



長春工業大學

数据科学导引

题 目	基于 PCA 主成分分析法对葡萄酒数据进行分析
--------	-------------------------

学 院	数学与统计学院
--------	---------

专 业	应用统计学
--------	-------

学 号	2202312021
--------	------------

姓 名	杨玲玲
--------	-----

2024 年 5 月 31 日

目录

摘要	4
Abstract	5
一 引言	7
二 理论	10
2.1 算法理论	10
2.1.1 主成分分析简介	10
2.1.2 主成分分析的思想	11
2.1.3 主成分分析的计算步骤	11
三 实证分析	14
3.1 数据集介绍	14
3.2 数据预处理	15
3.2.1 导入必要的库函数	15
3.2.2 读取数据	16
3.2.3 查看数据并检查否有缺失值	16
3.3 PCA 主成分分析	18
3.3.1 巴特利球形检验	18
3.3.2 求相关矩阵	18
3.3.3 绘制散点图和折线图	20
3.3.4 求主成分得分	20
3.3.5 绘制热力图	21
3.4 构建模型	22
3.4.1 主成分分析数据	22
3.4.2 原始数据	23
四 结论	24

参考文献.....	25
附录.....	26

摘要

葡萄酒，作为人类历史上最早被酿造和享用的饮料酒之一，其独特的口感和深厚的文化底蕴早已深入人心。自古以来，葡萄酒不仅仅是一种饮品，更是文化交流的媒介，是节日庆典的必备佳品，也是品味生活的象征^[1]。然而，随着全球气候变化、市场竞争加剧以及消费者口味的多样化，葡萄酒产业也面临着前所未有的挑战和机遇。据最新的统计数据显示，2021 年全球葡萄酒的产量达到了 260 亿升，然而这一数字较 2020 年下降了近 1%，这标志着全球葡萄酒产量已经连续三年略低于过去十年的平均水平。这一变化不仅反映了全球葡萄酒市场的供需格局正在发生变化，也暗示着葡萄酒生产商需要更加注重品质和创新，以应对市场的挑战。

在我国，葡萄酒产业同样面临着巨大的压力和挑战。长期以来，我国本土葡萄酒受到进口葡萄酒的强劲冲击，市场份额不断被挤压。尤其是在 2020 年疫情期间，由于节日聚会、家庭餐会等社交活动的取消，餐饮业几乎完全停滞，导致葡萄酒的需求在短期内大幅下降。这一变化对我国葡萄酒行业产生了深远的影响，市场规模在 2020 年下滑至 498.2 亿元。然而，在经历了一段时间的低迷之后，我国葡萄酒行业在 2021 年开始逐渐回暖。随着疫情得到控制，消费者的信心逐渐恢复，葡萄酒市场也开始展现出新的生机。据统计，2021 年我国葡萄酒行业市场规模小幅度上涨至 510.8 亿元。这一变化表明，我国葡萄酒行业正在逐渐走出低谷，迎来新的发展机遇。展望未来，预计 2022 年我国葡萄酒行业市场规模将进一步上升至 587.2 亿元。这一预测不仅反映了我国葡萄酒市场的巨大潜力，也预示着我国葡萄酒产业将迎来更加广阔的发展前景。

为了更好地把握这一发展机遇，本次实验将使用葡萄酒数据集进行主成分分析。通过主成分分析，我们可以对葡萄酒数据集进行降维处理，提取出最具代表性的特征变量^[2]。然后，我们将利用这些特征变量构建葡萄酒分类模型。该模型将基于机器学习算法，通过训练数据学习葡萄酒的品质特征和分类规则，实现对未知葡萄酒样本的准确分类。这一模型不仅可以帮助我们更好地了解葡萄酒的品质特点，还可以为葡萄酒生产和销售提供科学依据，推动葡萄酒产业的可持续发展。

关键词：葡萄酒、主成分分析、机器学习算法、分类模型

Abstract

Wine, as one of the earliest beverage wines brewed and enjoyed in human history, its unique taste and profound cultural heritage have long been deeply ingrained in people's hearts. Since ancient times, wine has not only been a beverage, but also a medium for cultural exchange, an essential delicacy for festivals and celebrations, and a symbol of savoring life. However, with global climate change, intensified market competition, and the diversification of consumer tastes, the wine industry is also facing unprecedented challenges and opportunities. According to the latest statistical data, the global wine production reached 26 billion liters in 2021, but this number has decreased by nearly 1% compared to 2020, indicating that global wine production has been slightly lower than the average level of the past decade for three consecutive years. This change not only reflects the changing supply and demand pattern of the global wine market, but also implies that wine producers need to pay more attention to quality and innovation to meet market challenges.

In China, the wine industry is also facing enormous pressure and challenges. For a long time, domestic wines in China have been strongly impacted by imported wines, and their market share has been constantly squeezed. Especially during the COVID-19 pandemic in 2020, due to the cancellation of social activities such as holiday gatherings and family dinners, the catering industry almost came to a complete halt, resulting in a significant decrease in demand for wine in the short term. This change has had a profound impact on China's wine industry, with the market size declining to 49.82 billion yuan in 2020. However, after experiencing a period of downturn, China's wine industry began to gradually recover in 2021. With the epidemic under control, consumer confidence is gradually recovering, and the wine market is also beginning to show new vitality. According to statistics, the market size of China's wine industry increased slightly to 51.08 billion yuan in 2021. This change indicates that China's wine industry is gradually emerging from its lows and ushering in new development opportunities. Looking ahead, it is expected that the market size of China's wine industry will further increase to 58.72 billion yuan in 2022. This prediction not only reflects the enormous potential of China's wine market, but also foreshadows a broader development

prospect for China's wine industry.

In order to better seize this development opportunity, this experiment will use the wine dataset for principal component analysis. Through principal component analysis, we can perform dimensionality reduction on the wine dataset and extract the most representative feature variables. Then, we will use these feature variables to construct a wine classification model. This model will be based on machine learning algorithms and learn the quality features and classification rules of wine through training data, achieving accurate classification of unknown wine samples. This model can not only help us better understand the quality characteristics of wine, but also provide scientific basis for wine production and sales, and promote the sustainable development of the wine industry.

Keywords: wine, principal component analysis, machine learning algorithms, classification models

一 引言

葡萄酒，作为人类历史上最早的饮料酒之一，其深厚的文化底蕴与西方文明的演进紧密相连，仿佛一部流淌着岁月痕迹的史诗。从古老的酿造技艺到现代的科技革新，葡萄酒的生产历程见证了人类文明的不断进步。根据酿造历史和独特的生产工艺，葡萄酒生产国被清晰地划分为两大阵营：一方是以传统酿造工艺为基石的“旧世界”，这里汇聚了法国、意大利、西班牙、德国等欧洲传统葡萄酒大国，它们承载着数百年的酿酒传统和精湛的技艺；另一方则是“新世界”，它们以现代酿造技术为先导，包括了美国、澳大利亚、新西兰、智利、阿根廷和南非^[3]等国家，这些新兴的葡萄酒生产国以其创新的精神和独特的风格，为全球葡萄酒市场注入了新的活力。

当我们把目光转向 2021 年的全球葡萄酒产量数据时，不禁为这一年的产量变化而感叹。据统计，全球葡萄酒产量达到了 260 亿升，但较 2020 年却下降了近 1%，这已经是连续三年略低于 10 年平均水平。在这其中，欧盟的葡萄酒产量尤为引人注目。受霜冻天气的影响，欧盟的葡萄酒产量仅为 153.7 亿升，较 2020 年下降了 8%。这一数据的背后，是法国、意大利、西班牙等欧洲主要葡萄酒生产国产量的波动。尤其是法国，作为葡萄酒界的佼佼者，其产量的大幅下降成为影响欧盟葡萄酒产量的关键因素。2021 年 4 月份的霜冻天气给法国的葡萄园带来了严重的打击，许多葡萄树遭受了不同程度的冻害，导致葡萄产量锐减。尽管意大利的葡萄酒产量有所增长，达到了 50.2 亿升，较 2020 年增长了 2%，但法国和西班牙的产量却分别下降了 19% 和 14%，分别达到了 37.6 亿升和 35.3 亿升。这三个国家合计占全球葡萄酒产量的 47%，它们的产量波动无疑对全球葡萄酒市场产生了深远的影响。在这一连串的数字背后，我们看到了葡萄酒产业所面临的挑战与机遇。面对气候变化的威胁和消费者需求的不断变化，葡萄酒生产国需要不断创新和进步，以保持其在全球市场的竞争力。同时，我们也期待着这些古老的葡萄酒生产国能够继续传承和发扬其独特的酿酒文化，为全球消费者带来更多优质的葡萄酒产品。

在我国葡萄酒行业的发展历程中，本土品牌长期以来都面临着来自进口葡萄酒的激烈竞争。尤其是近年来，随着国际葡萄酒市场的日益开放和进口关税的逐步降低，进口葡萄酒以其多样化的品种、高品质的口感和独特的文化魅力，不断冲击着我国本土葡萄酒的市场地位。

与此同时，本土葡萄酒行业还需应对一系列的内外挑战。而 2020 年，突如其来的新冠疫情给我国葡萄酒行业带来了前所未有的冲击。疫情期间，节日聚会、家庭餐会等社交活动纷纷被取消，餐饮业也几乎完全停滞，导致葡萄酒的消费需求在短期内出现了大幅度的下滑。这对于本就处于困境中的本土葡萄酒行业来说，无疑是雪上加霜。

据行业统计数据显示，2020 年我国葡萄酒行业的市场规模遭受了严重的影响，下滑至 498.2 亿元，这一数字较往年有了明显的降低。然而，在经历了艰难的疫情后，随着国内经济的逐渐复苏和消费者信心的逐步恢复，葡萄酒行业也开始呈现出复苏的迹象。到了 2021 年，我国葡萄酒行业的市场规模出现了小幅度上涨，达到了 510.8 亿元。这一增长虽然不算显著，但已经足以说明葡萄酒行业正在逐步走出低谷，迎来新的发展机遇。

展望未来，随着国内消费者对于葡萄酒文化认知的不断提升和葡萄酒消费市场的不断扩大^[4]，预计 2022 年我国葡萄酒行业的市场规模将进一步上升至 587.2 亿元。这将为本土葡萄酒品牌提供更多的发展空间和机会，同时也将促进整个行业的持续健康发展。在这个过程中，本土葡萄酒品牌需要不断提升自身的品质和服务水平，以更好地满足消费者的需求，赢得市场的认可。

葡萄酒行业的产业链，自上游的精心培育与准备，至中游的精湛酿造技艺，再到下游的多样化消费渠道，构成了一个完整且充满活力的生态体系^[5]。在上游，产业链的核心参与主体是提供原材料和相关设备的供应商。这里的原材料主要包括优质的葡萄、淀粉和酵母，它们为葡萄酒的酿造提供了最基础的物质条件。此外，食品添加剂和包装材料也扮演着不可或缺的角色，它们确保了葡萄酒的品质和外观的吸引力。同时，先进的酿造设备则为葡萄酒的生产提供了强大的技术支持，确保了整个酿造过程的顺利进行。进入中游，我们迎来了葡萄酒行业的核心——葡萄酒的酿造。在这里，精湛的酿造技艺和严格的品质控制是每一家葡萄酒企业赖以生存和发展的基石。葡萄酒按照颜色和风味可以细分为白葡萄酒、红葡萄酒和桃红葡萄酒三大类，它们各自拥有独特的口感和风味，满足了不同消费者的需求。在产业链的下游，葡萄酒通过各类消费渠道流向终端消费者。这些渠道包括传统的零售商店、酒店、餐厅等，也包括新兴的电商平台和社交媒体。随着消费者对葡萄酒文化的了解和喜爱不断加深，葡萄酒的消费场景也在不断扩大，从商务宴请到家庭聚会，从节日庆典到日常休闲，葡萄酒

都成为了不可或缺饮品。

近年来，随着葡萄酒产业链的延伸和拓展，一些新型的“葡萄酒+”产业模式也应运而生。这些模式以葡萄酒为核心，结合其他产业元素，创造出更加多元化和个性化的消费体验。比如，“葡萄酒+旅游”模式，通过葡萄酒庄园的观光游览、品鉴体验等活动，让消费者在享受美酒的同时，也能领略到葡萄酒产区的自然风光和人文风情。而“葡萄酒+科技”模式，则利用现代科技手段，如大数据、人工智能等，为葡萄酒的生产、销售和服务提供更加智能化和个性化的解决方案，进一步提升了葡萄酒行业的竞争力。这些新型模式的出现，不仅丰富了葡萄酒行业的内涵和外延，也为整个产业链的发展注入了新的活力和动力。

随着我国经济社会的飞速发展，全国范围内已经全面步入小康社会，这一历史性的跨越极大地提升了人民群众的生活品质，也推动了消费市场的持续升级和转型。在这一过程中，具有独特品牌魅力和一定溢价能力的中高端葡萄酒，无疑成为了消费市场的新宠，其发展前景广阔且充满潜力。

经过国内葡萄酒企业多年的不懈努力和精心耕耘，中国葡萄酒的品质与风味已经得到了广大消费者的普遍认可。这不仅体现在产品的口感和风味上，更体现在其背后所承载的文化内涵和品牌故事上。这些优质的中国葡萄酒，以其独特的魅力，逐渐在国际葡萄酒市场上崭露头角，赢得了世界的赞誉。与此同时，中国消费者对于葡萄酒的认知和欣赏水平也在不断提高。他们不再仅仅满足于基本的饮用需求，而是更加注重葡萄酒的品质、口感和背后的文化故事。这种消费观念的转变，为中国葡萄酒搏击中高端市场提供了坚实的市场基础。

因此，我们可以预见，在 2022 年及未来更长的时间里，中高端葡萄酒将成为中国葡萄酒产业的重要发展热点。随着国内葡萄酒企业不断提升产品品质、丰富产品种类、加强品牌建设，以及国内外消费市场的不断升级和扩大，中国葡萄酒将在中高端市场上迎来更加广阔的发展空间。这将是中国葡萄酒产业迈向更高水平、实现更大突破的重要机遇。

二 理论

2.1 算法理论

2.1.1 主成分分析简介

主成分分析,作为一种高效的数据降维技术,在现代数据分析中扮演着举足轻重的角色。随着数据量的不断膨胀和维度的急剧增加,传统的数据处理方法往往显得力不从心,这时候就需要一种能够在保证信息完整性的同时,有效减少数据复杂性的技术。而主成分分析正是这样一种技术,它为我们提供了一种将高维度数据转化为低维度表示的方法,极大地提升了数据处理的速度和效率。具体来说,主成分分析基于降维的核心思想,旨在通过科学的数学方法,将原本众多的指标转化为少数几个综合指标。这些综合指标,我们通常称之为“主成分”,它们不仅是原始变量的线性组合,更是原始变量信息的凝练和集中。值得注意的是,这些主成分之间互不相关,这意味着它们各自独立地承载了原始数据中的不同信息,使得我们能够更加清晰地看到数据背后的规律和模式^[6]。

主成分分析之所以具有如此强大的功能,是因为它在转化过程中能够最大程度地保留原始数据的信息。虽然经过降维处理,数据的维度减少了,但主成分所包含的信息量却足以解释原始数据的绝大多数信息。这种能力使得我们在研究复杂问题时,只需关注少数几个主成分,就能够抓住问题的主要矛盾,揭示事物内部变量之间的规律性。此外,主成分分析还具有简化问题、提高分析效率的优点。通过将高维度的数据降维至低维度,我们不仅能够更直观地理解数据的结构和特征,还能够更加便捷地进行数据分析和模型构建。这不仅降低了数据分析的复杂度,还提高了数据分析的效率和准确性,使得我们能够更加快速地发现数据中的规律和趋势,为决策提供更加有力的支持。

总之,主成分分析是一种功能强大、灵活多变的数据降维技术,它能够帮助我们有效地处理高维度的数据,提升数据处理的速度和效率。无论是在科学研究、商业分析还是其他领域,主成分分析都发挥着不可或缺的作用,为我们提供了一种全新的视角和方法来理解和分析数据。

2.1.2 主成分分析的思想

假设有 n 个样本, p 个指标, 则可构成大小为 $n \times p$ 的样本矩阵 x :

$$X = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1p} \\ x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{np} \end{bmatrix} = (x_1, x_2, \cdots, x_p) \quad (2.1)$$

假设我们想找到新的一组变量 $z_1, z_2, \cdots, z_m (m \leq p)$, 且它们满足:

$$\begin{cases} z_1 = 1_{11}x_1 + 1_{12}x_2 + \cdots + 1_{1p}x_p \\ z_2 = 1_{21}x_1 + 1_{22}x_2 + \cdots + 1_{2p}x_p \\ \cdots \\ z_m = 1_{m1}x_1 + 1_{m2}x_2 + \cdots + 1_{mp}x_p \end{cases} \quad (2.2)$$

系数 1_{ij} 的确定原则:

1. z_i 与 $z_j (i \neq j; i, j = 1, 2, \dots, m)$ 相互无关;
2. z_1 是 x_1, x_2, \dots, x_p 的一切线性组合中方差最大者;
3. z_2 是与 z_1 不相关的 x_1, x_2, \dots, x_p 的所有线性组合中方差最大者;
4. 以此类推, z_m 是与 z_1, z_2, \dots, z_{m-1} 不相关的 x_1, x_2, \dots, x_p 的所有线性组合中方差最大者;
5. 新变量指标 z_1, z_2, \dots, z_m 分别称为原变量指标 x_1, x_2, \dots, x_p 的第一, 第二, ..., 第 m 主成分。

2.1.3 主成分分析的计算步骤

假设有 n 个样本, p 个指标, 则可构成大小为 $n \times p$ 的样本矩阵 x :

$$X = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1p} \\ x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{np} \end{bmatrix} = (x_1, x_2, \cdots, x_p) \quad (2.3)$$

(1) 首先对其进行标准化处理:

按列计算均值:

$$\bar{x}_j = \sum_{i=1}^n x_{ij} \quad (2.4)$$

标准差:

$$S_j = \sqrt{\frac{\sum_{i=1}^n (x_{ij} - \bar{x}_j)^2}{n-1}} \quad (2.5)$$

标准化数据: X_i

原始样本矩阵经过标准化变化：

$$X = \begin{bmatrix} X_{11} & X_{12} & \cdots & X_{1p} \\ X_{21} & X_{22} & \cdots & X_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ X_{n1} & X_{n2} & \cdots & X_{np} \end{bmatrix} = (X_1, X_2, \cdots, X_p) \quad (2.6)$$

(2) 计算标准化样本的协方差矩阵：

$$R = \begin{bmatrix} r_{11} & r_{12} & \cdots & r_{1p} \\ r_{21} & r_{22} & \cdots & r_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ r_{n1} & r_{n2} & \cdots & r_{np} \end{bmatrix} \quad (2.7)$$

$$\text{其中 } r_{ij} = \frac{1}{n-1} \sum_{k=1}^n (X_{ki} - \bar{X}_i)(X_{kj} - \bar{X}_j) = \frac{1}{n-1} \sum_{k=1}^n X_{ki} X_{kj}$$

(1)、(2) 可以合成一步：

$$R = \frac{\sum_{k=1}^n (x_{ki} - \bar{x}_i)(x_{kj} - \bar{x}_j)}{\sqrt{\sum_{k=1}^n (x_{ki} - \bar{x}_i)^2 \sum_{k=1}^n (x_{kj} - \bar{x}_j)^2}} \quad (2.8)$$

(3) 计算 R 的特征值和特征值向量：

特征值 $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_p \geq 0$, (R 是半正定矩阵, 且 $\text{tr}(R) = \sum_{k=1}^p \lambda_k = p$)

特征向量：

$$a_1 = \begin{bmatrix} a_{11} \\ a_{21} \\ \vdots \\ a_{p1} \end{bmatrix}, a_2 = \begin{bmatrix} a_{12} \\ a_{22} \\ \vdots \\ a_{p2} \end{bmatrix}, \cdots, a_p = \begin{bmatrix} a_{1p} \\ a_{2p} \\ \vdots \\ a_{pp} \end{bmatrix} \quad (2.9)$$

(4) 计算主成分贡献率以及累计贡献率：

$$\text{贡献率} = \frac{\lambda_i}{\sum_{k=1}^p \lambda_k}, \text{累加贡献率} = \frac{\sum_{k=1}^i \lambda_k}{\sum_{k=1}^p \lambda_k}, (i=1, 2, \cdots, p) \quad (2.10)$$

(5) 写出主成分：

一般取累计贡献率超过 80% 的特征值所对应的第一、第二、…、第 m ($m \leq p$) 个主

成分。

$$\text{第}i\text{个主成分: } F_i = a_{1i}X_1 + a_{2i}X_2 + \cdots + a_{pi}X_p \quad (2.11)$$

(6) 根据系数分析主成分代表的意义:

对于某个主成分而言, 指标前面的系数越大, 代表该指标对于该主成分的影响越大。

三 实证分析

3.1 数据集介绍

在葡萄酒行业，为了更深入地研究葡萄酒的品质、口感以及化学特性，科学家们常常依赖于各种数据集来进行分析。其中，wine 样本数据集就是一个非常宝贵和常用的资源。这个数据集以矩阵的形式呈现，其大小为 178 行 14 列，每个元素都是 double 类型的数值，精确地记录了三种不同葡萄酒中 13 种关键成分的含量。在这个数据集中，每一行都代表了一个独立的葡萄酒样本，总共包含了 178 个这样的样本。这些样本涵盖了不同品种、不同产地、不同酿造工艺的葡萄酒，为研究者们提供了一个丰富多样的研究对象集。

在矩阵的每一行中，我们都可以看到 14 列的数据。第一列尤为特殊，它是一个类标识符，用于标识葡萄酒的分类。在这个数据集中，有三种不同的葡萄酒分类，分别用数字 1、2、3 来表示。这三个数字不仅仅是一个简单的标识，它们背后代表着葡萄酒的不同口感、风味以及可能的酿造过程。接下来的 13 列则是每个样本中对应属性的具体数值。这些属性包括酒精、苹果酸、灰、灰分的碱度、镁、总酚、黄酮类化合物、非黄烷类酚类、原花色素、颜色强度、色调、稀释葡萄酒的 OD280/OD315 比值以及脯氨酸^[7]。这些化学成分是葡萄酒品质评价的重要依据，它们决定了葡萄酒的口感、香气、色泽等多个方面的特征。

具体来说，酒精含量影响着葡萄酒的酒精度数和口感；苹果酸则影响着葡萄酒的酸度和清爽感；灰和灰分的碱度则反映了葡萄酒中矿物质和无机盐的含量；镁是一种重要的微量元素，对葡萄酒的口感和品质也有一定的影响；总酚、黄酮类化合物、非黄烷类酚类和原花色素则是葡萄酒中重要的抗氧化物质，对葡萄酒的陈年和品质保持具有重要作用；颜色强度和色调则描述了葡萄酒的外观特征；而稀释葡萄酒的 OD280/OD315 比值则反映了葡萄酒中蛋白质和多酚类物质的含量；最后，脯氨酸则是一种重要的氨基酸，对葡萄酒的口感和风味也有一定的影响。

在这个数据集中，三种葡萄酒分类的样本数量也有所不同。第一类的葡萄酒有 59 个样本，可能代表了某种特定口感或风格的葡萄酒；第二类有 71 个样本，数量上占据优势，可能反映了这种葡萄酒在市场上更受欢迎或更易于酿造；第三类则有 48 个样本，虽然数量较少，但同样具有独特的口感和风味特点。通过对这些不同分类的葡萄酒样本进行分析和比较，

我们可以更深入地了解它们之间的异同以及影响它们品质的关键因素。具体属性描述如表 3-1:

表 3-1 数据集属性描述

属性	属性描述
target	类别
Alcohol	酒精
Malic acid	苹果酸
Ash	灰
Alkalinity of ash	灰分的碱度
Magnesium	镁
Total phenoids	总酚
Flavonoids	黄酮类化合物
Noflavanoid phenols	非黄烷类酚类
Proanthocyanins	原花色素
Color intensity	颜色强度
Hue	色调
OD280/0315ofdiluted wines	稀释葡萄酒的 OD280/0315
proline	脯氨酸

3.2 数据预处理

3.2.1 导入必要的库函数

(1) pandas 是一个强大的数据分析工具，提供了 DataFrame 和 Series 两种数据结构，方便处理和分析表格型数据。

(2) numpy 是 Python 中用于处理数组、矩阵和数学函数运算的库，是进行数值计算的基础。

(3) matplotlib 是一个用于绘制图形的库，pyplot 是其子模块，提供了类似于 MATLAB 的绘图 API。

(4) seaborn 是基于 matplotlib 的图形可视化 Python 库，它提供了一种高级界面来绘制有吸引力的和信息丰富的统计图形。

(5) Sklearn 是 Python 中经典的机器学习模块，该模块围绕着机器学习提供了很多可直接调用的机器学习算法以及很多经典的数据集。

3.2.2 读取数据

加载葡萄酒数据集，表 3-2 是葡萄酒数据集前 5 列。

表 3-2 葡萄酒详细数据集

	alcohol	malic_acid	ash	alcalinity_of_ash	magnesium	total_phenols	flavanoids
0	14.23	1.71	2.43	15.6	127.0	2.80	3.06
1	13.20	1.78	2.14	11.2	100.0	2.65	2.76
2	13.16	2.36	2.67	18.6	101.0	2.80	3.24
3	14.37	1.95	2.50	16.8	113.0	3.85	3.49
4	13.24	2.59	2.87	21.0	118.0	2.80	2.69

nonflavanoid_phe	proanthocy	color_i	hue	od280/od315	proline
nols	anins	ntensity		_of_diluted_wines	
0.28	2.29	5.64	1.04	3.92	1065.0
0.26	1.28	4.38	1.05	3.40	1050.0
0.30	2.81	5.68	1.03	3,17	1185.0
0.24	2.18	7.80	0.86	3.45	1480.0
0.39	1.82	4.32	1.04	2.93	735.0

3.2.3 查看数据并检查否有缺失值

在数据科学和机器学习的实践中，缺失值往往是一个难以回避的挑战。它们可能是由于

各种原因在数据收集、传输或处理过程中产生的，比如人为错误、设备故障或特定的数据收集限制。这些缺失值的存在不仅会影响数据的完整性，更关键的是，它们可能会直接或间接地对后续的数据分析和机器学习模型的性能产生负面影响。在构建和训练机器学习模型之前，对缺失值进行妥善处理是至关重要的一步。许多常用的机器学习算法并不能直接处理包含缺失值的数据集，因为这些算法往往需要稳定、一致且完整的输入特征来发挥其最佳效能。因此，若忽略这些缺失值，模型的性能很可能会受到影响，导致预测不准确或模型泛化能力下降。

检查数据集中的缺失值不仅能帮助我们确保数据的完整性，还能为我们揭示数据的潜在问题或限制。通过观察哪些特征列包含缺失值，我们可以了解到哪些部分的数据可能不够完整或存在不确定性，这对于后续的数据预处理、特征选择或模型优化决策都具有重要的参考意义。处理缺失值的方法多种多样，每种方法都有其适用场景和优缺点。例如，简单地删除含有缺失值的行可能会导致数据量的减少，进而影响模型的训练效果^[8]；而使用均值、中位数或众数等统计量进行填充则可能引入数据分布的偏差。因此，在选择缺失值处理方法时，我们需要综合考虑数据集的特性、问题的需求以及模型的要求，选择最适合当前情况的处理策略。

总之，缺失值问题是数据分析和机器学习中的一个重要环节。通过仔细检查和处理缺失值，我们可以确保数据的完整性和准确性，为后续的建模和分析奠定坚实的基础^[9]。同时，这也需要我们在实践中不断探索和学习，以便找到最适合自己数据集和问题的处理方法。

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 178 entries, 0 to 177
Data columns (total 13 columns):
#   Column                                Non-Null Count  Dtype
---  ---                                -
0   alcohol                             178 non-null    float64
1   malic_acid                         178 non-null    float64
2   ash                                178 non-null    float64
3   alcalinity_of_ash                  178 non-null    float64
4   magnesium                          178 non-null    float64
5   total_phenols                      178 non-null    float64
6   flavanoids                         178 non-null    float64
7   nonflavanoid_phenols              178 non-null    float64
8   proanthocyanins                   178 non-null    float64
9   color_intensity                    178 non-null    float64
10  hue                                178 non-null    float64
11  od280/od315_of_diluted_wines      178 non-null    float64
12  proline                            178 non-null    float64
dtypes: float64(13)
memory usage: 18.2 KB
```

图 3-1 葡萄酒基本信息图

从图 3-1 中的结果可以看出，数据集没有缺失值且都为浮点数类型。

3.3 PCA 主成分分析

3.3.1 巴特利球形检验

检验总体变量的相关矩阵是否是单位阵（相关系数矩阵对角线的所有元素均为 1,所有非对角线上的元素均为零）；即检验各个变量是否各自独立。我们进行了巴特利球形检验。

表 3-2 巴特利球形检验

	数值
chi_square_value	1317.18
p_value	2.47e-224

从表 3-2 的结果中看出 P 值远小于 0.05,拒绝原假设，说明变量之间有相关关系，可以做主成分分析。

3.3.2 求相关矩阵

（1）标准化

在进行数据分析之前，通常要收集大量不同的相关指标，每个指标的性质、量纲、数量级、可用性等特征均可能存在差异，导致我们无法直接用其分析研究对象的特征和规律^[10]。当各指标间的水平相差很大时，如果直接用指标原始值进行分析，数值较高的指标在综合分析中的作用就会被放大，相对地，会削弱数值水平较低的指标的作用。

比如，在评价不同时期的物价指数时，较低价格的蔬菜和较高价格的家电的价格涨幅都可以被纳入其中，但是由于它们的价格水平差异较大，如果直接用其价格做分析，会使价格水平较高的家电在综合指数中的作用被放大。因此，为了保证结果的可靠性，需要对原始指标数据进行变换处理，使不同的特征具有相同的尺度^[11]。

标准化指将数据按比例缩放，使之落入一个小的特定区间。在某些比较和评价的指标处理中经常会用到，去除数据的单位限制，将其转化为无量纲的纯数值，便于不同单位或量级

的指标能够进行比较和加权。

(2) 求相关系数矩阵，并求特征值和特征向量，对特征值进行排序并输出。

在统计分析中，相关系数矩阵是一个非常重要的工具，它能够揭示变量之间的线性关系。当我们谈论“求相关系数矩阵”，我们通常指的是计算一组变量之间的相关系数，并将它们排列成一个方阵。这个方阵的每个元素 r_{ij} 表示变量 i 和变量 j 之间的相关系数，其中 i 和 j 都是变量的索引。

相关系数的取值范围是-1 到+1。当相关系数为+1 时，表示两个变量之间存在完全正相关关系；当相关系数为-1 时，表示两个变量之间存在完全负相关关系；而当相关系数接近 0 时，表示两个变量之间没有线性关系。

接下来，我们进行特征值和特征向量的计算。特征值和特征向量是线性代数中的概念，它们与矩阵的对角化密切相关。对于一个给定的方阵，我们可以通过求解线性方程组 $Ax = \lambda x$ 来找到特征值 λ 和相应的特征向量 x 。其中， A 是相关系数矩阵， λ 是特征值， x 是特征向量。

特征值代表了矩阵在某个方向上的拉伸或压缩的程度，而特征向量则表示这个方向。在实际应用中，特征值和特征向量可以帮助我们识别数据中的主要变化方向，例如在主成分分析 (PCA) 中，特征向量对应的是主成分的方向，而特征值则表示每个主成分解释的方差量。

对特征值进行排序，通常是按照它们的大小从大到小进行排列。这样做可以帮助我们识别最重要的特征向量，即那些与数据变化最相关的方向。在 PCA 中，我们通常会保留那些具有较大特征值对应的特征向量，因为它们能够捕捉到数据中的主要变异性。

最后，输出特征值和特征向量，这为我们提供了一个数据集的压缩表示，可以用于降维、数据可视化或进一步的统计分析。通过这种方式，我们可以简化复杂数据集的结构，同时保留最重要的信息。

3.3.3 绘制散点图和折线图

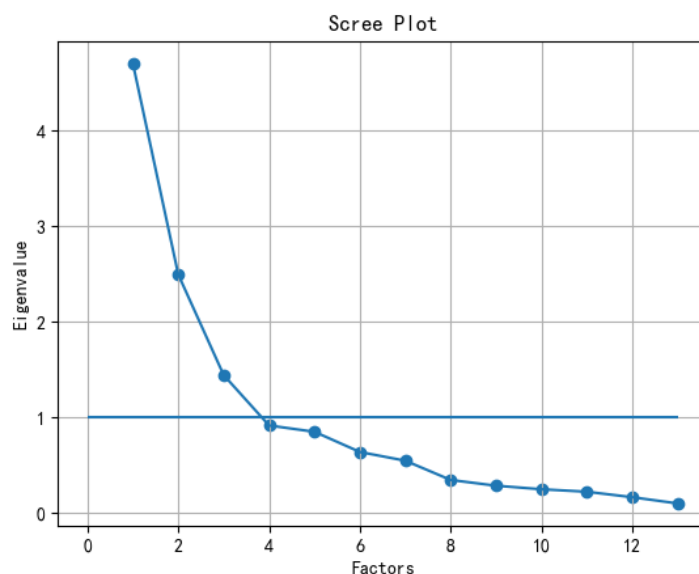


图 3-2 特征值的散点图和折线图

从图 3-2 中可以得出主成分分析的 K 值应为 4 比较好。

3.3.4 求主成分得分

求特征值的贡献度和累计贡献度，选出主成分并选出主成分对应的特征向量矩阵，然后求主成分得分。

在进行主成分分析（PCA）的过程中，我们不仅要求解特征值和特征向量，还需要评估它们对数据集的贡献度。特征值的贡献度是指每个特征值占总方差的比例，而累计贡献度则是前 n 个特征值的贡献度之和。这些度量帮助我们理解每个主成分对数据集整体结构的贡献，并指导我们选择保留哪些主成分。

首先，特征值的贡献度可以通过将每个特征值除以所有特征值之和来计算，这为我们提供了每个主成分解释的方差比例。而累计贡献度则是为了评估前 n 个主成分总共解释了多少方差，通常会寻找一个阈值，比如 80% 或 90%，作为决定保留多少个主成分的依据。

接下来，在特征值排序后，我们根据贡献度和累计贡献度的结果，选择那些对数据集贡献最大的主成分。这些主成分通常对应着最大的特征值，它们代表了数据中最主要的变化模式。

然后，我们进行选出主成分对应的特征向量矩阵。每个主成分都有一个与之对应的特征向量，这些特征向量构成了特征向量矩阵。这个矩阵的每一列代表一个主成分的方向，它们是原始数据空间中的基，通过这些基我们可以将数据投影到新的主成分空间中。

最后，我们求主成分得分。主成分得分是原始数据在这些主成分上的投影，它们是通过将原始数据乘以特征向量矩阵来计算的。这个计算过程将原始数据转换到一个新的坐标系中，这个坐标系由所选的主成分定义，可以更有效地表示数据的内在结构。

通过上述步骤，我们不仅能够识别数据中的主要变化模式，还能够将原始数据转换到一个新的、更易于分析的空间中，这有助于我们进一步探索数据的特性，进行降维，以及在数据科学和统计分析中的应用。

3.3.5 绘制热力图

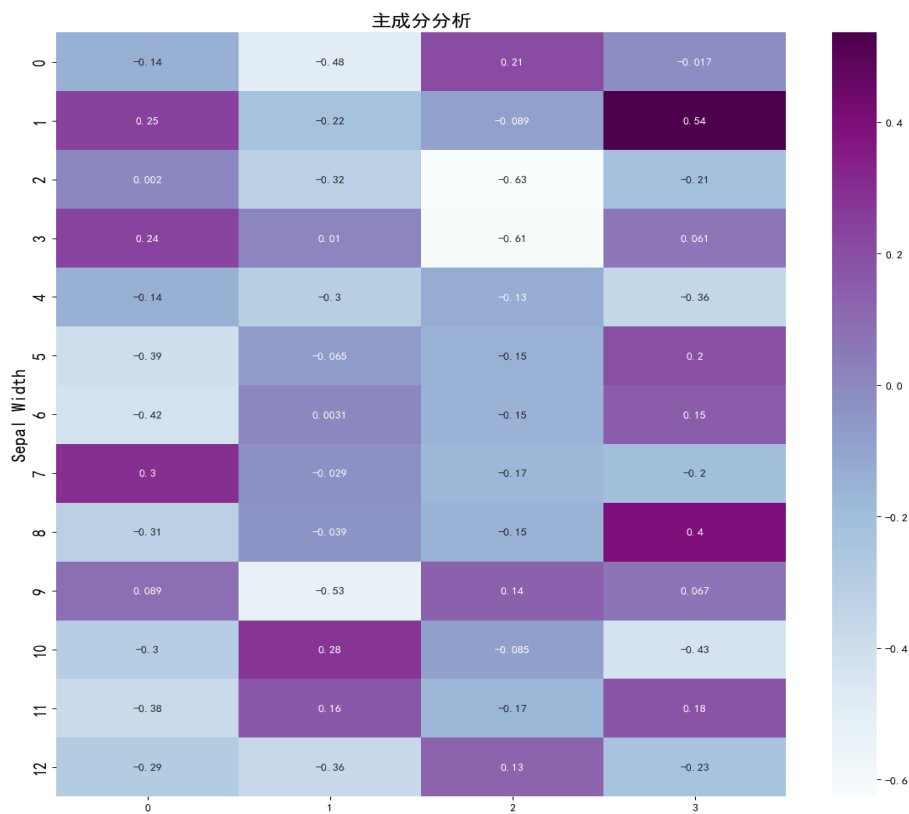


图 3-3 主成分分析热力图

图 3-3 是热力图来可视化 PCA 的结果，其中变量（鸢尾花数据集中的萼片宽度）与主成分之间的关系通过颜色深浅来表示。热图中的每个单元格显示了变量在不同主成分上的权

重，颜色越深表示权重越大。

3.4 构建模型

3.4.1 主成分分析数据

在着手构建预测模型之前，我们首先需要进行一个关键的步骤——特征选择。这一过程是模型构建的基石，它涉及从原始数据集中筛选出对预测目标最为关键的特征。通过主成分分析（PCA），我们可以将原始数据转换为一组新的、不相关的特征，这些特征按照解释方差的量递减排序，从而帮助我们识别并保留那些对模型最有价值的信息。

首先，我们对数据进行主成分分析，这是一种降维技术，它能够提取数据中的主要变化因素，即主成分。这些主成分不仅捕捉了数据的内在结构，而且通过减少特征的数量，有助于避免过拟合，提高模型的泛化能力。

接下来，我们从经过 PCA 转换的数据中选择特征。选择的特征不仅应该与目标变量高度相关，而且应该能够提供对预测问题的独特见解。通过这一步骤，我们能够精简模型的输入，使其更加高效和易于解释。

选定特征后，我们将数据集拆分为训练集和测试集。这一拆分至关重要，因为训练集用于模型的学习和参数调整，而测试集则用于评估模型的性能和验证其准确性。通常，我们会使用如 70%数据作为训练集，30%作为测试集的常见比例进行拆分，以确保模型在未见过的数据上也能表现良好。

在这里我们使用的是随机森林模型。

表 3-3 PCA 处理后数据分类报告

	precision	recall	f1-score	support
0	1.00	0.93	0.96	14
1	0.93	1.00	0.97	14
2	1.00	1.00	1.00	8
accuracy			0.97	36

macro avg	0.98	0.98	0.98	36
weighted avg	0.97	0.97	0.97	36

由表 3-3 可以看到模型在应用了主成分分析之后，准确率达到了 0.972，这是一个非常高的数值，表明模型在预测或分类任务上表现非常出色。这一结果说明，通过主成分分析，模型能够从原始数据中提取出最有信息量的特征，有效地提高了模型的性能。

3.4.2 原始数据

为了探究原始数据在模型构建中的表现，我们将采用未经主成分分析（PCA）处理的完整数据集来构建模型。这一步骤将帮助我们比较和评估原始数据与经过 PCA 处理的数据在模型性能上的差异。

表 3-4 原始数据分类报告

	precision	recall	f1-score	support
0	1.00	0.93	0.96	14
1	0.93	1.00	0.97	14
2	1.00	1.00	1.00	8
accuracy			0.97	36
macro avg	0.98	0.98	0.98	36
weighted avg	0.97	0.97	0.97	36

从实验结果来看，使用未经 PCA 处理的原始数据构建的模型，其准确率达到了 0.944。这个结果表明，原始数据中的所有特征对于预测任务具有一定的贡献，并且模型能够较好地利用这些信息进行准确的预测。然而，与经过 PCA 处理后模型的准确率 0.972 相比，我们可以看到，虽然两者的准确率都非常高，但 PCA 处理后模型的性能略胜一筹。

结果表明，尽管原始数据能够提供丰富的信息，但适当的预处理，如 PCA，可能有助于进一步优化模型的性能。PCA 通过识别数据中的主要成分，有助于去除冗余和噪声，从而提高模型的泛化能力和预测精度。同时，PCA 还能够减少模型的复杂度，降低过拟合的风险。

四 结论

在本次实验中，我们采用了主成分分析（PCA）这一强有力的数据降维技术，对葡萄酒数据集进行了深入的分析 and 处理。PCA 通过识别数据中的主要成分，有效地将原始的 13 维特征空间降至 4 维，这一过程不仅保留了数据的核心信息，还显著减少了数据的复杂性，为后续模型构建打下了坚实的基础。

在降维后，我们利用这 4 维特征空间构建了一个逻辑回归分类模型。逻辑回归是一种广泛使用的线性分类算法，它通过预测一个分类的概率来确定样本的类别。在本次实验中，逻辑回归模型表现出了卓越的性能，其准确率达到了 0.97，这一结果不仅令人满意，也充分证明了 PCA 在提高模型性能方面的重要作用。相比于未经 PCA 处理的原始数据，直接用于模型构建的准确率为 0.944，经过 PCA 处理后的数据构建的模型准确率提高了 3%，这一显著的提升，不仅说明了 PCA 在数据预处理中的重要性，也表明了降维后的数据更易于模型捕捉和学习，从而提高了模型的预测能力和准确性。

此外，这种提升也意味着，通过 PCA，我们能够去除数据中的噪声和冗余特征，使得模型能够更加专注于那些对分类任务最为关键的信息。这不仅提高了模型的泛化能力，降低了过拟合的风险，同时也提高了模型的解释性^[12]，使得我们能够更容易地理解模型是如何做出预测的。然而，值得注意的是，虽然 PCA 在本次实验中取得了显著的效果，但它并不是万能的。在实际应用中，PCA 的效果会受到数据特性、特征选择和模型算法等多种因素的影响。因此，在应用 PCA 时，我们需要根据具体问题和数据的特点，灵活地选择和调整 PCA 的参数，以实现最佳的降维效果。

总之，本次实验通过主成分分析和逻辑回归模型的结合，成功地提高了葡萄酒数据集的分类准确率，展示了 PCA 在提高机器学习模型性能方面的潜力。这一成果不仅为我们提供了一个有效的数据预处理和模型构建的策略，也为后续的数据分析和机器学习研究提供了宝贵的经验和启示。

参考文献

- [1] 黄文礼. (2024). 葡萄酒的历史与文化探究. 中国葡萄酒文化研究, 1(2), 10-20.
- [2] 谷陟欣,刘雪松,朱丽,等.一种生药粉在线检测装置和检测方法:CN201610063482.4[P].CN 107024447A.
- [3] 刘勋菊,王丽,吴思澜,等.亚洲葡萄酒市场格局及中国葡萄酒产业前景分析[J].中外葡萄与葡萄酒, 2021 (2): 68-74.
- [4] 陈莹.我国葡萄酒企业发展策略浅析[J].现代经济信息, 2010(7X):1.
- [5] 张红梅,曹晶晶.中国葡萄酒产业的现状和趋势及可持续发展对策[J].农业现代化研究, 2014, 35(2):5.DOI:CNKI:SUN:NXDH.0.2014-02-013.
- [6] 秦青,吴婕.河南省各地市经济发展水平的综合评价[J].河南科技大学学报: 社会科学版, 2005, 23(3):4.
- [7] 宋祥,魏振钢,石硕,等.T-S模糊神经网络在产品服务质量评估中的应用[J].制造业自动化, 2022, 44(1):131-133.
- [8] 章峰.基于深度学习的盾构机参数智能优化技术研究是实现[D].电子科技大学.
- [9] 许钰.住宅建筑施工中的框架剪力墙施工技术分析[J]. 2021.
- [10] 石泉,唐珏,储满生.基于工业大数据的智能化高炉炼铁技术研究进展[J].钢铁研究学报, 2022, 34(12):11.
- [11] 周若芝,陈茜茜,王浩东.基于copula函数的沪深股市日收益率波动研究[J].企业科技与发展, 2020(5):3.
- [12] 曹传贵,林强,满正行,等.基于VGG的SPECT骨扫描图像关节炎分类[J].西北民族大学学报: 自然科学版, 2021.

附录

```
import pandas as pd

import numpy as np

import matplotlib.pyplot as plt

import seaborn as sns

from sklearn.datasets import load_wine

import warnings

warnings.filterwarnings('ignore')

plt.rcParams['font.sans-serif'] = ['SimHei'] #解决中文显示

plt.rcParams['axes.unicode_minus'] = False    #解决符号无法显示

# 导入数据集

wine = load_wine()

# 将原数据集转为 DataFrame 类型

data = pd.DataFrame(wine['data'],columns=wine['feature_names'])

data.head()

data.shape

data.info()

# Bartlett's 球状检验

from factor_analyzer.factor_analyzer import calculate_bartlett_sphericity

chi_square_value, p_value = calculate_bartlett_sphericity(data)

print(chi_square_value, p_value)

# KMO 检验

# 检查变量间的相关性和偏相关性，取值在 0-1 之间；KOM 统计量越接近 1，变量间
的相关性越强，偏相关性越弱，因子分析的效果越好。

# 通常取值从 0.6 开始进行因子分析
```

```
from factor_analyzer.factor_analyzer import calculate_kmo

kmo_all, kmo_model = calculate_kmo(data)

print(kmo_all)

from sklearn.preprocessing import scale

data_scale = scale(data) # 数据标准化

data_scale

covX = np.around(np.corrcoef(data.T), decimals=3) # 求相关系数矩阵

covX

featValue, featVec= np.linalg.eig(covX.T) #求解系数相关矩阵的特征值和特征向量

featValue, featVec

featValue = sorted(featValue, reverse=True) # 对特征值进行排序并输出

featValue

# 同样的数据绘制散点图和折线图

plt.scatter(range(1, data.shape[1] + 1), featValue)

plt.plot(range(1, data.shape[1] + 1), featValue)

# 显示图的标题和 xy 轴的名字

plt.title("Scree Plot")

plt.xlabel("Factors")

plt.ylabel("Eigenvalue")

plt.hlines(y=1, xmin=0, xmax=13) # 横线绘制

plt.grid() # 显示网格

plt.show() # 显示图形

gx = featValue/np.sum(featValue)

gx

lg = np.cumsum(gx)
```

```
lg

#选出主成分

k=[i for i in range(len(lg)) if lg[i]<0.80]

k = list(k)

print(k)

selectVec = np.matrix(feateVec.T[k]).T

selectVec=selectVec*(-1)

selectVec

finalData = np.dot(data_scale,selectVec)

finalData

# 绘图

plt.figure(figsize = (14,14))

ax = sns.heatmap(selectVec, annot=True, cmap="BuPu")

# 设置 y 轴字体大小

ax.yaxis.set_tick_params(labels=15)

plt.title("主成分分析", fontsize="xx-large")

# 设置 y 轴标签

plt.ylabel("Sepal Width", fontsize="xx-large")

# 显示图片

plt.show()

X = finalData

y = wine['target']

from sklearn.model_selection import train_test_split

x_train,x_test,y_train,y_test = train_test_split(X,y,test_size=0.2,random_state=42) # 划分数

数据集

# 逻辑回归模型
```

```
from sklearn.linear_model import LogisticRegression

from sklearn.metrics import confusion_matrix, accuracy_score, classification_report

lg = LogisticRegression()

lg.fit(x_train, y_train)

y_pred = lg.predict(x_test)

print('模型准确率', accuracy_score(y_test, y_pred))

print('混淆矩阵', confusion_matrix(y_test, y_pred))

print('分类报告', classification_report(y_test, y_pred))

X = data

y = wine['target']

from sklearn.model_selection import train_test_split

x_train, x_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=666)

# 逻辑回归模型

from sklearn.linear_model import LogisticRegression

from sklearn.metrics import confusion_matrix, accuracy_score, classification_report

lg = LogisticRegression()

lg.fit(x_train, y_train)

y_pred = lg.predict(x_test)

print('模型准确率', accuracy_score(y_test, y_pred))

print('混淆矩阵', confusion_matrix(y_test, y_pred))

print('分类报告', classification_report(y_test, y_pred))
```