

四川大學

本科生课程论文



题目 优化生鲜商超收益：蔬菜商品的销售规律与补货决策分析

| | |
|------|---------------|
| 课程名 | 数据挖掘导引 |
| 任课教师 | 段磊 |
| 学院 | 计算机学院 |
| 专业 | 计算机科学与技术 |
| 学生姓名 | 杨一舟 |
| 学号 | 2022141460176 |
| 年级 | 2022 级 |

2024 年 5 月 21 日

优化生鲜商超收益： 蔬菜商品的销售规律与补货决策分析

计算机学院 2022141460176 杨一舟

摘 要

随着经济的发展，如今新鲜的生鲜果蔬倍受消费者喜爱。但是在生鲜商超中，一般蔬菜类商品隔日就无法再售。故对于每日的生鲜进货量与定价进行规划安排就尤其重要。本文提出了基于商超生鲜蔬菜销售情况的历史数据，通过曲线拟合，k-means 聚类等算法建立的模型，分析了蔬菜各品类及单品销售量的分布规律及相互关系，分析了各蔬菜品类的销售总量与成本加成定价的关系，并通过历史数据给出各蔬菜品类未来一周的日补货总量和定价策略，为生鲜商超提供进货量安排与定价决策参考。

关键词：k-means，LSTM，线性回归，时间序列分析，曲线拟合

一、引言

1.1 问题背景

生鲜商超是人民生活的重要组成部分，也在社会经济中扮演者重要角色。但生鲜蔬菜往往保鲜期较短，大部分种类均不能隔日再售。所以为了保证商超利润，同时确保人们能够买到足量而新鲜的蔬菜，商超需要根据各种商品的销售情况进行“成本加成定价”方法补货。无论是供给侧还是需求侧，进行市场调研，然后形成合理的销售组合都是十分重要的。蔬菜的销售量、销售时间、销售成本与运输损耗等数据则是补货决策与定价决策的关键所在。

1.2 问题提出

数据给出了商超的商品信息、流水明细、批发价格与损耗率等数据，我们可根据这些数据对以下问题提供解决的方案：

问题一：对所给商品数据进行分析处理，建立模型得出蔬菜各个品类直接以及各个单品之间的分布规律和关联关系。

问题二：通过分析各个品类销量与定价的关系，预测未来销售情况及能使收益最大的补货与定价决策。

问题三：对问题二给出的决策进行进一步细化，给出能使商超收益最大的单日单品补货计划。

二、 方法分析

2.1 问题一分析

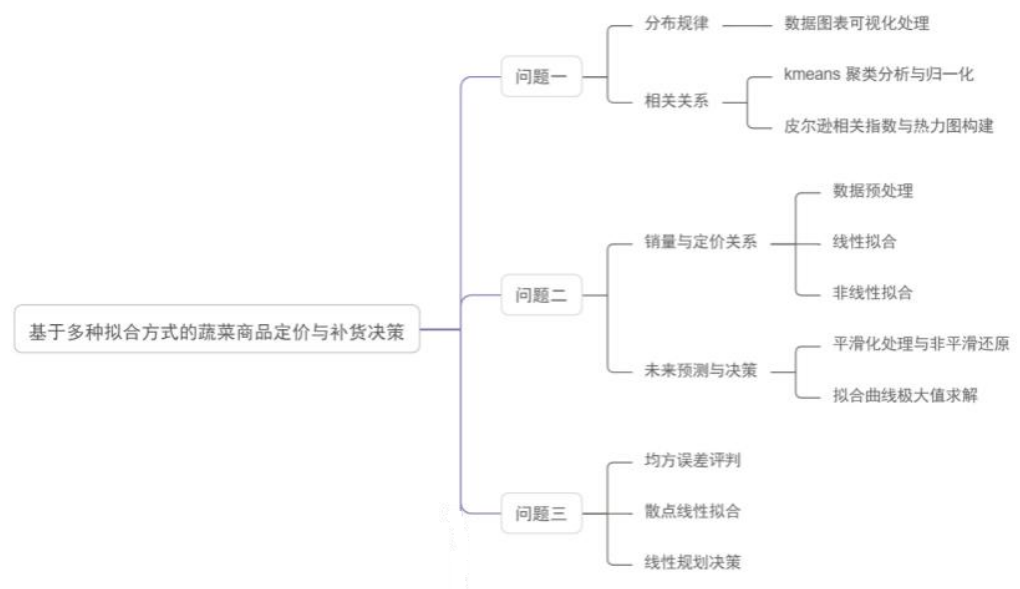
问题一要求分别在商品的不同品类与不同单品之间找出一定的分布规律与相互关系，需要寻找数据的极值与变化趋势。对于数据所给的多种品类，我们首先进行数据可视化处理，然后采用 **k means** 聚类算法来建立模型，通过数据的多次迭代后收敛，发现其中的关联关系。最后我们构建热力图并计算皮尔逊相关指数，进一步验证品类间的关系。

2.2 问题二分析

问题二为预测性问题，要求根据“成本加成定价”方法与历史销售数据对未来一周的销售进行预测并给出补货与定价决策。本文以历史销售价格、进货成本、销售量为参数进行线性拟合，预测未来一周的预计销售量，并基于此做出补货与定价的决策。

2.3 问题三分析

问题三为单目标多决策变量问题，需要首先根据往年的销售与进价数据选出 7 月 1 日的上市菜品单品，然后通过这些数据对单品的 7 月 1 日的销售量和进价进行预测。通过预测得到的数据和题目中给出的陈列量与可售单品总数的约束条件，最后确定一批待选单品名单，并对这些单品的可能利润值进行预测排序。



三、 模型假设

- 1、在预测时间内市场情况，顾客需求等不确定因素保持相对稳定
- 2、未来销售走向可在一定程度上由历史销售记录反映
- 3、各单品的最佳加价策略可由其对应品类的最佳加价策略表示
- 4、蔬菜销量会受到如周末等节假日的影响而相较于平时出现明显波动
- 5、各蔬菜单品的损耗率不变

四、 符号说明

| 符号 | 说明 |
|-----------------------------|---|
| $J(c, \mu)$ | k means 算法中的损失函数， 定义为各个样本距离所属中心点的误差平方和 |
| $\rho(X, Y)$ | 皮尔逊相关指数 |
| S_i | 一周中星期 i 的销售额 |
| δ | 加权调整率， 定义为周末总销售额与工作日总销售额的比值 |
| G | 单日总利润未来预测值 |
| w | 损耗率 |
| $P_{\text{售}}、P_{\text{进}}$ | 售价、进价 |
| M | 均方误差 |
| $Y_i、\hat{Y}_i$ | 实际数据点、拟合预测数据点 |

五、模型的建立与求解

5.1 分布规律与关联关系

在建立模型前，我们对数据进行了一定的预处理：首先我们对于附件 2 中的数据进行了异常值处理，对与异常值，我们采用统计学中四分位距(IQR)法，将各单品的销量与价格数据从低到高排序，再将数据分为 4 个相等的部分，并指定第一、二、三、四分位数是 Q_1 、 Q_2 、 Q_3 、 Q_4 ，IQR 为 Q_3 与 Q_1 的差，由此计算数据集的上下限分别为 $Q_1-1.5IQR$ 、 $Q_3+1.5IQR$ ，将大于上限或低于下限的数据视为异常值，并进行平均值填充处理，即选用异常值相邻两数据的平均值代替异常值。

5.1.1 蔬菜各品类及单品销售量的分布规律

通过统计并分析各品类蔬菜销量按月份分布与按品类分布的情况，我们得到了以下数据图：

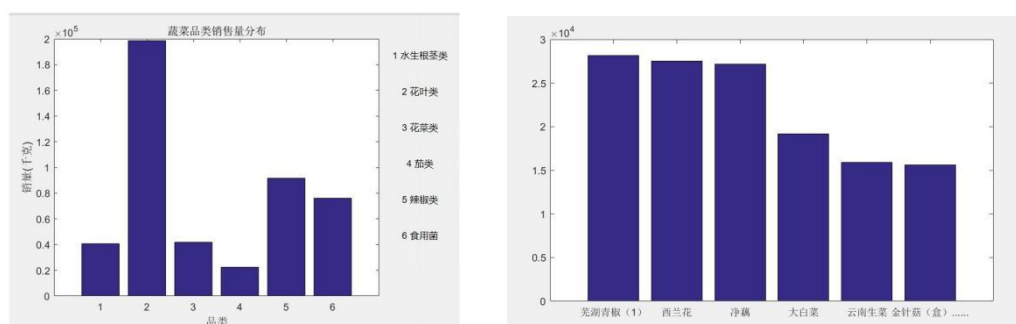


图 1 蔬菜品类及单品销售量分布条形图

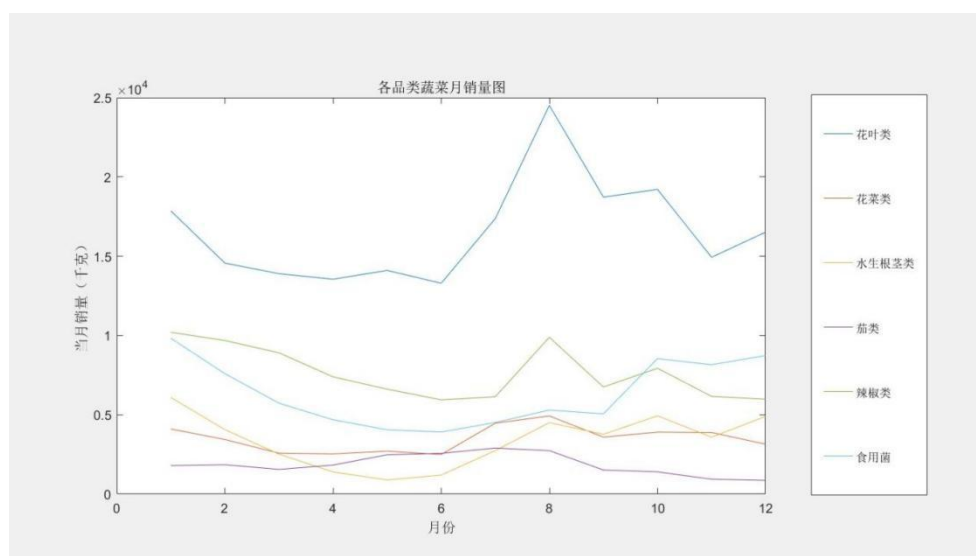


图 2 各品类蔬菜月销量折线图

由图 1 与图 2 可以看出：

对于各个品类，不同品类的蔬菜销量差距较大，花叶类蔬菜销售最多，其次是辣椒类、食用菌类，而茄类蔬菜销售最少；

由蔬菜单品总销量排名较高的几种单品的条形图，可以看出芜湖青椒、西兰花等是销量最高的几种蔬菜单品。

按照月份，蔬菜销量在不同月份的波动较大，表现出明显的季节性，销售旺季主要集中在 8 月前后，而 5 月前后为蔬菜销售淡季。

5.1.2 蔬菜各品类及单品销售量的相互关系

为了刻画蔬菜类商品不同品类或不同单品之间可能存在的关联关系，我们基于不同类型商品的销售次数与总销售量，对不同单品进行分类分析处理。本问题选用 k means 聚类算法来构建模型，使用 MATLAB 对所给数据进行聚类划分处理，寻找对于不同品类蔬菜最合适的分类数量（即 k 值）来体现关联关系，其中 k 值用公式（1）来进行确定：

$$D = \text{类内平均距离} / \text{类间平均距离} \quad (1)$$

首先随机给定一些质心点位置，同时以商品品类的销售次数为横坐标，以总销售量为纵坐标确定各个数据点 x_i ，接着通过迭代寻找 k 个簇使得聚类结果对应的损失函数最小，对于每一个数据点 x_i ，去计算它离哪一个质心 u_j 最近并将其分配到最近的质心，如公式（2）所示：

$$c_i^t = \arg \min_k \|x_i - \mu_k^t\|^2 \quad (2)$$

然后对于每一个簇中心 k，重新计算该簇的中心，如公式（3）所示：

$$\mu_k^{(t+1)} = \arg \min_{\mu} \sum_{i: c_i^t = k} \|x_i - \mu\|^2 \quad (3)$$

其中损失函数定义为各个样本距离所属中心点的误差平方和，如公式（4）所示：

$$J(c, \mu) = \sum_{i=1}^M \|x_i - \mu_{c_i}\|^2 \quad (4)$$

本方法通过调整每个数据点所属的类来减少 J ，再固定每个样本的类别，调整中心点继续减小 J 。两个过程交替迭代，使 J 单调递减直到最（极）小值，此时中心点和样本划分的类别数量同时收敛，从而可以得到各个品类之间存在的一定的联系。

经过迭代后得到收敛的数据，再进行归一化之后得到以下聚类分布图

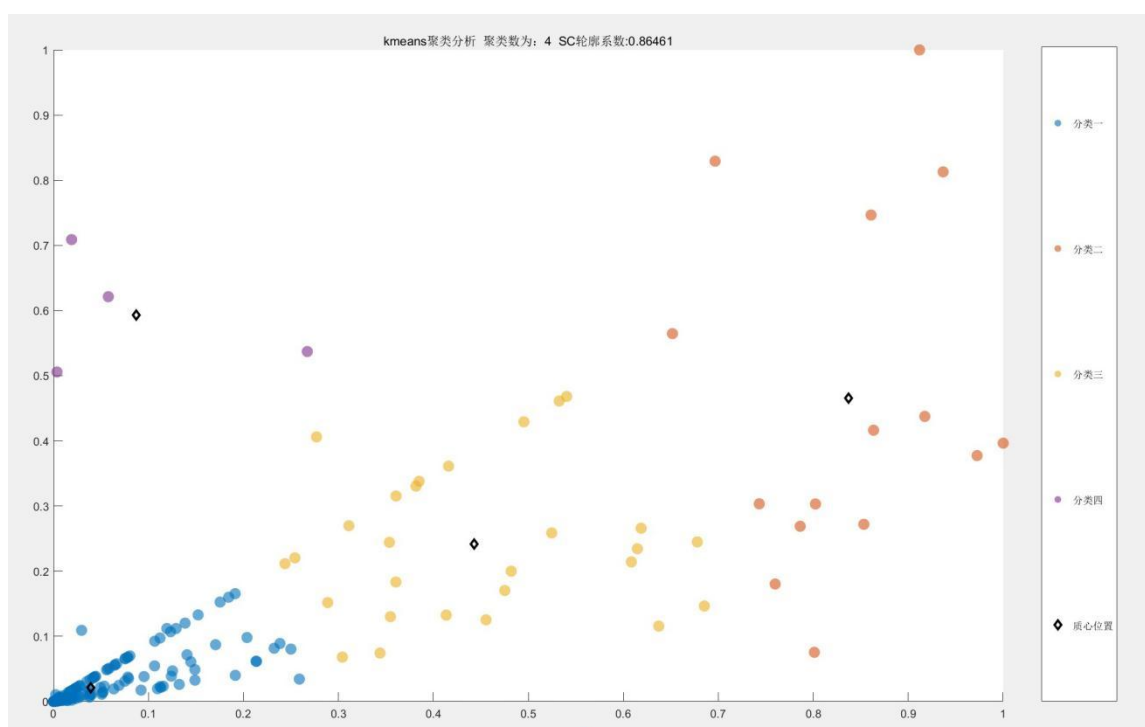


图 3 k means 聚类分析图

由图 3 可知，对于所有单品，所给数据大致可分为四个簇，且 SC 轮廓系数达到 0.86461，说明每个簇内部关联性较强，聚类效果较好。

且大部分数据点分布于分类一（蓝色）区域内，说明大部分品类的需求量都较小，顾客多以少次、少量的方式进行购买。

同时我们以销量，销售次数和销售时间为相关系数的评价指标，构建了包含皮尔逊相关指数的热力图，相关指数计算方法如公式（5）：

$$\rho(X, Y) = \frac{cov(X, Y)}{\sigma_X \sigma_Y} = \frac{E[(X - \mu_X)(Y - \mu_Y)]}{\sigma_X \sigma_Y} \quad (5)$$

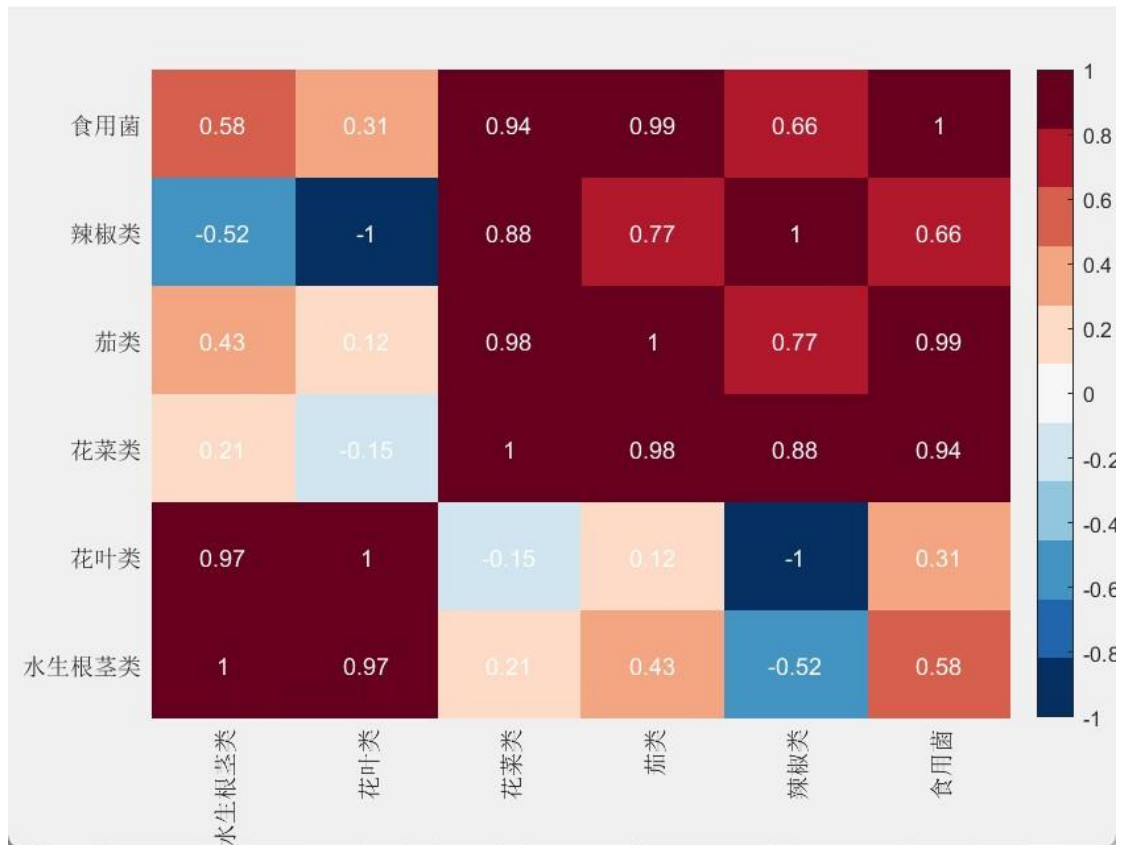


图 4 各品类间相关情况热力图

由图 4 可以看出，不同品类间的关联度差异较大。

辣椒类与花叶类的皮尔逊相关指数达到-1，有较高的负相关性，而茄类与食用菌类的相关指数达到 0.99，有较高的正相关性。同时花叶类与水生根茎类、茄类与花菜类、辣椒类与花菜类、食用菌类与花菜类的关联度也较高。这意味着其中一类蔬菜的销量增加或减少时，另一类蔬菜的销量大概率也会对应地增加或减少，且两类蔬菜在销售时间上存在较高的相关性。

花叶类与花菜类的相关指数仅为-0.15，关联度最低，同时花菜类与水生根茎类、花叶类与茄类的关联度也较低。这意味着这两类蔬菜的销售较为独立，相互之间的影响与联系均较少。

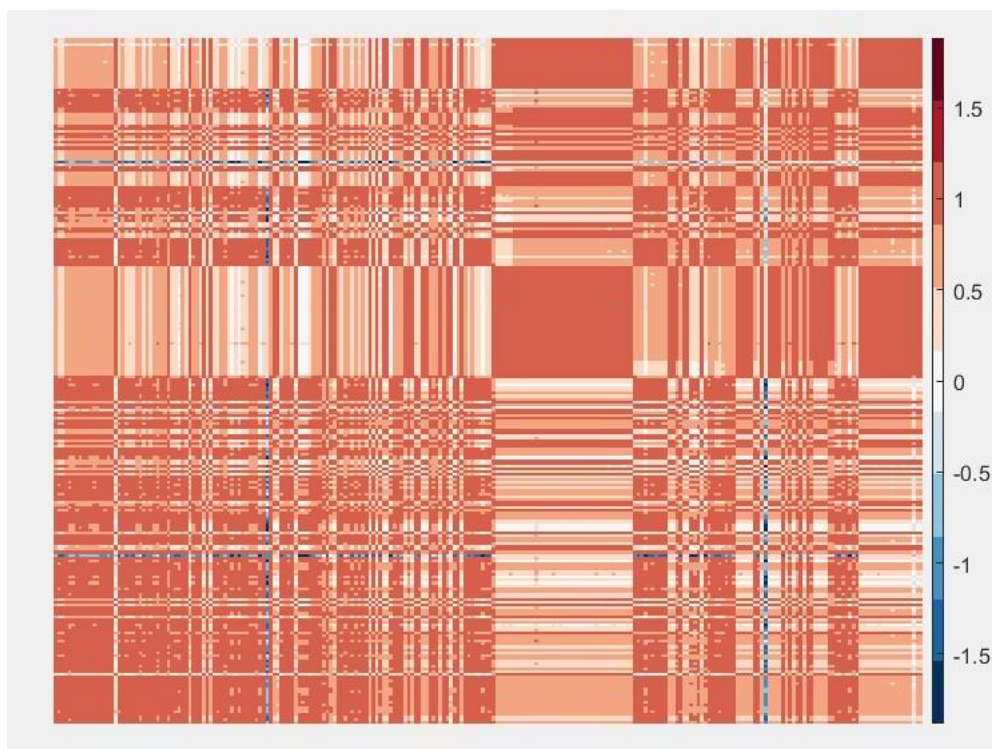


图 5 各单品间相关情况热力图

由图 5 的数据中可以分析得到，红珊瑚(粗叶)102900011033562、绿牛油(102900011033586)、红橡叶 102900011033562 两两之间的皮尔逊相关系数最高，分别为 0.991，0.987 和 0.990，说明其在销售量及购买时间上均存在较高的相关性。

5.2 补货与定价决策

5.2.1 各蔬菜品类的销售总量与成本加成定价的关系

依据对各个品类数据的初步分析，我们猜测单日销售总量与成本加成定价之间可能存在线性拟合关系。分别以各个品类的成本加成率与销量为横纵坐标，我们得到了各个品类的数据散点图。其中成本加成率由以下方法确定：

$$\text{成本加成率} = \frac{\text{未打折售价}}{\text{进价成本}} - 100\%$$

同时采用最小二乘法对各个品类进行线性拟合，得到了如图 6 所示的各个散点图及拟合直线。

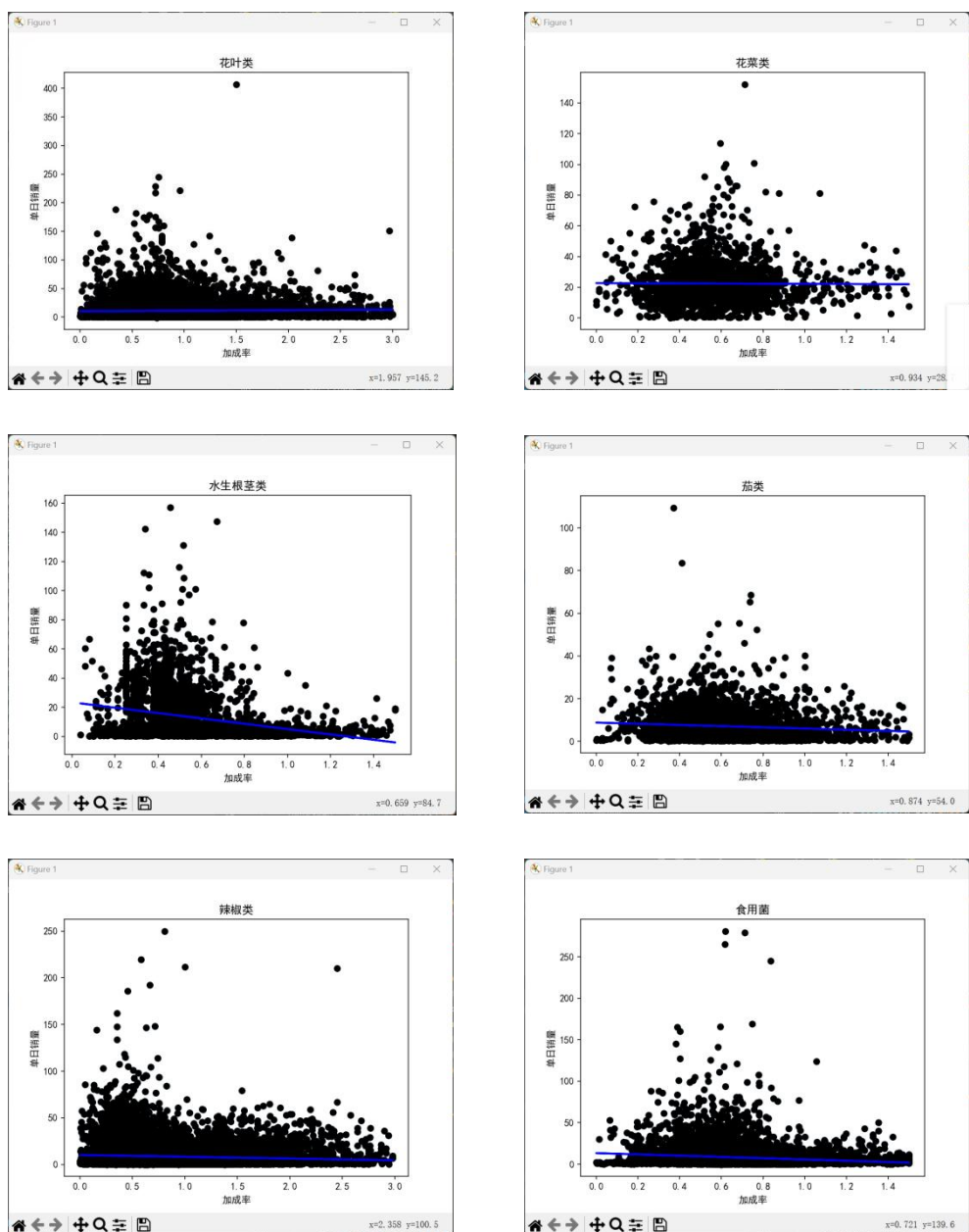


图 6 各品类散点图及其拟合直线

在进行拟合之后，我们发现直线不能很好地拟合散点，即直线的拟合效果并不理想，无法有效确定使利润最大的加成率，故最终未采用此方案。

在线性拟合的尝试效果不甚理想之后，我们尝试将单日总利润作为坐标纵轴，使用加成率作为坐标横轴，进行非线性模型拟合，从而反映销售总量与成本加成定价之间的关系。其中单日总利润由以下方法确定：

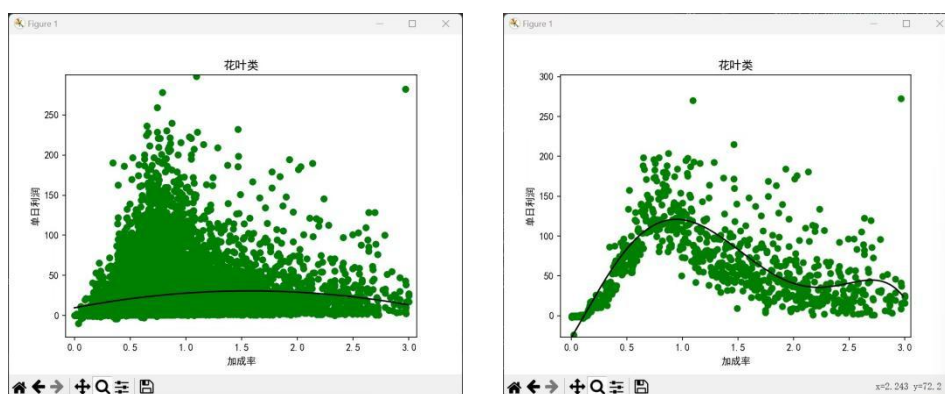
$$\text{单日总利润} = \text{总销售额} - \frac{\text{单位进货成本} \times \text{单日销售量}}{1 - \text{损耗率}\%}$$

对于处理之后以利润为纵轴的散点图，我们转而尝试非线性曲线拟合处理，且发现其中三次曲线拟合效果相对较好，得到如图 6 中左侧的拟合曲线。但是此曲线拟合程度也并不理想，不能完全准确地反映散点图的数据特点。

与此同时，我们发现散点图总是在低于某个曲线的范围内波动，这也与现实生活中人们的购物情况相符合，即在通常情况下人们的购物量趋向于在一个较合理范围内稳定波动。此二者逻辑相洽，于是我们尝试拟合此极大值曲线。

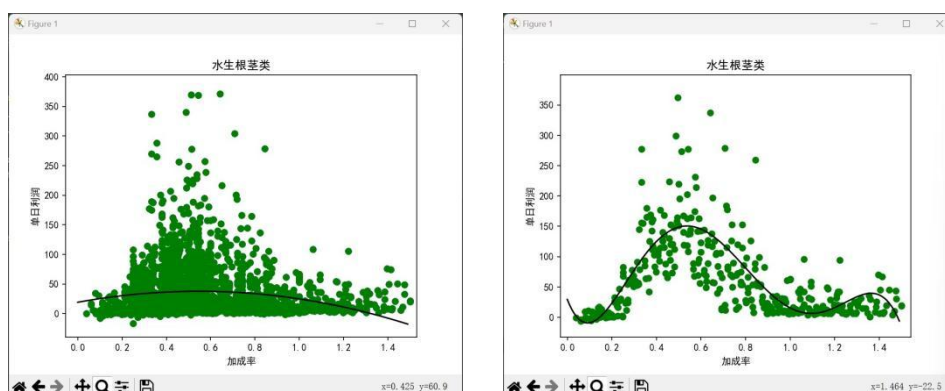
我们对散点数据再次进行处理，在横轴上以每 0.01 的加成率为一个小区间，在区间内保留利润最高的 3 个散点。由此我们得到如图 7 右侧的图像与拟合曲线及其曲线方程。可以发现，此时的曲线拟合程度较好，可由此刻画销售总量与成本加成定价之间的关系。

注：左侧为第二次预处理之前图像，右侧为第二次预处理之后图像，图下侧所示方程为右图曲线对应方程



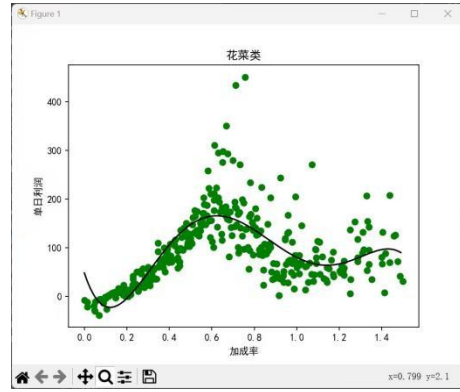
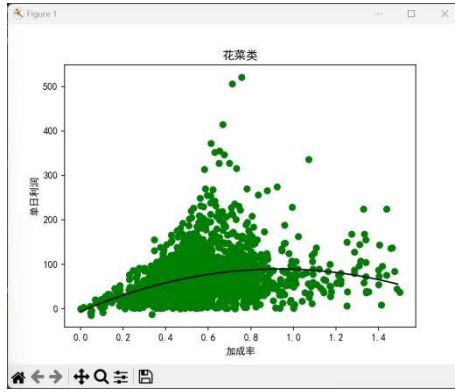
$$f_1(x) = -23.64x^5 + 163.5x^4 - 341.7x^3 + 121.4x^3 + 228.4x^2 - 27.49$$

极大值点 (0.94517, 120.972)



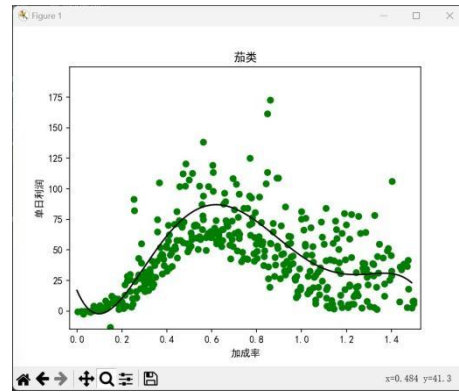
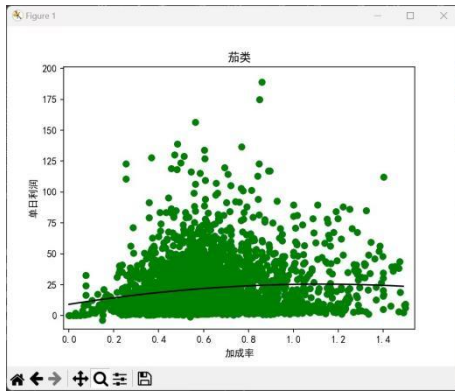
$$f_2(x) = -2604x^5 + 10050x^4 - 13390x^3 + 6883x^2 - 948.2x + 29.46$$

极大值点 (0.53805, 151.06744)



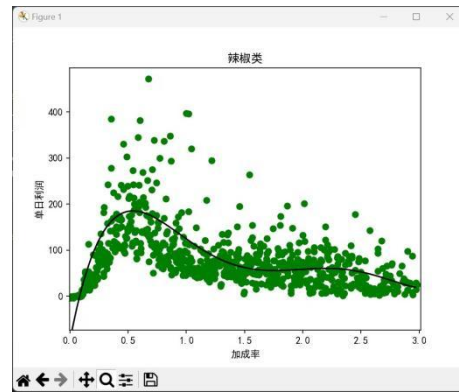
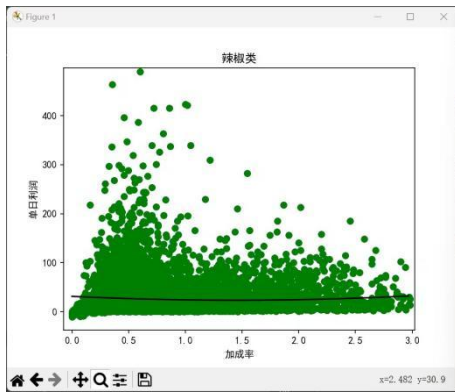
$$f_3(x) = -2163x^5 + 8965x^4 - 13050x^3 + 7614x^2 - 1330x + 48.22$$

极大值点 (0.62463, 166.84442)



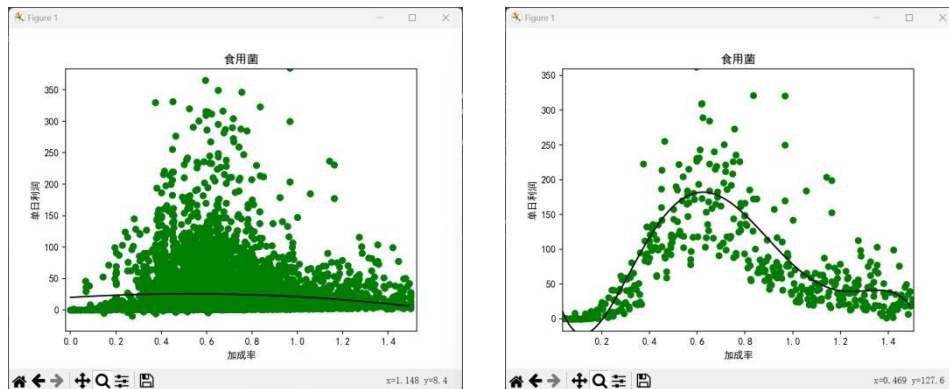
$$f_4(x) = -840.9x^5 + 3492x^4 - 5075x^3 + 2885x^2 - 432.5x + 16.87$$

极大值点 (0.6201, 87.15144)



$$f_5(x) = 41.46x^5 - 392x^4 + 1367x^3 - 2124x^2 + 1327x - 96.16$$

极大值点 (0.53967, 184.8867)



$$f_6(x) = -2190x^5 + 9160x^4 - 13510x^3 + 7946x^2 - 1382x + 50.65$$

极大值点 (0.62221, 181.32664)

图 7 各品类非线性拟合模型

5.2.2 各蔬菜品类未来一周的补货总量预测与最佳定价策略

通过分析图 7 曲线得到的极大值点，我们可以得出各个品类的最佳定价策略如下表所示：

| 品类 | 加成率（基于成本的倍数） |
|-------|--------------|
| 花叶类 | 94.517% |
| 水生根茎类 | 53.805% |
| 花菜类 | 62.463% |
| 茄类 | 62.010% |
| 辣椒类 | 53.967% |
| 食用菌类 | 62.221% |

其中定价可以由成本与成本加成率按以下方法确定：

$$\text{售价} = \text{成本} \times (1 + \text{成本加成率})$$

对于补货总量的预测，我们首先对历史数据进行了可视化分析，如图 8 所示：

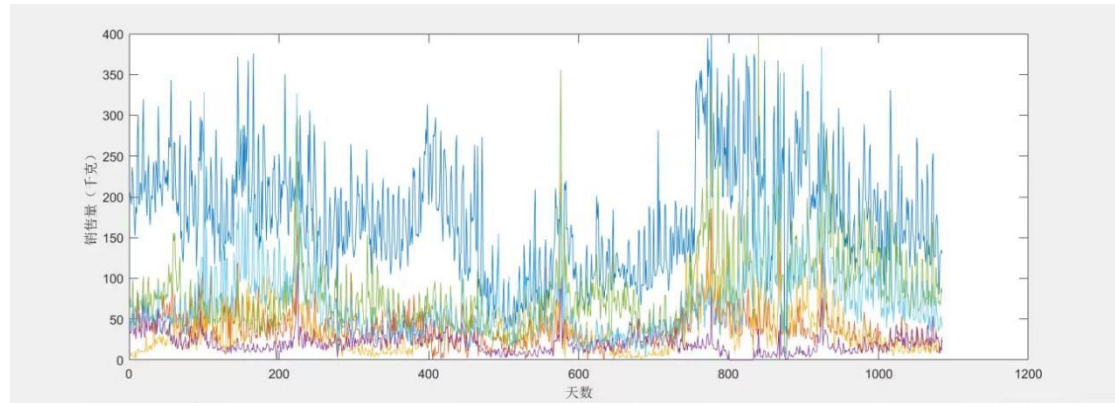


图 8 历史销售数据折线图

通过对此折线图的分析，我们发现销售量对时间出现了明显的周期性波动，同时根据此数据与日期的对应分析，我们发现销售量的峰值往往出现在周末，而谷值往往出现在工作日。

为了提高拟合效果，使得曲线更平滑更便于拟合，我们定义加权调整率 δ 为周末销售额与工作日销售额的比值，通过以下公式来确定 δ 的值：

$$\delta = \frac{S_6 + S_7}{\sum (S_1 + S_2 + S_3 + S_4 + S_5)}$$

对于周末销售量 S_6 、 S_7 ，我们通过除以调整率来进行平滑化处理，即：

$$S_i = \frac{S_i}{\delta}$$

数据处理样例如下所示：

| | | | | | | | | | | | | | |
|---------------------|----------|---------|--------|---------|---------|---------|---------------------|---------|--------|--------|--------|--------|--------|
| 2021-06-19 00:00:00 | 151.0035 | 44.969 | 9.057 | 7.256 | 54.473 | 27.4985 | 2021-06-19 00:00:00 | 203.369 | 44.511 | 10.313 | 30.488 | 87.633 | 39.484 |
| 2021-06-20 00:00:00 | 149.9415 | 44.072 | 8.745 | 38.649 | 50.974 | 26.5925 | 2021-06-20 00:00:00 | 202.307 | 43.614 | 10.001 | 61.881 | 84.134 | 38.578 |
| 2021-06-26 00:00:00 | 128.2535 | 45.75 | 10.91 | 18.8285 | 42.3295 | 31.976 | 2021-06-26 00:00:00 | 197.69 | 71.821 | 20.953 | 35.449 | 82.439 | 47.402 |
| 2021-06-27 00:00:00 | 144.1105 | 30.471 | 9.939 | 15.1555 | 39.0115 | 33.714 | 2021-06-27 00:00:00 | 213.547 | 56.542 | 19.982 | 31.776 | 79.121 | 49.14 |
| 2021-07-03 00:00:00 | 155.2325 | 47.0725 | 15.59 | 43.0985 | 60.8155 | 48.5605 | 2021-07-03 00:00:00 | 197.925 | 48.497 | 29.657 | 47.285 | 73.552 | 66.849 |
| 2021-07-04 00:00:00 | 154.6295 | 44.2545 | 10.147 | 35.0545 | 44.8465 | 26.1775 | 2021-07-04 00:00:00 | 197.322 | 45.679 | 24.214 | 39.241 | 57.583 | 44.466 |

对于平滑化处理后的数据，我们对其进行非线性拟合后可得到效果较好的预测折线图。然后我们对折线图中周末的销售量再重新乘以调整率 δ 进行非平滑还原处理，得到如图 8 所示的折线图，能够更准确刻画 7.1 至 7.7 的销售数据，并由此来确定日补货总量：

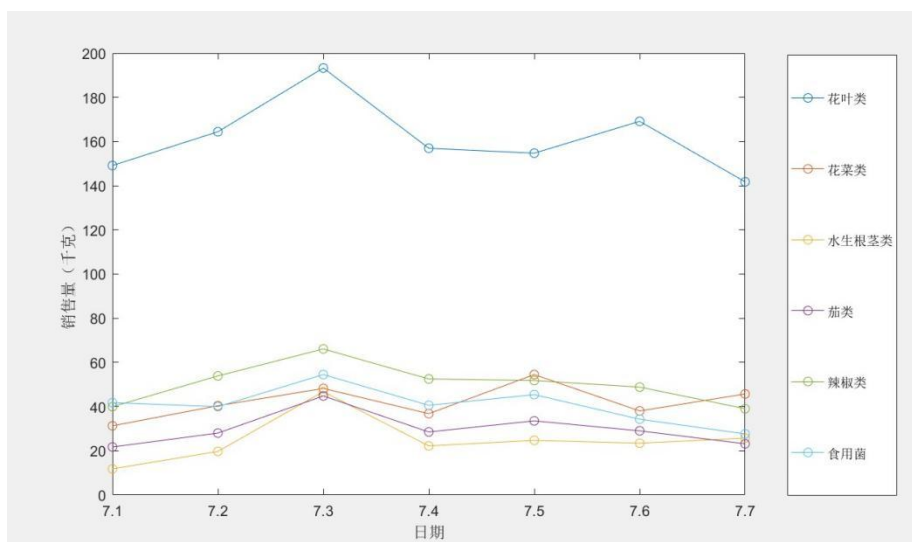


图 9 7.1—7.10 销售量预测折线图

5.3 细化单日单品补货计划

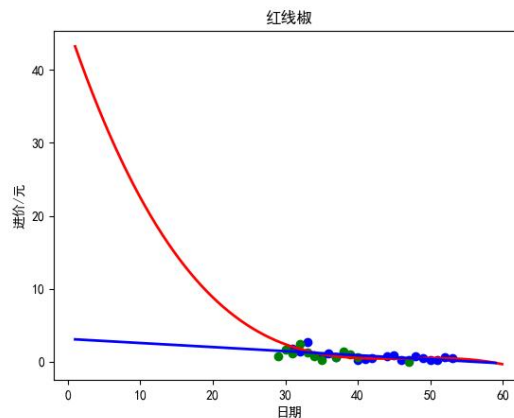
问题三要求我们在进一步考虑可售单品总数、单品订购量与可售品种等限制之下，对未来单日做出单品补货量与定价的最优决策。

在数据预测方法方面，对于未来销售量与未来进价的预测，我们尝试了均值代替以及多种方式拟合。虽然进价的价格波动较小，但是对比线性拟合与均值的准确效果，最终选择了拟合方式来预测数据。

在拟合方式中，我们发现线性拟合与三次曲线拟合是效果较好两种拟合方式。为了确定最终拟合方式，我们采用均方误差算法（MSE）来构建损失函数如下所示：

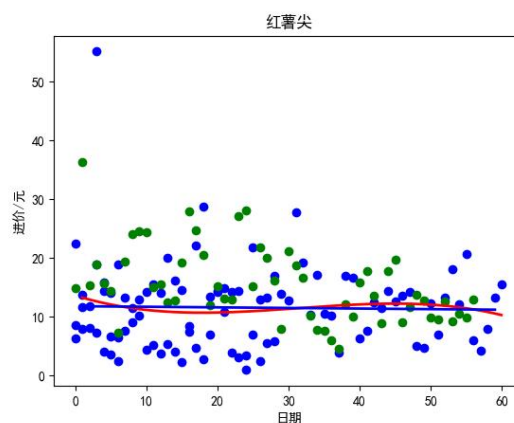
$$M = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

我们使用 MSE 损失函数观察并评判之后发现，在大部分情况下二者拟合效果相似，而在某些极端情况下线性拟合效果比三次曲线更好，如图 9、图 10 所示，所以我们选择了线性拟合作为最终拟合方式。



| | 红线椒 |
|-------|-------------|
| 线性损失 | 1.518439209 |
| 非线性损失 | 145.807092 |

图 10 极端情况下两种拟合方式效果对比



| | 红薯尖 |
|-------|-------------|
| 线性损失 | 30.16197928 |
| 非线性损失 | 29.64066489 |

图 11 普遍情况下两种拟合方式效果对比

在数据内容选取方面，考虑到蔬菜销量的季节性变化，所以我们选取了 6.1-7.31 的数据来进行拟合，以便减小季节性差异对于特定时间节点预测效果的影响，更好地对 7.1 日的进货方案进行预测规划。

由于每个单品数据较少，数据点分散，无法准确拟合，所以我们选用大品类的加价率代表单品的加价率。在对所有品类进行进价与销量的拟合分析后，我们得到各个品类的进价与销量预测如图所示：

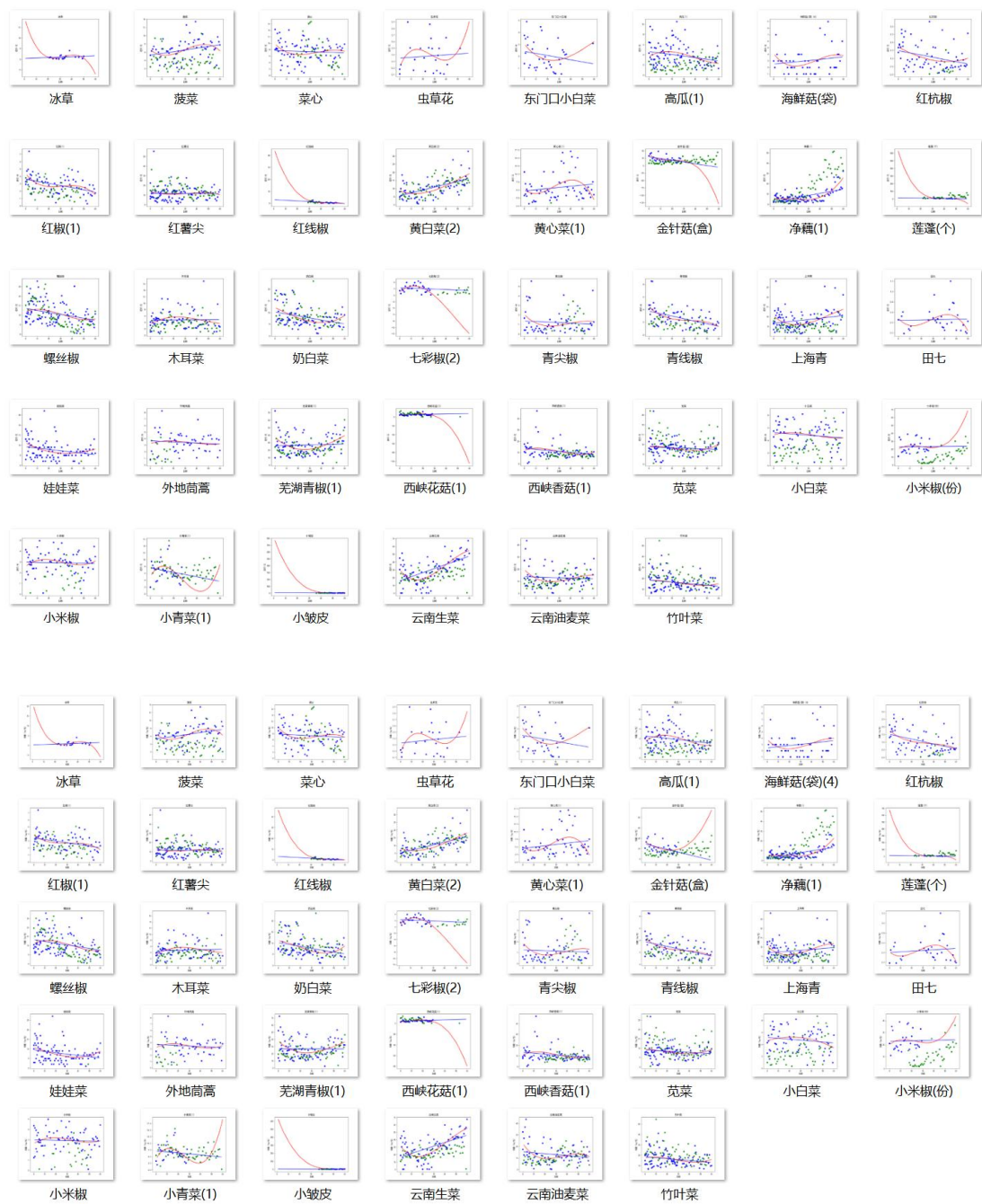


图 12 各品类进价及销量预测拟合分析

在获取销售量与进价的预期数据后，根据对题意信息的分析我们建立了总利润关于预计销售量、售价、预计进价与损耗率的目标函数如下所示：

$$G = \sum_{i=1}^n (S_i \times (P_{\text{售}} - P_{\text{进}}) \times (1 - w\%))$$

$$P_{\text{售}} = P_{\text{进}} \times (100\% + \text{加价率})$$

$$\begin{cases} 27 \leq n \leq 33 \\ S_i \geq 2.5 \end{cases}$$

对此函数进行线性规划分析，我们可以得出使得利润 G 最大的决策，即最合适的定价与进货量。最终决策结果如下表所示：

表 1 使利润最大的各品类补货量及定价决策

| | 预测利润 | 预测进价 | 加价率 | 定价 | 补货量（预测销量） | 损耗率 |
|-----------|---------|---------|--------|---------|-------------|--------|
| 云南生菜 | 95.067 | 4.66093 | 94.517 | 9.0663 | 25.46289461 | 0.1525 |
| 西峡香菇(1) | 70.6894 | 11.9545 | 62.221 | 19.3927 | 11.0275408 | 0.1382 |
| 小白菜 | 43.2779 | 9.65447 | 94.517 | 18.7796 | 5.132260305 | 0.0759 |
| 芜湖青椒(1) | 33.571 | 3.0547 | 53.967 | 4.70323 | 21.59511161 | 0.057 |
| 菠菜 | 31.7375 | 6.20428 | 94.517 | 12.0684 | 6.641505584 | 0.1851 |
| 净藕(1) | 27.5489 | 11.1514 | 53.805 | 17.1515 | 4.860748527 | 0.0554 |
| 黄白菜(2) | 26.6457 | 3.66899 | 94.517 | 7.13681 | 9.104986912 | 0.1561 |
| 竹叶菜 | 26.3807 | 2.17686 | 94.517 | 4.23437 | 14.84333818 | 0.1362 |
| 红薯尖 | 24.6228 | 2.45354 | 94.517 | 4.77255 | 11.59401935 | 0.0842 |
| 云南油麦菜 | 24.0692 | 2.77607 | 94.517 | 5.39992 | 10.52097458 | 0.1281 |
| 螺丝椒 | 23.4109 | 4.51523 | 53.967 | 6.95196 | 10.69639425 | 0.1018 |
| 娃娃菜 | 23.0752 | 4.21396 | 94.517 | 8.19686 | 5.940881349 | 0.0248 |
| 外地茼蒿 | 22.7614 | 9.34698 | 94.517 | 18.1815 | 3.489198152 | 0.2616 |
| 上海青 | 21.1035 | 3.72193 | 94.517 | 7.23978 | 7.010600877 | 0.1443 |
| 黄心菜(1) | 20.2112 | 3.42683 | 94.517 | 6.66577 | 6.983065988 | 0.1064 |
| 西峡花菇(1) | 19.3418 | 15.1529 | 62.221 | 24.5812 | 2.29984557 | 0.108 |
| 金针菇(盒) | 18.9596 | 1.57753 | 62.221 | 2.55908 | 19.40329349 | 0.0045 |
| 菜心 | 18.7901 | 3.37373 | 94.517 | 6.56248 | 6.828072002 | 0.137 |
| 小米椒(份) | 16.1982 | 1.53664 | 53.967 | 2.36592 | 21.56657427 | 0.0943 |
| 木耳菜 | 14.0728 | 2.92595 | 94.517 | 5.69148 | 5.507806028 | 0.0761 |
| 茼菜 | 12.0722 | 2.14581 | 94.517 | 4.17396 | 7.305232819 | 0.1852 |
| 小米椒 | 11.0427 | 7.17286 | 53.967 | 11.0438 | 3.030269973 | 0.0586 |
| 奶白菜 | 10.6972 | 2.1346 | 94.517 | 4.15217 | 6.288013989 | 0.1568 |
| 红线椒 | 8.70468 | 7.28518 | 53.967 | 11.2168 | 2.398997079 | 0.0771 |
| 高瓜(1) | 8.66634 | 5.54098 | 53.805 | 8.52231 | 4.108656901 | 0.2925 |
| 莲蓬(个) | 8.33402 | 1.38733 | 53.805 | 2.13378 | 12.71188895 | 0.1217 |
| 红椒(1) | 6.39216 | 4.97878 | 53.967 | 7.66568 | 2.696067784 | 0.1176 |
| 小青菜(1) | 5.98191 | 1.96298 | 94.517 | 3.81833 | 3.595564168 | 0.1033 |
| 青线椒 | 5.60698 | 4.72179 | 53.967 | 7.27 | 2.386508555 | 0.078 |
| 冰草 | 5.00219 | 7.91251 | 94.517 | 15.3912 | 0.786988524 | 0.1501 |
| 红杭椒 | 3.81162 | 8.33279 | 53.967 | 12.8298 | 0.941671993 | 0.0999 |
| 虫草花 | 3.39293 | 16.7188 | 62.221 | 27.1214 | 0.369295445 | 0.1168 |
| 东门口小白菜 | 3.39067 | 2.90194 | 94.517 | 5.64476 | 1.713130992 | 0.2784 |
| 海鲜菇(袋)(4) | 2.62443 | 1.79754 | 62.221 | 2.91598 | 2.640369116 | 0.1113 |
| 田七 | 1.47581 | 3.80488 | 94.517 | 7.40114 | 0.494127537 | 0.1695 |
| 青尖椒 | 1.47313 | 3.2923 | 53.967 | 5.06905 | 0.888845559 | 0.0672 |

5.4 模型的分析与检验

1、在第二问的模型中，为了改变得到拟合的最大值，我们多次调整采样区域宽度和采样点数量，最后取得了一个较为合适的值。

2、在第三问的模型中，我们在对非上市蔬菜进行预处理后，尝试了对不同长度的时间范围区间的数据进行分析拟合。

首先尝试了对多年的数据利用循环神经网络进行预测。发现在宏观上数据呈现周期性波动，在而微观上数据波动随机性大，且样本量少，容易出现过拟合情况，所以我们选择对短区间内进行线性拟合与非线性拟合。

3、多次尝试后，发现线性拟合和三次曲线拟合有较好的数据。为了决定最终方法，我们抽出了一些测试点使用 **MSE** 损失函数来计算拟合损失，发现在大部分情况下，二者表现均较好。但是在某些极端情况下，三次曲线拟合会出现过拟合的情况，所以我们最后选取了线性拟合。

六、 模型的评价、改进与推广

6.1 模型的优点

1、第二问中的模型符合消费者购物现实情况，创新性地使用最大值采样的方法对销售数据进行预处理后再拟合，在一定程度上准确且巧妙描述了问题，结果相对与直接拟合更合理直观。

2、模型对大多数数据都能够很好的解释，具有普遍性

3、第三问中的模型以每天为单位，同时考虑了周末与工作日的差别，从微观层面上预测数据，相比与宏观分析更精确，更有说服力。

4、模型结果大多能以图表的形式展示，直观易懂，便于理解操作。

6.2 模型的缺点

1、在第二问中对于单个商品的数据无法表现出很好的适用性，日销售额数据量少，对日补货量的预测拟合程度不强。

2、由于每个单品数据较少，数据点分散，无法准确拟合，所以我们在第三问中选用大品类的加价率代表单品的加价率，不能非常精确地代表每个单品的加价率数据。

3、每日销售额与进价存在诸多变数，如天气、节假日、特殊活动等因素，无法给出补货量更精细的进一步预测与决策建议。

6.3 模型的改进与推广

1、可以在第二问模型获取了足够多的单品数据后，对单品再进行拟合，调整单品的加价策略，同时采用使用机器学习等方法。

2、对于日补货量，可以在获取更多数据后，利用循环神经网络进行模拟预测。

七、 参考文献

- [1]冯海康. 基于 k-means 和 LSTM 的量化策略研究[D]. 广州大学,2023.DOI:10.27040/d.cnki.ggzdu.2023.002181.
- [2]刘恒宇. 果蔬产品库存、生产与销售策略研究[D].北京交通大学,2019.
- [3]徐琪,张阳奎.基于 ARIMA 模型与随机森林模型的零售服装动态销售预测模型探析[J].中国管理信息化,2022,25(06):100-104.
- [4]Zeng Q ,Gong Z ,Wu S , et al.Measuring cyclists' subjective perceptions of the street riding environment using K-means SMOTE-RF model and street view imagery[J].International Journal of Applied Earth Observation and Geoinformation,2024,128103739-.
- [5]Kulisz M ,Kłosowski G ,Rymarczyk T , et al.The use of the multi-sequential LSTM in electrical tomography for masonry wall moisture detection[J].Measurement,2024,234114860-.

八、 附录

数据挖掘相关代码及可视化成果展示

问题一： C1.m

问题二： 问题 2： 单日利润线性预测.py
 问题 2： 最大值采样预测.py
 问题 2： 销售价格汇总.py

问题三： 问题 3： 进价预测.py
 问题 3： 销量预测.py
 问题 3： 利润汇总.py

问题二、三所有图表输出数据： 文件夹 data