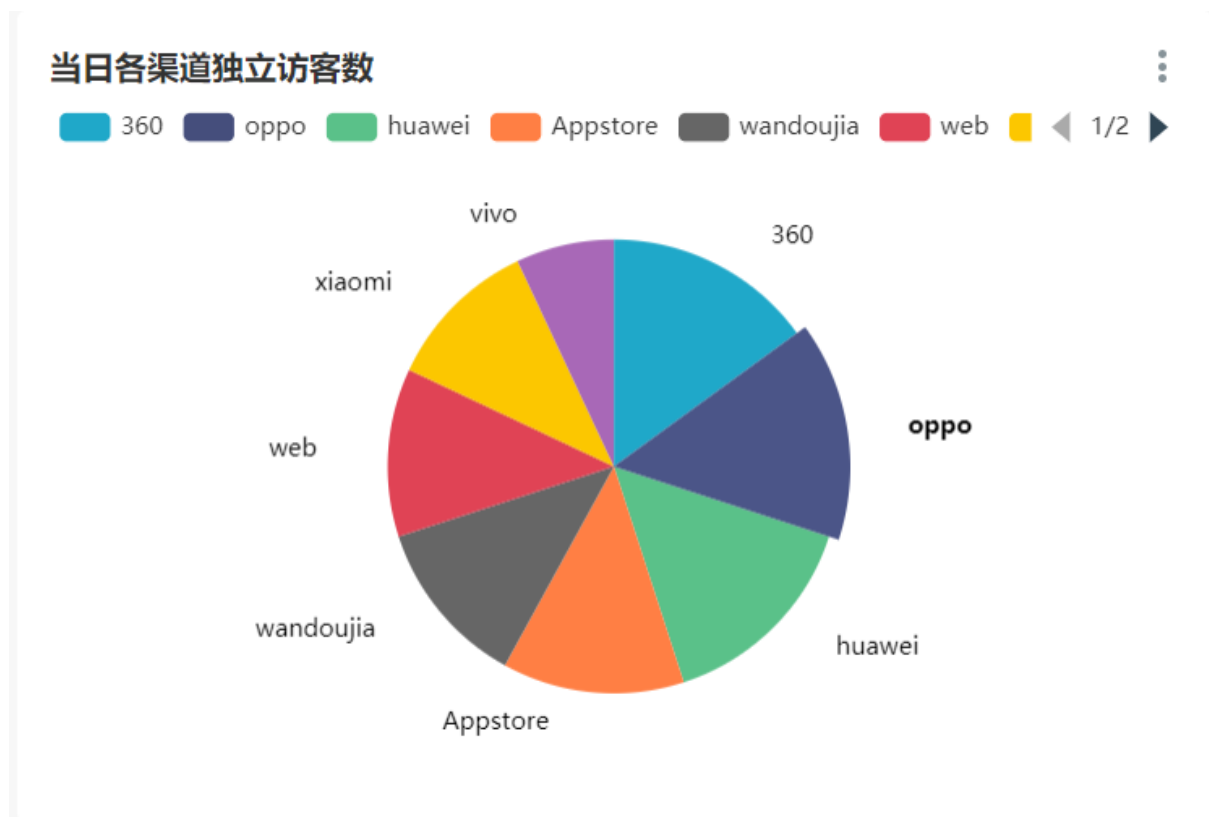
## 用户行为漏斗分析



显示主页浏览人数，浏览商品详情人数，加购人数，下订单人数和支付成功人数，最终数据呈现漏斗分布。

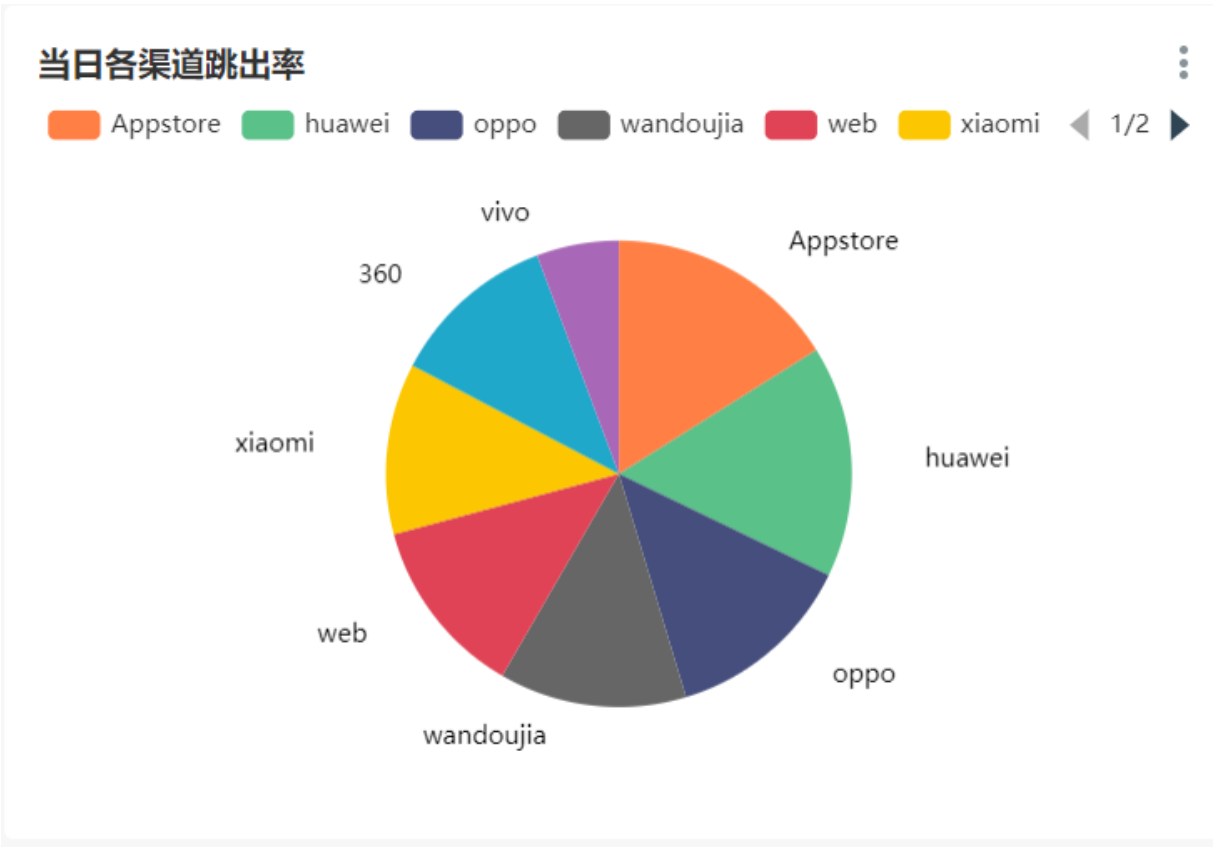
## 4.4 流量主题

### 4.4.1 当日各渠道独立访客数



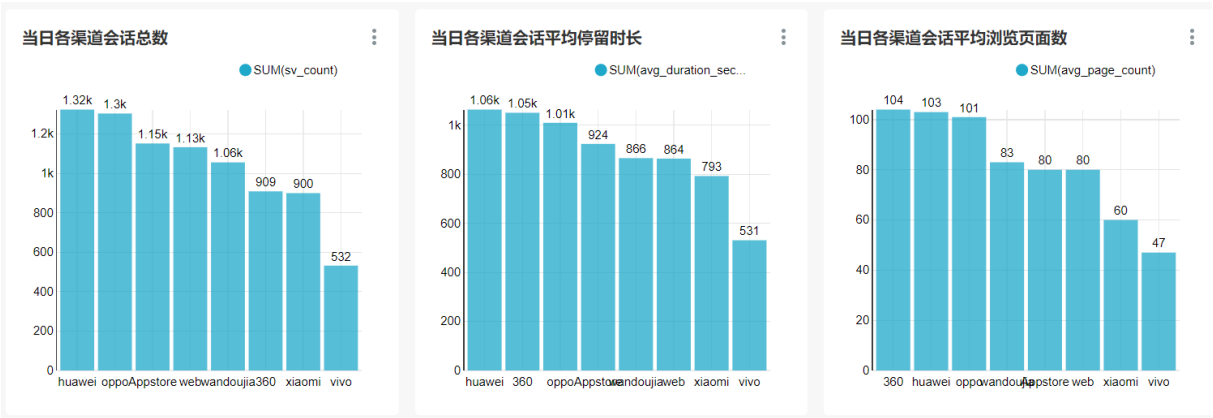
显示了当天通过各个渠道访问的可客户数量，最终数据使用饼状图呈现。

4.4.2 当日各渠道跳出率



显示了当天通过各个渠道的跳出率，最终数据使用饼状图呈现。

4.4.3 当日各渠道会话信息

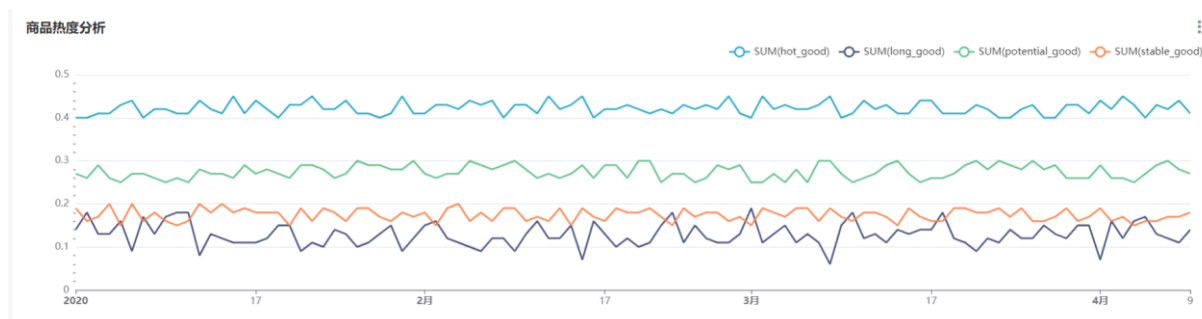


显示了当天各渠道的会话总数、会话平均停留时长以及会话平均浏览页面数。



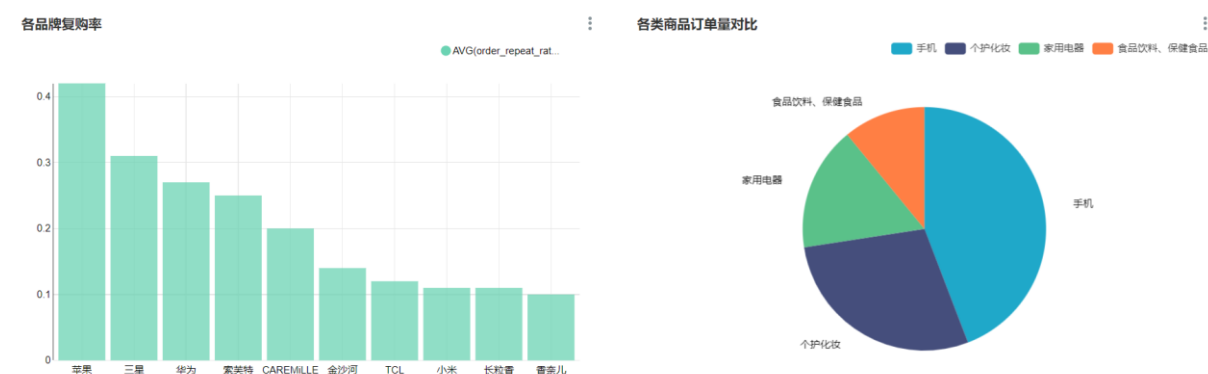
## 4.5 商品主题

### 4.5.1 商品热度分析



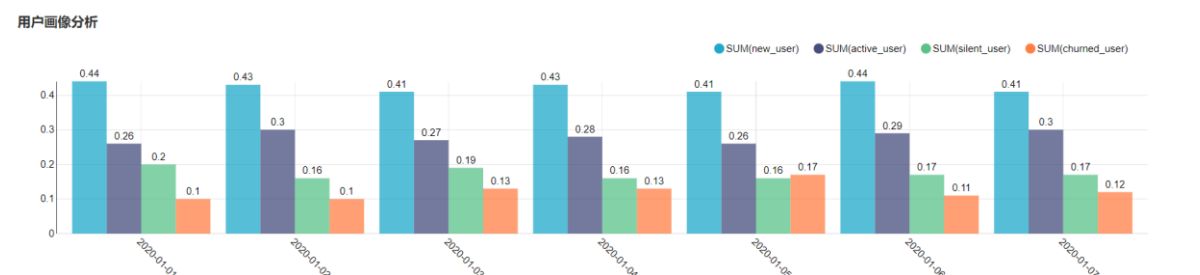
显示了每天的热销商品、潜力商品、稳定商品和长尾商品的分别占比，使用折线图呈现。

### 4.5.2 品牌商品分析



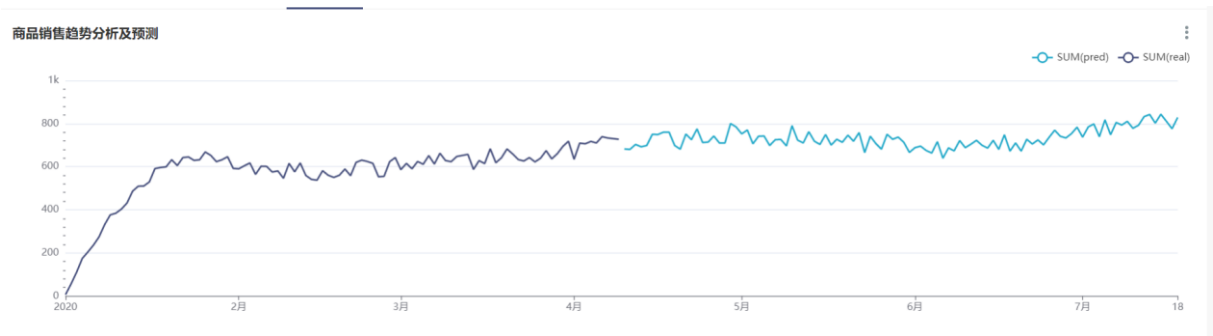
显示了各品牌的复购率和各类商品订单对比，分别用柱状图和饼状图呈现。

### 4.5.2 用户画像分析

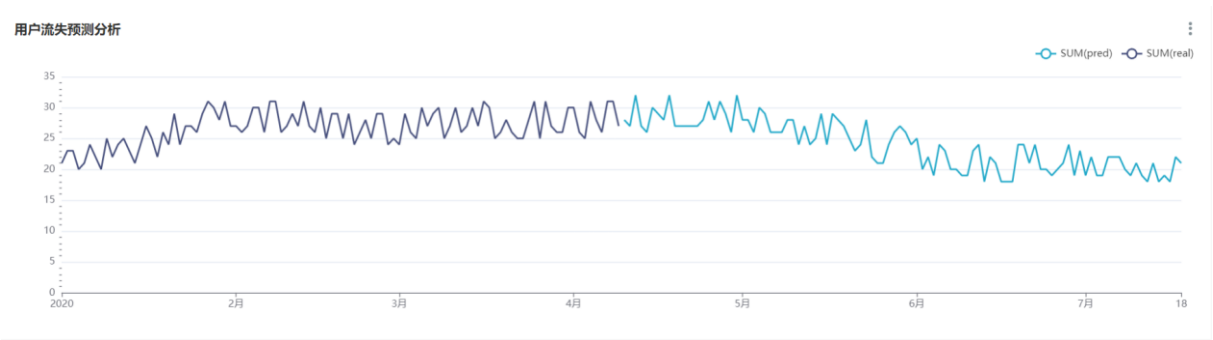


显示了一周时间内新用户、活跃用户、沉默用户、流动用户的占比，采用柱状图呈现。

4.5.3 商品趋势分析



人工智能模块通过分析前 100 天的数据来预测后 100 天的商品销售趋势，结果通过折线图展示。



人工智能模块通过分析前 100 天的数据来预测后 100 天的用户流失趋势，结果通过折线图展示。

六、项目总结

在中国电子的指导下，我们成功地完成了一个综合实战项目，将云计算、大数据、人工智能和软件开发等领域相结合，同时还开发了一个用户精细化运营系统。以下是我们小组对项目的总结：

- 1.项目背景和目标：该项目的背景是满足不断增长的数据处理需求，提升运营效率和精细化运营能力。我们的目标是建立一个集云计算、大数据和人工智能于一体的综合解决方案，同时开发一个用户精细化运营系统，以提高客户满意度和业务成果。
- 2.团队协作和沟通：项目中，团队合作和有效的沟通是成功的关键。我们建立了跨部门的合作机制，明确分工和责任，并定期召开会议以共享进展、解决问题和调整项目方向。通过密切的合作和开放的沟通，我们能够高效地推进项目，并充分发挥每个团队成员的优势。
- 3.技术选型和方案设计：在项目中，我们进行了全面的技术评估和选型。针对云计算、大

数据和人工智能等领域的需求，我们选择了合适的技术框架和工具，包括云平台、数据存储和处理工具、机器学习算法等。我们的方案设计考虑了系统的可扩展性、性能和安全性，并保证了技术组件之间良好集成。

4.数据处理和分析：在项目中，我们处理了大量的数据，包括结构化和非结构化数据。我们建立了高效的数据采集、存储和处理机制，利用大数据技术进行数据清洗、转换和分析。通过应用人工智能算法和模型，我们能够从数据中挖掘有价值的信息，并为业务决策提供支持。

5.用户精细化运营系统：我们开发了一个用户精细化运营系统，旨在通过个性化推荐、行为分析和智能营销等手段，提升用户体验和增加业务价值。该系统集成了大数据和人工智能技术，能够实时监测用户行为、分析用户需求，并通过精确的营销策略与用户互动。这为企业提供了更精细化、个性化的服务，增强了客户忠诚度和竞争力。

6.测试和上线：在项目开发过程中，我们进行了全面的测试，包括单元测试、集成测试和系统测试等。我们保证了软件的质量和稳定性，并及时修复和优化存在的问题。最终，我们成功地将项目上线，并进行了充分的用户培训和支持，以确保用户能够顺利使用和享受系统带来的好处。

7.成果和经验教训：通过这个综合实战项目，我们获得了很多宝贵的经验和成果。我们提高了对云计算、大数据和人工智能等领域的理解和应用能力。同时，我们也面临了一些挑战和教训，例如需求变更、技术难题和项目管理等方面。我们总结了这些经验教训，并将其作为今后类似项目的参考和借鉴。

总之，通过该综合实战项目及用户精细化运营系统的开发，我们不仅成功地将云计算、大数据和人工智能等技术应用于实际业务中，还提升了运营效率和用户体验。这个项目为企业带来了实际的业务价值和竞争优势，也为我们团队的技术成长和发展提供了宝贵的机会。