

# 基于 XXXXXXXX 的葡萄酒评价

August 24, 2022

## 摘要

葡萄酒备受大部分人的热爱，质量较好的葡萄酒往往带来更好的感官体验。葡萄酒的评价有着不同的指标，但其好坏却主要来自于评酒师的个人主观评价，所以会导致评价结果的差异性和不够客观性，存在着不同大众的口味喜爱偏差。本文为了解决该问题，通过探索酿酒葡萄对葡萄酒质量的影响，通过建立模型的方式更合理的对酿酒葡萄质量分类，以及如何依赖不同的酿酒葡萄和葡萄酒之间的关系来客观的给出葡萄酒评价结果，通过客观修正的数据建立了数学模型，以客观的方式来给出葡萄酒的评价。

针对问题一，由于数据较多，且分布不均的原因，首先通过对数据的筛选处理，发现两组对红葡萄酒和白葡萄酒的打分情况是两两相互比较和配对进行 K-S 检验判断数据是否符合正态分布，再对样本 T 检验进行显著性差异判断，通过方差齐性检验来假设检验，最后依据显著性差异水平来判断打分组别的稳定性和可靠性。

针对问题二，需要对酿酒葡萄进行分类，需要考虑到所酿葡萄酒的好坏。由于评酒员的打分存在主观性的干扰，我们将打分数据重新加权处理求和，得到一份较为中肯的综合评价结果。对酿酒葡萄和葡萄酒理化指标进行标准化处理，去除均值，后利用主成分分析对酿酒葡萄和葡萄酒理化指标提取主成分。将酿酒葡萄的数据和葡萄酒理化指标数据进行关联，采用改进的 K-means++ 分类模型，通过对主成分聚类分析将葡萄酒质量分为六大类，由于部分葡萄酒得分评价数值较低，以此为下限进行分类，借鉴罗伯特帕克葡萄酒评分体系依据酿酒葡萄和葡萄酒的相关性，以酿酒葡萄的理化指标含量和分布来进行排名，分出酿酒葡萄的等级和优劣。

针对问题三，探求酿酒葡萄和葡萄酒的理化指标之间的关系。首先对数据进行筛选处理出关键数据，剔除部分多余数据。提取葡萄酒某指标以及对应的酿酒葡萄理化指标进行相关性分析，通过 Pearson 系数和显著性水平来判断其相关性的强弱。后获取到相关性较强的理化指标数据进行共线性诊断后，进行德宾沃森残差分析后通过多元回归的方法来进行拟合，以 R 方的值来判断拟合的契合度高低，以高的模型数据为准得到葡萄酒和对应酿酒葡萄理化指标之间的系数，建立其函数关系。

针对问题四，借鉴问题二中使用到的主成分分析法先对酿酒葡萄的理化指标进行分类，提取到八大类主成分，分析各主成分中相关较高的酿酒葡萄和葡萄酒的理化指标，从而得到对葡萄酒质量贡献较大的因素指标。对评分数据进行加权求和、去极值处理后。以评分标准数据作为因变量，酿酒葡萄理化指标主成分作为自变量来进行回归分析，得出回归的评分表。通过误差分析来分析以理化指标得出的评价和已给出的评价之间的误差，来判断是否能以葡萄酒和酿酒葡萄的理化指标作为葡萄酒质量的评价标准。

**关键字：**K-S 检验聚类分析主成分分析 K-means++ 分类模型相关性分析多元回归 Pearson 系数

# 1 问题背景与重述

## 1.1 问题背景

葡萄酒是当今世界上最畅销的酒类之一，在各种场合都有葡萄酒的身影。然而葡萄酒的酿造是取决于多种因素，各种因素的叠加会导致葡萄酒的品质差异明显。在不同的原材料已经酿制方法的差别下葡萄酒会继续细分，例如红葡萄酒和白葡萄酒等。对各种葡萄酒的鉴别是必不可少的一个步骤，而采用人工品尝打分和采用仪器进行理化指标的检验已成为最为科学的鉴别方法。最后经过安全检查、筛选分级的葡萄酒方可上市成为饮用酒。

## 1.2 题目所给信息及参数

此次比赛是根据 10 位品酒员为 27 款红酒和 28 款白葡萄酒的打分，以及上述葡萄酒的指标情况和芳香物质为基础进行数学分析和建模，并探讨品酒员的打分是否合理以及论证是否科学分级。红葡萄酒和白葡萄酒的市场在国际上的价值非常之高，葡萄酒依旧是未来的主力酒类，对此分析依旧存在价值。现在根据三个数据文件，并对三个数据进行分析处理后描述统计，完成数学建模和预测。

数据一：葡萄酒品尝评分表；数据二：指标总表；数据三：芳香物质。

## 1.3 所需解决的问题

1. 根据附件所提供的两组品酒员对 27 款红葡萄酒和 28 款白葡萄酒的打分判断两组结果是否有显著性差异，并判断哪一组的更加可信。
2. 根据酿酒葡萄的理化指标和葡萄酒的质量，使用无监督方法计算相似度，通过相似度进行分级。
3. 分析酿酒葡萄和葡萄酒的理化指标之间是否具有相关性，以及其之间具有什么样的联系。
4. 通过主成分分析对数据进行降维，得到贡献值大的特征进行表示，建立其函数关系，根据建立的函数进行预测，将预测结果与评分标准进行误差分析，判断建立的模型是否合理。

# 2 问题分析

针对不同的国家，地区和相对应的医疗水平进行对应的数据指标分析。主要分析感染率，病人接触率，治愈率，以及传染期接触数。模型构建还需要考虑到新冠肺炎的无症状感染者这一特殊的情况，根据这些指标进行相轨线分析，合理的进行疫情的分析 and 未来疫情走向以及各地区、国家的防疫政策研究。

## 2.1 问题一的分析

第一，根据附件 1 中给出两组品酒员的打分情况判断两组的打分是否有显著性差异。对于此问题，分析附件一所提供的数据，研究发现两组对红葡萄酒和白葡萄酒的打分情况是两两相互比较和配对，适合于先进行单样本 K-S 检验判断数据是否符合正态分布，再进行两配对样本 T 检验进行显著性差异判断的办法。

第二，判断两组品酒员的打分情况哪一组更加可信。对于此问题，分析两组品酒员的打分情况，检验两组中打分的更稳定的一方，越稳定的分数即代表品酒员偏好更少，更加可信。提取两组品酒员对于白葡萄酒和红葡萄酒的分数的标准差，根据标准差的大小进行可信性的判断。

## 2.2 问题二的分析

基于改进的 K-means++[?] 进行分类模型, 为了降低数据数, 首先对红、白葡萄和葡萄酒理化指标采用主成分分析法提取出主成分, 但是经过 Bartlett 球形度检验 [?] 等发现不适合进行主成分分析, 进行标准化处理, 去除极值, 使评价分数更加客观, 然后借鉴 Robert Parker 葡萄酒评分体系 [?], 对这些主成分聚类分析得出 6 种聚类并依据判别标准 (聚类后葡萄酒样本的平均值), 最终确定红、白葡萄的分级。

## 2.3 问题三的分析

首先要参照附件所给的数据来进行分析酿酒葡萄和葡萄酒之间的相关性, 附件的信息内容过多, 需要进行合理的过滤和筛选数据, 但要尽可能保证其数据的完整性和真实性。我们采取相关性分析, 依据相关性皮尔森系数来判断葡萄酒的数据和酿酒葡萄的理化指标之间的相关性显著程度。

在进行了相关性分析后, 可以确定下一些具有显著相关的数据流, 要进一步解决其葡萄酒某指标与这些理化指标的关系, 则要进行其关系的拟合, 从而得出实质性的结论来判断酿酒葡萄和葡萄酒之间的关系, 以及其关系的可靠性。

## 2.4 问题四的分析

# 3 符号说明

符号	基本说明
$S(t)$	表示 $t$ 时刻 易感人群 的总人数
$I(t)$	表示 $t$ 时刻 感染人数 的总人数
$R(t)$	表示 $t$ 时刻 退出者 (治愈 + 死亡) 的总人数
$Y(t)$	表示 $t$ 时刻 疑似者 (实际被感染 + 实际未被感染) 的总人数
$N(t)$	表示 $t$ 时刻 所有未隔离患者 的总人数
$\kappa$	疑似人群中被确定未感染人数占疑似人群总数的比例
$\lambda$	隔离者被确诊人数占隔离者总人数的比例
$\theta$	被隔离的人数占未隔离总人数的比例
$\omega$	被确诊且隔离人数占未隔离总人数的比例
$\rho$	感染者平均每天对任何状态的人的接触率
$\xi$	退出率 (死亡率 + 治愈率)

# 4 模型假设

- 1) 考虑到目前已处于疫情控制阶段, 且人们自我隔离意识较好, 不妨设每个病人的有效日接触率为定值
- 2) 所有人口都为易感染者, 不考虑个别免疫体质
- 3) 疑似病例一旦确诊即被隔离, 不会再传播给他人

- 4) 隔离人群中一旦被确定未感染，则立刻结束隔离，即恢复易感者身份
- 5) 在考虑隔离者与未隔离者时将确诊感染者和潜伏期患者都定义为感染者

## 5 模型建立与求解

### 5.1 问题二的求解

题目要求“参考问题 1”，确定至少选择多少家供应商供应原材料才可能满足生产的需求。因此，本文首先将供应商压缩为排名前 50 的供应商，继而针对这 50 家供应商，建立 0-1 整数规划模型，以满足生产需要的最少的供应商数目为目标函数，从而求解出选用哪几家供应商供应原材料，并且这些供应商的总数为多少。

#### 5.1.1 目标函数的构建

为尽可能选择少的供应商以满足生产的需要，本文从 402 家供应商中选出供货能力排名前 50 家供应商作为研究对象。所以，本文构建的目标函数为：

$$\sum_{i=1}^{50} x_{ij} \quad (1)$$

其中  $x_{ij}$  代表第  $i$  个供应商在第  $j$  周是否向该企业供货。若  $x_{ij}$  为 1 则代表供货，反之则代表没有供货。

#### 5.1.2 约束条件的构建

##### 1. 稳定供货能力

为了求解满足企业生产需要的最少的供应商数量，本文当中假设每个供应商的供货量达到自己的最大供货水平。但是从实际情况来看，对于一家供应商来说，24 周每周向企业提供的供应量不可能都是历史最大供货量，供货量也可能有起伏。因此稳定的供货能力不能简单的视为最大的供货量，实际中稳定供货能力应小于或等于历史最大供货量。定义供货系数为  $B, B \in (0, 1)$ ，使得：

$$\text{稳定的供货能力 } S_{ij} = B \cdot \text{最大供货能力} \quad (2)$$

##### 2. 供给需求关系

由于在第  $j$  周所有供货商的供货经过转运到达企业后应满足企业本周的生产需要和库存需要，因此存在供给需求的等式关系，即该生产企业在第  $j$  周实际的原材料接收量与第  $j-1$  周企业仓库中剩余的原材料库存量之和必须大于等于企业本周的产能所需原材料量和两周生产需求的原材料库存量。为了便于计算，我们将等式两边的原材料接收量与原材料库存量均转化为对应的企业产能，则供给需求关系式如下：

$$\sum_{i=1}^{50} x_{ij} S_{ij} (1 - L_{ij}) f_{\text{type}}(i) + R_{j-1} = P + S_{2w} \quad (3)$$

其中， $f_{\text{type}}(i)$  为对原材料加工的转换函数，是计算将不同种类的原材料转换成成品的函数。 $S_{ij}$  表示第  $i$  个供货商在第  $j$  周给该企业的供应量； $L_{ij}$  表示第  $i$  个供货商在第  $j$  周给该企业的供货过程中的损耗率； $R_j$  表示在第  $j$  周后剩余库存转换为相应产能的产量， $P$  表示企业的每周产能， $S_{2w}$  表示满足企业两周用量的原材料库存数量。

3. 损耗率在考虑第  $i$  个供货商在第  $j$  周发出的供货量在转运过程中的损耗率  $L_{ij}$  时，由于无法确定转运公司的选择，所以采用附件二中关于转运商 240 周的损耗率的平均值为损耗率，统计出以 24 周为周期每周的平均损耗率，计算公式如下：

$$L_{ij} = L = \frac{\sum_{k=1}^8 L_{jk}}{8} \quad (4)$$

这种情况下第  $i$  个供货和第 1 个在第  $j$  周的损耗率相同，则：

$$L_{ij} = L_{1j} = \frac{\sum_{k=1}^8 L_{jk}}{8} \quad (5)$$

#### 4. 库存转换为产能

根据题意，该企业要尽可能保持不少于两周生产需求的原材料库存量，因此，有关系式如下：

$$S_{2w} = 2P \quad (6)$$

其中  $P$  为库存， $S_{2w}$  为产能。

#### 5. 库存剩余量的计算迭代公式

第  $j$  周的原材料库存剩余量  $R_j$ ，应等于第  $j$  周的原材料接收量与第  $j-1$  周的原材料库存剩余量之和再减去第  $j$  周的产能所耗的原材料。本文将等式两边的原材料接收量与原材料库存剩余量均转化为相应的产能，所以库存剩余量的迭代关系式如下：

$$\begin{cases} R_j = \sum_{i=1}^{50} x_{ij} S_{ij} (1 - L_{ij}) f_{type}(i) + R_{j-1} - P \\ R_{j-1} = \sum_{i=1}^{50} x_{i(j-1)} S_{i(j-1)} (1 - L_{i(j-1)}) f_{type}(i) + R_{j-2} - P \\ \cdot \\ \cdot \\ R_1 = \sum_{i=1}^{50} x_{i1} S_{i1} (1 - L_{i1}) f_{type}(i) + R_0 - P \end{cases} \quad (7)$$

#### 6. 初始库存量

在知道库存剩余量的迭代关系后，确定初始的库存剩余量  $R_0$  至关重要。考虑到企业依旧持有库存，所以在上一个为期 24 周的生产结束时企业不会将仓库中所有的库存量全部生产用尽，因此不能简单的认为  $R_0 = 0$ 。

根据图 6 的订货量曲线不难发现订货量大致呈现随周期变化的规律，且周期长度为 24 周。所以在以 24 周为周期的订购计划的开始，初始库存剩余量大致可用以往周期的初始库存剩余量来计算。因此不妨设  $R_0 = S_{2w}$ 。

##### 5.1.3 基于 0-1 整数规划的最少供应商模型

在满足企业的生产需求的情况下，尽可能用最少的供应商向该企业提供生产的原材料，结合已知条件和假设，基于 0-1 整数规划的最少供应商模型确立为：

$$\sum_{i=1}^{50} x_{ij} \quad (8)$$

$$\text{s.t.} \begin{cases} \sum_{i=1}^{50} x_{ij} S_{ij} (1 - L_{ij}) f_{type}(i) + R_{j-1} = P + S_{2w} \\ R_{j-1} = \sum_{i=1}^{50} x_{i(j-1)} S_{i(j-1)} (1 - L_{i(j-1)}) f_{type}(i) + R_{j-2} - P \\ L_{ij} = \bar{L} \\ S_{ij} = B \cdot \text{最大供货能力} \\ R_0 = S_{2w} \\ S_{2w} = 2P \end{cases}$$

需要注意: 由于迭代, 该 0-1 整数规划模型对于每个  $j$  (即周) 都是一个全新的规划。

#### 5.1.4 问题二最少供应商模型求解

在满足生产需求的情况下的得到最少的供应商数为 24 个, 结果如下:

S229	S308	S356	S247
S361	S282	S268	S284
S151	S340	S306	S055
S108	S275	S194	S201
S330	S329	S143	S003
S229	S131	S352	S037

## 5.2 基于线性规划的最经济订购方案模型的建立

为了制定未来 24 周每周“最经济”的原材料订购方案, 本文针对上述模型求解出的 24 家原材料供应商, 建立线性规划模型, 在必须满足企业的生产需求的约束条件下, 以 24 周该生产企业需要给原材料供应商的采购费用和给物流公司 (即转运商) 的运输费用以及在仓库储存费用之和达到最小作为目标函数, 从而确定出 24 家原材料供应商每周接到的订货量 (总共 24 周)。

### 5.2.1 目标函数的构建

由于制定订购方案需要明确该生产企业需要订购的原材料供应商以及向这些供应商每周 (总共 24 周) 订购的订货量, 因此我们选择以第  $i$  家供应商在第  $j$  周接到的订货量为自变量, 即  $O_{ij}$ 。根据题目要求, 本文此处的供应商已经确定为上述 0-1 整数规划模型所求解得到的 24 家供应商, 因此  $i=1,2,\dots,24, j=1,2,\dots,24$ 。以 240 周内所有供应商提供的总供货量在采购和运输以及储存该三方面花费的成本最小构造目标函数, 表达式如下:

$$\min \sum_{i=1}^{24} \sum_{j=1}^{24} O_{ij} (1 - r_{os_{ij}}) (r_i p + c) \quad (9)$$

### 5.2.2 约束条件的构建

由于订购方案是为了满足企业的正常生产需求, 因此本文此处的约束条件与基于 0-1 整数规划的最少供应商模型相同, 只是变量的表示发生了变化, 因此此处不再多余赘述约束条件的构建。

### 5.2.3 基于线性规划的最经济订购方案模型

为了制定未来 24 周每周最经济的订购方案, 即要求订购的原材料总量在采购、运输及储存三方面的成本花费为最小, 在制定方案的同时也需要保证该订购方案满足企业的生产需求, 结合已知条件和假设, 该规划模型需要满足以下约束:

1. 该生产企业对一周的原材料接收量和上一周剩余的原材料库存量的总和不少于满足两周生产需要的原材料库存量与一周的产能需要的原材料的总和。
  2. 一周的原材料库存量是来源于该周的原材料接收量与上一周剩余的原材料库存量的总和中除去该周的产能需要的原材料量后剩余的数量。
  3. 该生产企业第一周的库存量为两倍的产能所需原材料。
  4. 该企业尽可能保持不少于满足两周生产需求的原材料库存量。
  5. 为了计算方便，本文将原材料的接收量和库存量转换为相应的等量产能。
- 因此，基于线性规划的最经济订购方案模型确立为：

$$\min \sum_{i=1}^{24} \sum_{j=1}^{24} O_{ij}(1 - r_{OS_{ij}})(r_1 r_1 + r_2) \quad (10)$$

$$\text{s.t.} \begin{cases} \sum_{i=1}^{24} x_{ij} O_{ij}(1 - r_{OS_{ij}}(1 - L_{ij}))f_{type}(i) + R_{i-1} = P + S_{2w} \\ R_{j-1} = \sum_{i=1}^{50} x_{i(j-1)} O_{i(j-1)}(1 - r_{OS_{ij}}(1 - L_{i(j-1)}))f_{type}(i) + R_{j-2} - P_i \\ L_{ij} = \bar{L} \\ r_{OS_{ij}} = r_{OS_{ij}} \\ R_0 = S_{2w} \\ S_{2w} = 2P \end{cases}$$

其中， $r_{OS_{ij}}$  表示以 240 周为时间总长，24 周为一个单位时间长度得到的第  $i$  家供应商在第  $j$  周的平均订购偏差率， $r_1$  是第  $i$  个供应商供应的原材料的采购单价相对 C 类原材料采购单价的比值，其可能取值为  $\{1.2, 1.1, 1\}$ ， $P$  表示 C 类原材料的采购单价， $C$  表示原材料运输和储存的单位费用。

#### 5.2.4 基于线性规划的最经济订购方案模型的求解

##### 方案实施效果分析：

根据本文模型建立的原理，每周 8 家转运商总转运量应大致等于该周向 24 家供应商订货总数，输出结果显示求解结果与预期吻合。此外，可知 24 周平均每立方米原材料需单位成本如下，可见订购方案较为经济。

1	2	3	4	5	6	7	8	9	10	11	12
1.60	1.61	1.57	1.57	1.62	1.56	1.57	1.58	1.58	1.62	1.52	1.55
13	14	15	16	17	18	19	20	21	22	23	24
1.65	1.65	1.57	1.57	1.65	1.57	1.57	1.65	1.57	1.57	1.65	1.60

### 5.3 基于线性规划的最少转运损耗方案模型的建立

根据上述基于线性规划的最经济订购方案模型，可以得到 24 家供应商每周接到的订购量，据此来制定每周内每家供应商的供货量分配给哪一家/几家转运商来运输，并以每家转运商的运输能力为 6000 立方米/周作为约束条件之一，建立线性规划模型，以 8 家转运商在 24 周的转运过程中的总损耗量最少为目标函数，从而确定出每周每家转运商转运哪一家/几家供货商的供货量。

由于过去的转运损耗不会影响本周的转运损耗，所以要使转运损耗最小只需要使每周的转运损耗都最小，为了简便计算可以一周为单位考虑。假设第  $j(j=1, 2, \dots, 24)$  周第  $i(i=1, 2, \dots, 24)$  家供货商由第  $k(k=1, 2, \dots, 8)$  家转运公司转运的原材料量为  $T_{ik}$ ，该转运公司的转运损耗率为  $L_k$ ，则第  $j$  周所有货

物的总损耗为  $\sum_i^{24} \sum_{k=1}^8 \mathbf{T}_{ik} \mathbf{L}_k$ ，所以为了使得转运损耗最小，有目标函数：

$$\min \sum_{i=1}^{24} \sum_{k=1}^8 \mathbf{T}_{ik} \mathbf{L}_k, \text{ for } j = 1, 2, \dots, 24 \quad (11)$$

### 5.3.1 约束条件的建立

#### 1. 转运损失率

由于同一个转运公司不同时间转运的损耗率不同，且由上文可知，整个订购运送过程具有周期性，近 5 年共有 10 个周期，周期为 24 周，故转运公司  $k$  在第  $j$  周的转运损失率可以用以往 10 个周期中第  $j$  周转运损失率的均值  $\mathbf{L}(k, j)$  表示，计算关系式如下：

$$\mathbf{L}_k = \mathbf{L}(k, j) \quad (12)$$

#### 2. 承接量

根据题意每家转运商的运输能力为 6000 立方米/周，所以第  $k$  家转运公司在 24 周内给第  $i$  家企业转运原材料的总量不大于 6000 (单位：立方米)，不等式如下：

$$\sum_{i=1}^{24} \mathbf{T}_{ik} \leq 6000, \text{ for } k = 1, 2, \dots, 8 \quad (13)$$

#### 3. 供货量

若第  $i$  家供货商由第  $k$  家转运公司转运的原材料量为  $\mathbf{T}_{ik}$ ，则该供货商由各个转运商转运的原材料量的总和应该为该供货商的总供货量，即：

$$\sum_{k=1}^8 \mathbf{T}_{ik} = \mathbf{S}_i, \text{ for } i = 1, 2, \dots, 24 \quad (14)$$

其中， $\mathbf{S}_i$  为第  $i$  家供货商在第  $j$  天的给定的总供货量。

### 5.3.2 基于线性规划的最经济转运方案模型

基于线性规划的最经济转运方案模型确立为：

$$\begin{aligned} & \min \sum_{i=1}^{24} \sum_{k=1}^8 T_{ik} L_k \\ \text{s.t. } & \begin{cases} L_k = \bar{L}(k, j) \\ \sum_{i=1}^{24} T_{ik} \leq 6000, \text{ for } k = 1, 2, \dots, 8 \\ \sum_{k=1}^8 T_{ik} = S_i, \text{ for } i = 1, 2, \dots, 24 \end{cases} \end{aligned} \quad (15)$$

该线性规划模型只针对求解一周的转运方案，为了获得 24 周的转运方案，则需要重复使用该模型，其中  $\mathbf{T}_{ik}, \mathbf{L}_k, \mathbf{S}_i$  都会随着周数的改变而改变。但因为两个约束条件分别是关于 24 周每周的约束条件和 8 家转运商每家的约束条件，所以该模型的可解性是一定被保证的。

## 5.4 基于线性规划的最少损耗转运方案模型的求解

方案实施效果分析：

24 周平均每周转运损耗率如下，可见转运方案转运效率良好：



1	2	3	4	5	6
0.1397	0.2122	0.1717	0.1859	0.2003	0.2567
7	8	9	10	11	12
0.1189	0.0280	0.1960	0.2439	0.1233	0.2115
13	14	15	16	17	18
0.2936	0.3483	0.3442	0.1412	0.4659	0.4881
19	20	21	22	23	24
0.1621	0.1621	0.1621	0.1621	0.1621	0.1621

```

1     import time
2     import numpy as np
3     import pandas as pd
4     import matplotlib.pyplot as plt

```