



超市客户聚类分析数学建模

通过科学的数学建模方法,将500名会员客户进行精准分群,为超市制定个性化营销策略提供数据支撑。

问题背景与研究目标

某超市为提升营销效果和客户满意度,收集了500名会员客户的消费行为数据。管理层希望通过科学方法将客户分群,针对不同群体制定个性化的营销策略。

研究目标包括:精准识别客户群体特征、优化券包发放策略、制定差异化推荐品类、确定合理触达频率。

500

样本客户

4

关键指标

客户数据关键指标



年龄

客户年龄范围18-70岁,反映不同年龄段的消费特征和购买习惯。



月收入

客户月收入水平2,272-17,484元,体现消费能力和购买力差异。



月消费金额

客户在超市的月消费金额439-3,026元,显示消费投入程度。



月购买频次

客户每月到超市购买次数3-20次,反映购物习惯和忠诚度。

K-means聚类方法选择

方法核心思想

K-means聚类将样本划分到"离自己最近"的聚类中心点,实现类内相似度最大化、类间差异最大化的目标。

方法优势

- 每类有明确的均值中心,便于构建客户画像
- 算法高效稳定,适合中等样本量分析
- 结果易于商业解释和应用



聚类数K的科学选择

()

肘部法分析

随着K增大,总平方误差(SSE)下降速度变缓,出现"拐点"处为最优K值。

([

业务导向考量

结合营销资源和管理复杂度,一般选择K=3-5个群体最为合理。

({

稳定性验证

多次运行算法,确保聚类结果的稳定性和可重复性。

数学模型构建

建立K-means聚类的数学模型,定义数据矩阵、特征向量和优化目标。

数据矩阵

$X \in \mathbb{R}^{n \times p}$,其中 $n=500$ (样本量), $p=4$ (特征数)

特征向量

x_1 :年龄、 x_2 :月收入、 x_3 :月消费金额、 x_4 :月购买频次

目标函数

最小化类内平方和(WCSS):

$$\min_{C_k, \mu_k} \sum_{k=1}^K \sum_{x_i \in C_k} \|x_i - \mu_k\|_2^2$$

其中 C_k 为第 k 类的样本集合, μ_k 为第 k 类的中心点。

数据预处理关键步骤



缺失值处理

检查数据完整性,删除少量缺失样本或用中位数填补,确保数据质量。



特征标准化

采用Z-score标准化消除量纲差异,使不同特征在同一尺度上比较。



数据验证

验证标准化后数据的分布特征,确保满足聚类分析的前提条件。

Z-score标准化公式

$$Z_{ij} = \frac{X_{ij} - \bar{X}_j}{S_j}$$

其中 \bar{X}_j 为第j个特征的均值, S_j 为标准差。标准化消除了年龄与月收入等不同量级特征的影响。

K-means算法迭代流程



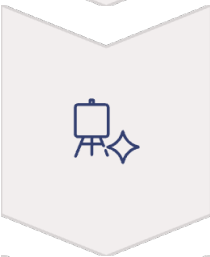
步骤1:初始化

设定聚类数K,随机选择K个初始中心点 $\mu_1^{(0)}, \mu_2^{(0)}, \dots, \mu_K^{(0)}$



步骤2:分配样本(E步)

将每个样本分配到距离最近的聚类中心,形成K个类别



步骤3:更新中心(M步)

重新计算每个类别的中心点,作为新的聚类中心



步骤4:收敛判断

当中心点变化小于阈值 ϵ 或达到最大迭代次数时停止

模型优化策略

K值选择优化

测试K值候选集 $K \in \{2, 3, 4, 5, 6\}$, 通过肘部法和轮廓系数选择最优K值。

稳定性增强

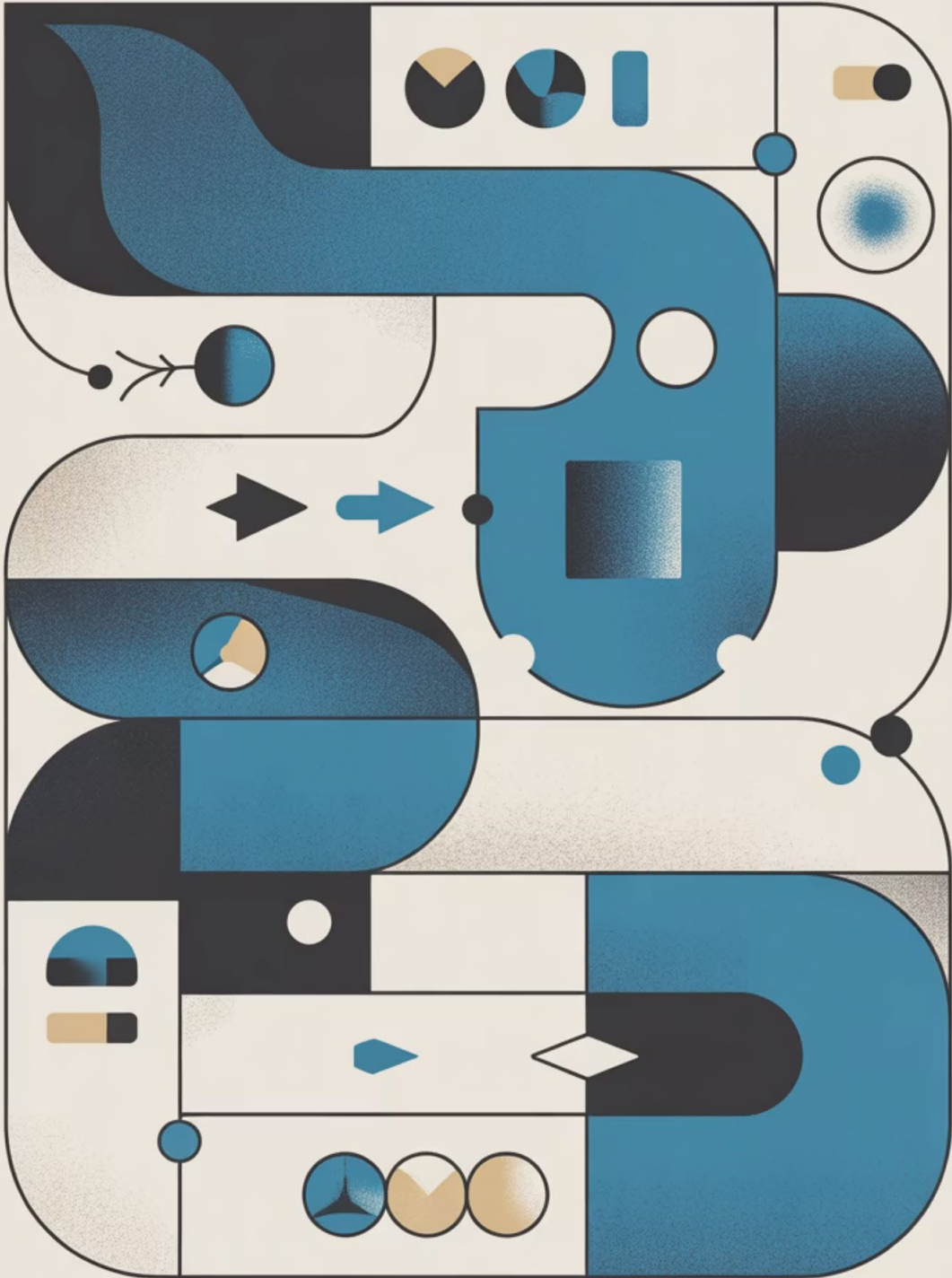
设置多次随机初始化($n_{init}=20-50$), 选择SSE最小的结果作为最终聚类, 避免局部最优解。

特征选择验证

测试不同特征组合 $p \in \{2, 3, 4\}$, 评估各特征对聚类效果的贡献度。

参数调优

调整收敛阈值和最大迭代次数, 平衡计算效率与聚类精度。



聚类质量评估指标

类内平方和(SSE)

衡量类内紧密度,SSE越小表示聚类效果越好。

$$SSE = \sum_{k=1}^K \sum_{x_i \in C_k} |x_i - \mu_k|^2$$

轮廓系数(Silhouette Score)

综合评估类内紧密度和类间分离度,取值范围[-1,1],越接近1越好。

$$s_i = \frac{b_i - a_i}{\max(a_i, b_i)}$$

Davies-Bouldin指数

评估聚类的分离度和紧凑度,DB指数越小表示聚类质量越高。

$$DB = \frac{1}{K} \sum_{k=1}^K \max_{j \neq k} \frac{\sigma_k + \sigma_j}{d_{kj}}$$

业务有效性验证



群体区分度

通过方差分析(ANOVA)检验各群体在关键指标上的差异显著性,确保群体间存在明显区别。



群体稳定性

评估类内标准差,确保群体特征明确,成员相似度高,聚类结果稳定可靠。



商业可解释性

验证聚类结果是否符合业务直觉和营销需求,能否转化为实际营销策略。

SPSS求解实施步骤

1

第一阶段

数据预处理

分析→描述统计→描述,生成标准化变量

2

第二阶段

K-means聚类

分析→分类→K均值聚类,设定K=3

3

第三阶段

结果统计分析

分析→比较均值→均值,生成群体画像

4

第四阶段

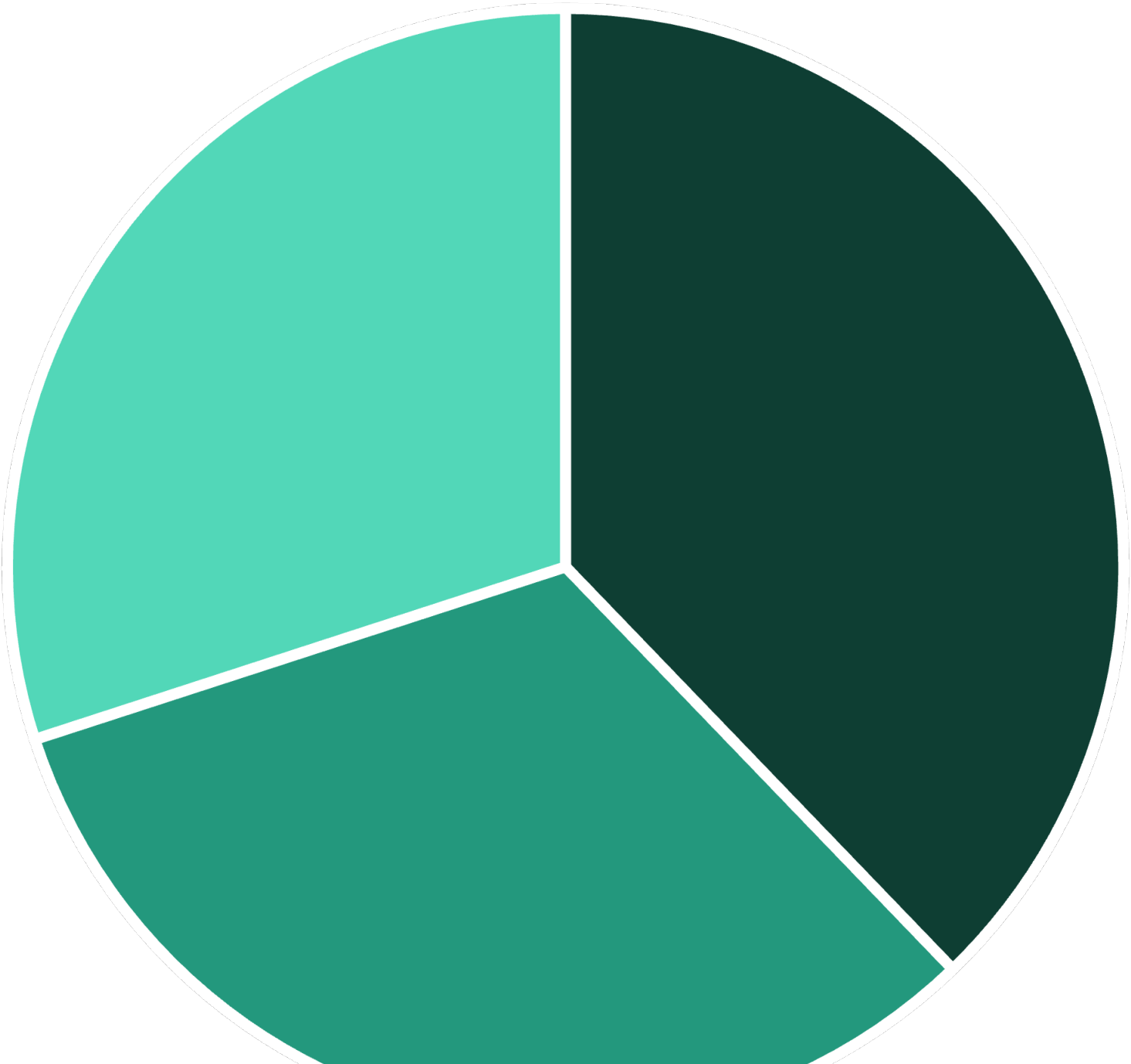
K值优化验证

测试K=2,4,5,6,比较SSE变化



聚类结果总览

基于SPSS分析,确定最优聚类数K=3,成功将500名客户划分为三个特征鲜明的群体。



三大客户群体关键指标对比

客户群体	人数占比	平均年龄	平均月收入	平均月消费	平均购买频次
群体1	189人(37.8%)	52.0岁	7,205元	854元	8.9次
群体2	161人(32.2%)	46.2岁	12,447元	1,870元	5.0次
群体3	150人(30.0%)	28.0岁	3,560元	1,247元	16.1次

三个群体在年龄、收入、消费金额和购买频次上呈现显著差异,为精准营销提供了清晰的客户画像。

群体1:理性消费型客户



数学特征

- 年龄中心: $\mu_{1,age} = 52.0$, 标准差适中
- 收入中心: $\mu_{1,income} = 7205$, 属中等收入
- 消费模式: 低金额(854元)+中频次(8.9次)
- 单次消费: 96元/次, 最为理性

商业画像

消费理性, 注重性价比, 购买行为规律, 生活习惯稳定。对价格敏感度中等, 不冲动消费。超市基础客户群, 提供稳定收入来源。

群体1营销策略建议

产品策略

主推性价比商品和自有品牌,突出质量可靠、价格实惠的特点。

商品组合

推广家庭装、经济装商品,满足家庭日常需求。

促销策略

设置定期促销活动,如周特价日,培养固定购物习惯。

服务增值

提供便民服务增强粘性,如免费送货、会员积分等。

群体2:品质追求型客户

数学特征

- 年龄中心: $\mu_{2,age} = 46.2$, 中年群体
- 收入中心: $\mu_{2,income} = 12447$, 高收入群体
- 消费模式: 高金额(1870元)+低频次(5.0次)
- 单次消费: 374元/次, 客单价最高

商业画像

追求高品质,对价格不敏感,注重品质和品牌。购买决策谨慎,客户忠诚度高。高价值客户,贡献最高客单价。

营销策略

推广高端、进口、有机优质商品,提供VIP服务和专属优惠,注重购物环境和服务体验,推出精品礼盒和限量商品。



群体3:价格敏感型客户

数学特征

年龄中心28.0岁,收入3,560元,月消费1,247元,购买频次16.1次,单次消费78元。

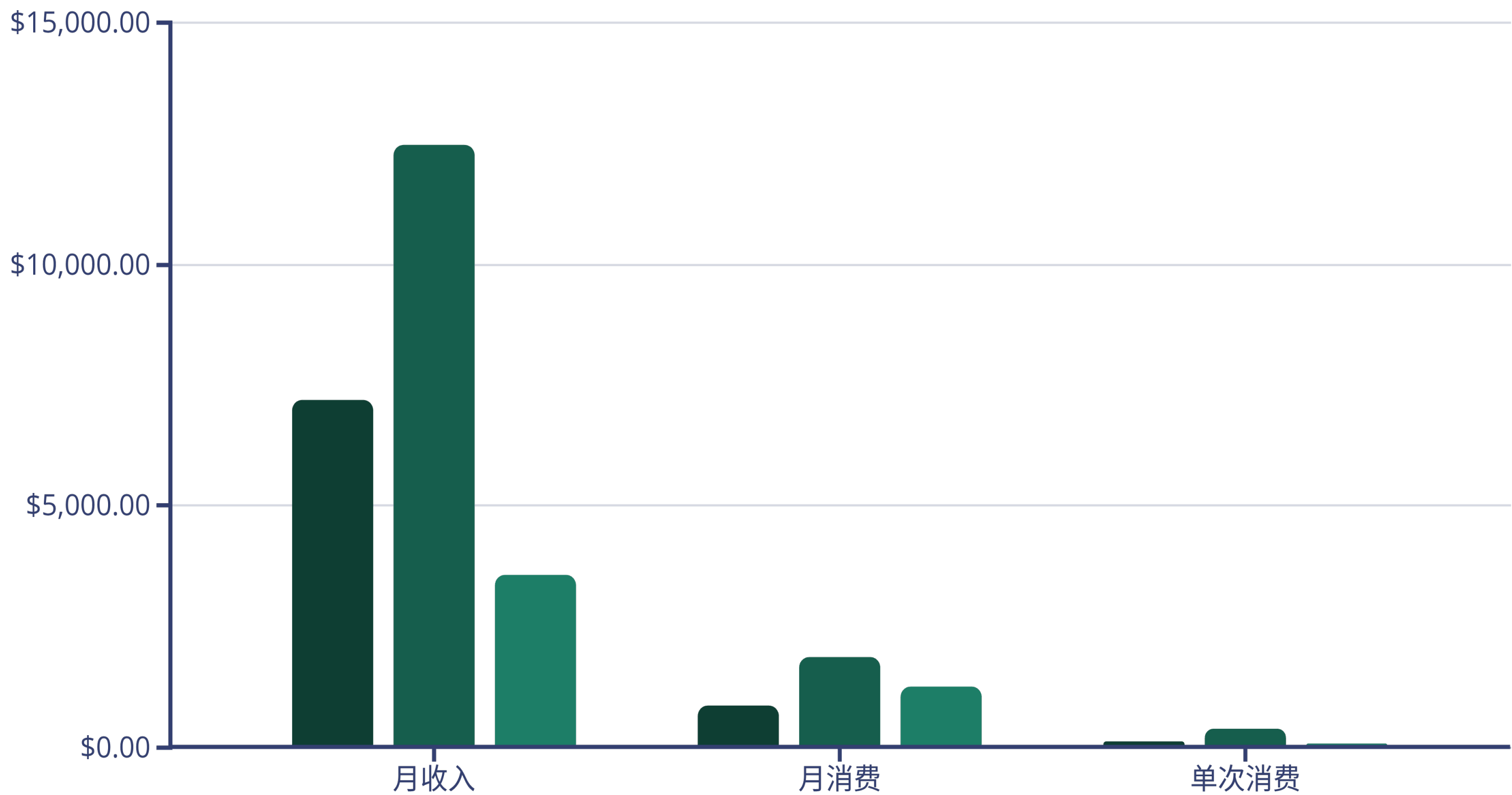
商业画像

价格极度敏感,购买频次最高,日常必需品主力购买者,消费占收入比35%,对促销活动反应强烈。

营销策略

推广特价商品和限时促销,增加小包装、经济型商品,设置日常必需品专区,推出积分兑换和满减活动。

三大客户群体消费行为对比



K值最优选择的综合判定

确定最佳聚类数K是K-means聚类分析中的关键步骤，它直接影响聚类结果的有效性和商业价值。我们通常会结合多种评估方法进行综合判断。

肘部法 (Elbow Method)

原理：随着聚类数K的增大，每个类别的聚合程度会提高，导致类内平方和(SSE)逐渐减小。当K增加到某个值时，SSE的下降幅度会急剧减缓，这个“拐点”即被认为是最佳K值。它平衡了模型复杂度和聚类效果。

$$SSE = \sum_{k=1}^K \sum_{x_i \in C_k} \|x_i - \mu_k\|_2^2$$

轮廓系数法 (Silhouette Method)

评估：轮廓系数结合了聚类的凝聚度（样本与同簇内其他样本的相似度）和分离度（样本与最近邻簇内样本的差异度）。该系数的平均值范围在-1到1之间，越接近1表示聚类效果越好，样本与自身簇的匹配度高且远离其他簇。

$$s_i = \frac{b_i - a_i}{\max(a_i, b_i)}$$

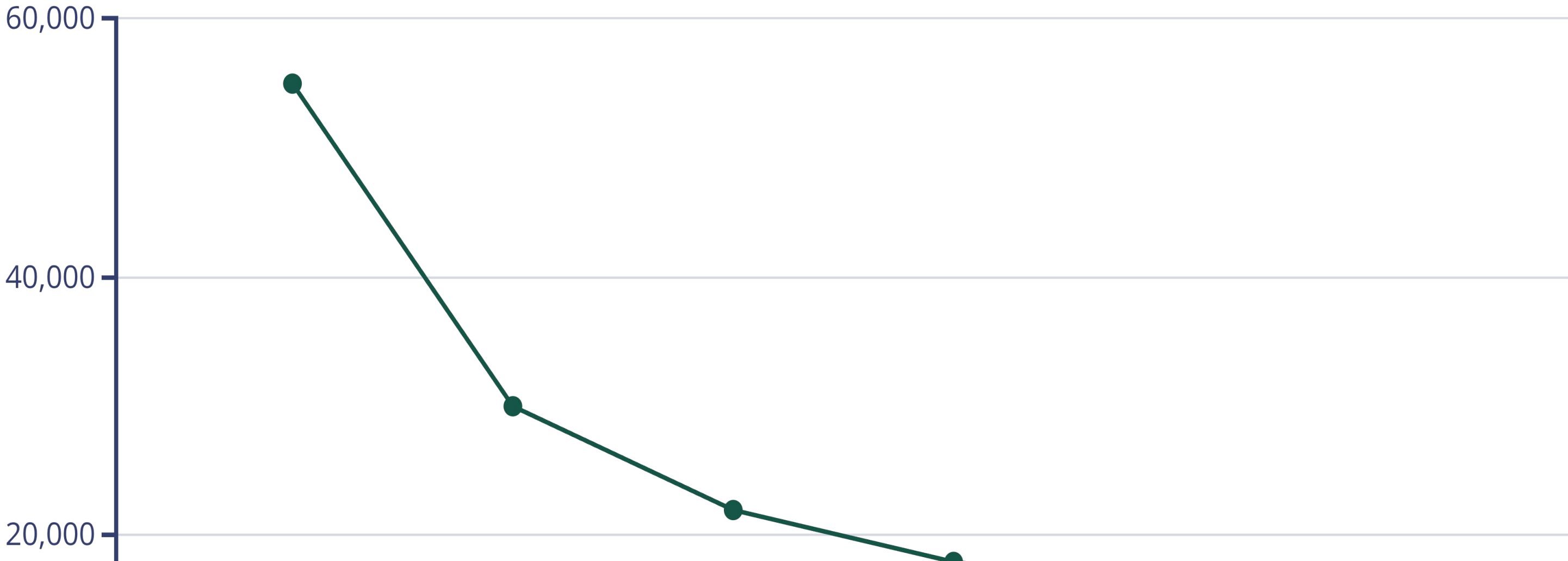
Gap统计量方法 (Gap Statistic Method)

理论：该方法通过比较数据集在不同K值下的总内部分散度与随机参考数据集的内部分散度。理想的K值是使“间隔”（Gap）最大化的那个点，表明当前聚类结构与随机分布相比有显著的改进。

当Gap统计量达到最大值时，通常意味着找到了最优的K值。

综合决策策略

由于各种方法都有其局限性，通常需要将上述多种评估指标的结果进行对比分析，并结合实际业务场景、领域知识和营销资源限制，最终确定一个既有统计学依据又符合商业逻辑的最优K值。





研究结论与应用价值

核心发现

通过K-means聚类分析,成功将500名客户划分为三个特征鲜明的群体:理性消费型(37.8%)、品质追求型(32.2%)和价格敏感型(30.0%)。

模型优势

- 数学模型严谨,聚类结果稳定可靠
- 客户画像清晰,商业解释性强
- 为精准营销提供科学依据

应用价值

- 优化券包发放策略,提高营销ROI
- 制定差异化推荐品类,提升客户满意度
- 确定合理触达频率,避免过度营销
- 指导商品组合和库存管理决策

未来展望

可进一步引入时间序列分析,动态追踪客户群体变化,实现更精准的客户生命周期管理。