

结构化数据医院——  
融合材料领域知识的  
结构化数据质量治理  
平台

V1.0

用户操作手册

文档编号：2025-R0374

最后修订日期：2025 年 7 月 14 日

版权所有 2025-。保留所有权利。

本文档中所提到的其他所有公司和产品名，均是其所有人的商标或者注册商标。

本文档内容如有变动，恕不另行通知。

# 目录

关于本文档	1
<b>第 1 章 软件概述</b>	<b>2</b>
1.1 模块介绍	2
1.2 用户层级、角色与相应权限	3
1.3 登录界面	3
1.4 主界面和导航栏	4
<b>第 2 章 数据挂号模块</b>	<b>5</b>
2.1 数据上传	5
2.2 数据模板预览和下载	7
<b>第 3 章 数据体检模块</b>	<b>9</b>
3.1 数据体检	9
3.1.1 可溯源性	9
3.1.2 时间敏感性	10
3.1.3 完整性	10
3.1.4 一致性	11
3.1.5 准确性	11
3.1.6 均衡性	13
3.1.7 规范性	13
3.1.8 冗余性	14
3.1.9 洞察力	15
3.2 元特征报告	15
3.3 体检报告	16
<b>第 4 章 数据诊治模块</b>	<b>18</b>
4.1 可溯源性诊治	18
4.2 时间敏感性诊治	19
4.3 完整性诊治	20
4.4 一致性诊治	24
4.5 准确性诊治	26
4.6 均衡性诊治	32
4.7 规范性诊治	33
4.8 冗余性诊治	36
4.8 洞察力诊治	41

# 关于本文档

## 主题

欢迎使用结构化数据医院——融合材料领域知识的结构化数据质量治理平台的用户操作手册，该手册包含了需要了解和使用结构化数据医院——融合材料领域知识的结构化数据质量治理平台的相关知识。

本文档包含了如下章节：

- **第 1 章 软件概述**
- **第 2 章 数据挂号**
- **第 3 章 数据体检**
- **第 4 章 数据上传**

序言中介绍了使用该文档的一些帮助信息。

## 读者

使用结构化数据医院——融合材料领域知识的结构化数据质量治理平台的最终用户和系统管理员。

# 第1章 软件概述

## 本章概述

本章介绍了结构化数据医院——融合材料领域知识的结构化数据质量治理平台的功能和操作界面特点。

## 内容

主题	页码
模块介绍	2
用户层级、角色与相应权限	3
登录界面	3
主界面和导航栏	4

## 1.1 模块介绍

在材料科学的研究中，数据不仅支撑着新材料的发现与性能预测，也逐步成为推动材料智能设计的核心要素。然而，目前材料数据普遍存在评测标准不一、质量参差不齐等问题，严重制约了基于机器学习的在材料领域中的运用。因此，如何在不依赖大量实验或高成本计算的前提下，提升材料数据的质量，已成为材料智能研究面临的关键挑战。因此，我们研发了“结构化数据医院——融合材料领域知识的结构化数据质量治理平台”。该平台包含数据挂号、数据体检、数据诊治三大核心功能模块：首先用户在数据挂号模块中填写并提交相关数据与信息；随后，在数据体检模块中，系统将会从可溯源性、时间敏感性、完整性、一致性、准确性、均衡性、规范性、冗余性、洞察力九大维度对数据进行质量评估，生成面向机器学习的数据画像报告和数据体检报告；最后，在数据诊治模块中用户可以分别从九大维度开展定向数据质量治理，系统性提升数据质量。平台基于 B/S 架构，采用 Vue.js、Element UI、Flask 与 Spring Boot 等技术构建，具备良好的扩展性与可维护性。同时，平台内置多种基于 Python 开发的数据质量治理算法，为用户提供了全面、高效的数据质量管理服务。平台旨在为材料领域的机器学习研究提供可靠、高质量的数据支撑。

具体来说，结构化数据医院——融合材料领域知识的结构化数据质量治理平台设计与实现了如图 1-1 所示的数据挂号、数据体检和数据诊治这三大模块。

### 1. 数据挂号

此模块是一个用于实现上传结构化数据相关信息的工具。用户在此依次填写数据基本信息、用户信息和来源信息，完成结构化数据的上传。系统支持上传.xlsx 和.csv 格式的结构化数据。同时，包含数据模板预览和下载以及上传数据的预览。

### 2. 数据体检

此模块是一个用于对数据进行全面数据质量体检的工具。平台首先计算 27 个元特征维度全面系统地刻画数据的特点，然后从可溯源性、时间敏感性、完整性、一致性、准确性、均衡性、规范性、冗余性和洞察力九大质量维度定性定量评估数据的质量，依次展示评估结果，形成数据画像报告和体检报告，并提供针对性诊治建议，引导用户前往相应科室进行质量提升。

### 3. 数据诊治

此模块是一个用于对数据进行全方位数据质量治理的工具。在数据诊治模块，平台将分

别从可溯源性、时间敏感性、完整性、一致性、准确性、均衡性、规范性、冗余性、洞察力九个治理维度对数据进行全面治理。

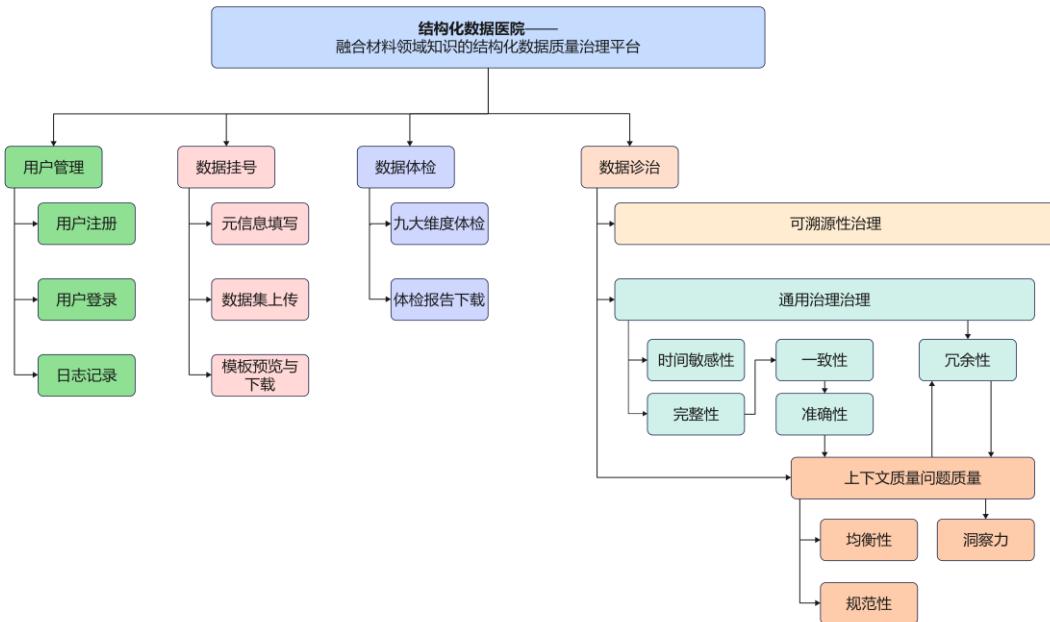


图 1-1 系统软件架构

## 1.2 用户层级、角色与相应权限

使用本系统的用户主要分为两个不同层级，即：管理员，普通用户，如表 1-1 所示。

表 1-1 用户层级

用户层级	用户名	密码
管理员	admin	admin123
普通用户	123@qq.com	111

不同层级角色的用户各自享有不同权限。例如对于同一功能，管理员可对其他用户信息进行编辑，而普通用户只可修改自身用户信息。

如在本手册阅读过程中，您发现所描述的内容与您所实际看到的内容不一致，可能是由于权限差异。

## 1.3 登录界面

如图 1-2 所示是结构化数据医院——融合材料领域知识的结构化数据质量治理平台的初始登陆页面。

对于用户，需选择对应账号及密码输入（表 1-1），以登录进入主界面。



图 1-2 登陆界面

## 1.4 主界面和导航栏

平台主界面由顶部区域、侧边导航栏模块构成，如图 1-3 所示。顶部区域包括平台 logo 和退出按钮。侧边导航栏包括数据挂号、数据体检和数据诊治三大模块，用户可根据研究需求选择对应模块开展具体工作。首先用户在数据挂号模块中填写并提交相关数据与信息；随后，在数据体检模块中，系统将会从可溯源性、时间敏感性、完整性、一致性、准确性、均衡性、规范性、冗余性、洞察力九大维度对数据进行质量评估，生成面向机器学习的数据画像报告和数据体检报告；最后，在数据诊治模块中用户可以分别从九大维度开展定向数据质量治理，系统性提升数据质量。



图 1-3 主界面

# 第 2 章 数据挂号模块

## 本章概述

本章主要介绍了结构化数据医院——融合材料领域知识的结构化数据质量治理平台的数据挂号部分的内容。我们提供了不同信息填写的表单以及模板预览与下载，让用户更好地上传数据。

## 本章内容

主题	页码
数据上传	5
数据模板预览与下载	7

## 2.1 数据上传

(1) 上传数据基本信息。在对应页面中输入数据的基本信息，包括数据集名称、样本数量、特征数量、目标属性、材料类别、获取方式、机器学习任务、时序特征名称（如果有时序特征）、数据集描述等，并上传测试数据集，如图 2-1 所示。

(2) 上传用户信息。点击“下一步”，进入用户信息上传页面，在对应页面中输入用户的信息，包括提交者、所在单位、邮箱、电话等，如图 2-2 所示。

(3) 上传数据来源信息。点击“下一步”，进入数据来源信息上传页面，在对应页面中输出数据来源信息，包括数据分类、数据来源等点击“提交”，存储所填写的信息，完成数据挂号，如图 2-3 所示。

数据集上传

数据基本信息

用户信息

来源信息

\* 数据集名称: 请输入数据集名称

\* 特征数量: 请输入特征数量

关键词: 请填写关键词, 多个用";"分隔

\* 样本数量: 请输入样本数量

\* 目标属性: 请输入目标属性

\* 材料类别: 请选择材料类别

\* 获取方式: 请选择获取方式

\* 机器学习任务:  回归  分类  聚类

是否含有时序特征:

目标属性阈值: 请输入目标属性的阈值(英文逗号)

\* 时序特征名称: 请输入时序特征名称, 用英文逗号

数据领域类型: 请选择数据领域类型

方向类型: 请选择方向类型

\* 数据集描述: 请输入数据集描述

文件上传: 将文件拖到此处, 或点击上传

请上传数据集, 支持的文件扩展名为.xlsx  
或.csv (第一列为索引, 最后一列为决策属性, 第一行为属性名称)

下载标准样式文件 上一步 下一步

图 2-1 数据基本信息上传页面

数据集上传

数据基本信息

用户信息

提交者: qs

所在单位: shu

校对者: qs

\* 邮箱: 123@qq.com

\* 电话: 19999999999

通讯地址: 上海大学

上一步 下一步

数据模板与预览

数据模板 上传数据预览

训练集模板 >

测试集模板 >

图 2-2 用户信息上传页面

图 2-3 数据来源信息上传页面

## 2.2 数据模板预览和下载

用户首先在数据模板与预览区域点击左边“下拉”标签即可看到训练集和测试集的数据模板，如图 2-4 所示。点击“下载模板”便可进行数据模板的下载，如图 2-5 所示。

图 2-4 展示了数据模型预览功能。在“数据集上传”模块下方，有一个名为“数据模型与预览”的子模块。该模块包含“数据模板”和“上传数据预览”两个部分。“数据模板”部分显示了训练集模板的预览，格式要求为 {1\_X, 2\_X, ..., n\_X} 表示特征，Y 表示目标属性（应为浮点或整数类型）。右侧有一个“下载模板”按钮。下方是测试集模板的预览。

图 2-4 数据模型预览

图 2-5 展示了数据模型下载功能。在“数据集上传”模块下方，有一个名为“数据模型与预览”的子模块。右侧有一个“近期的下载记录”窗口，显示了两个文件：“训练集模板 (1).xlsx”（完成）和“训练集模板.xlsx”（2小时前）。下方有一个“完整的下载记录”按钮。下方是训练集模板的预览。

图 2-5 数据模型下载

# 第3章 数据体检模块

## 本章概述

本章主要介绍了结构化数据医院——融合材料领域知识的结构化数据质量治理平台的数据体检部分的内容。我们首先计算 27 个元特征维度全面系统地刻画数据的特点，然后从可溯源性、时间敏感性、完整性、一致性、准确性、均衡性、规范性、冗余性和洞察力九大质量维度定性定量评估数据的质量，依次展示评估结果，形成数据画像报告和体检报告。

## 本章内容

主题	页码
数据体检	9
元特征报告	15
体检报告	16

## 3.1 数据体检

点击系统主页面中“数据体检”进入数据体检页面。平台将依次从可溯源性、时间敏感性、完整性、一致性、准确性、均衡性、规范性、冗余性和洞察力九大质量维度定性定量评估数据的质量。此处不需要人工操作，系统将根据上传的数据自动计算并展现结果。

### 3.1.1 可溯源性

平台在可溯源性维度展示了元信息完整率，包括基础信息完整率、材料知识完整率和机器学习任务信息完整率。此处以 nasicon 数据集为例进行测试，运行结果如下图 3-1 所示，在数据挂号部分我们上传的信息通过数据体检评估可知其基础信息完整率、材料知识完整率和机器学习任务信息完整率都是百分百。

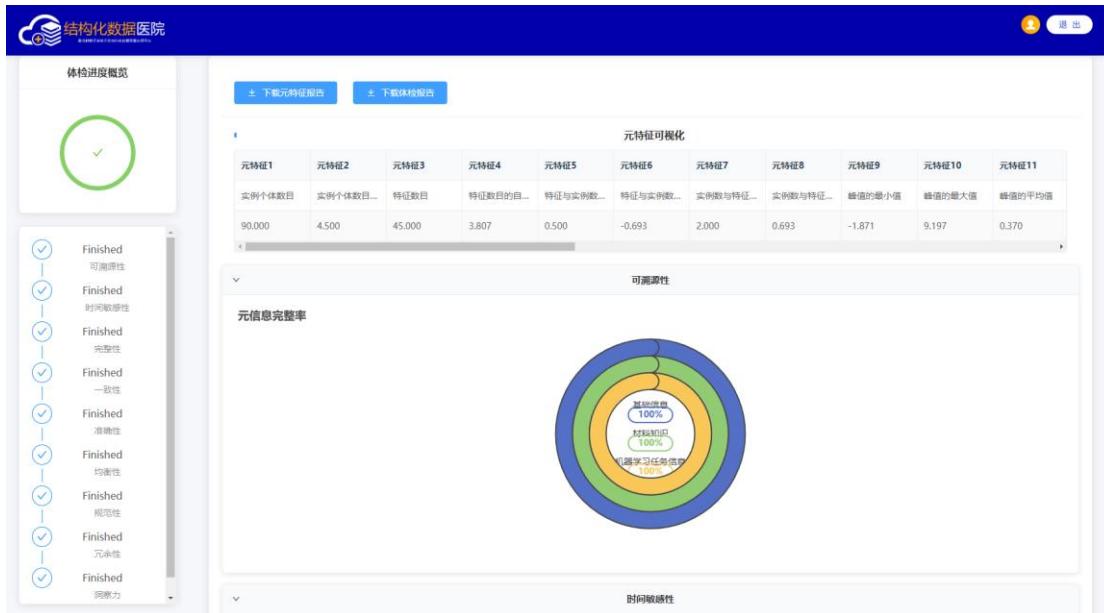


图 3-1 数据体检可溯源性

### 3.1.2 时间敏感性

平台在时间敏感性维度展示了时间序列的最小嵌入维。此处假设 nasicon 数据集的“a”特征列为其时间序列，并以此为例进行测试，运行结果如下图 3-2 所示，可知“a”的最小嵌入维为 4。

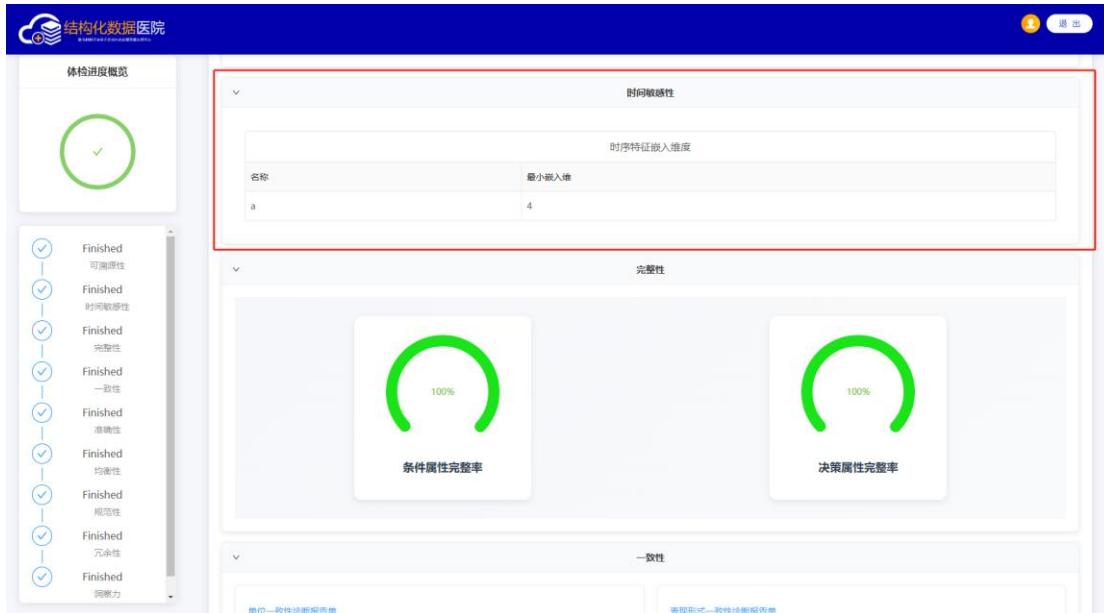


图 3-2 数据体检时间敏感性

### 3.1.3 完整性

平台在完整性维度展示了条件属性完整率和决策属性完整率。此处以 nasicon 数据集为例进行测试，运行结果如下图 3-3 所示，可知我们的条件属性完整率和决策属性完整率都为百分百。

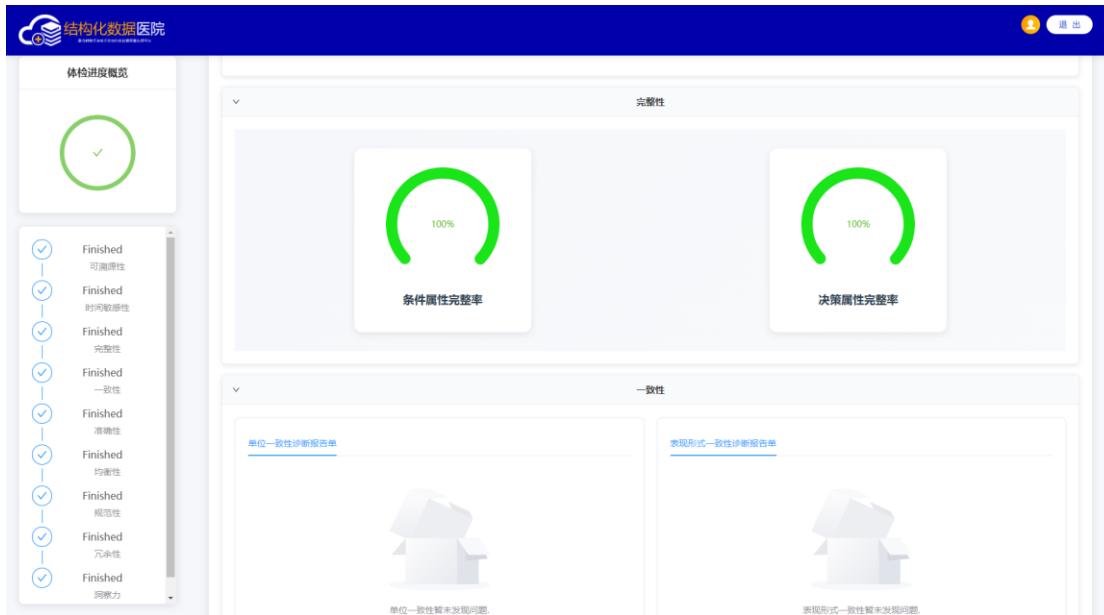


图 3-3 数据体检完整性

### 3.1.4 一致性

平台在一致性维度展示了单位一致性诊断报告单和表现形式一致性诊断报告单。此处以 nasicon 数据集为例进行测试，运行结果如下图 3-4 所示，可知单位一致性和表现形式一致性都不存在问题。

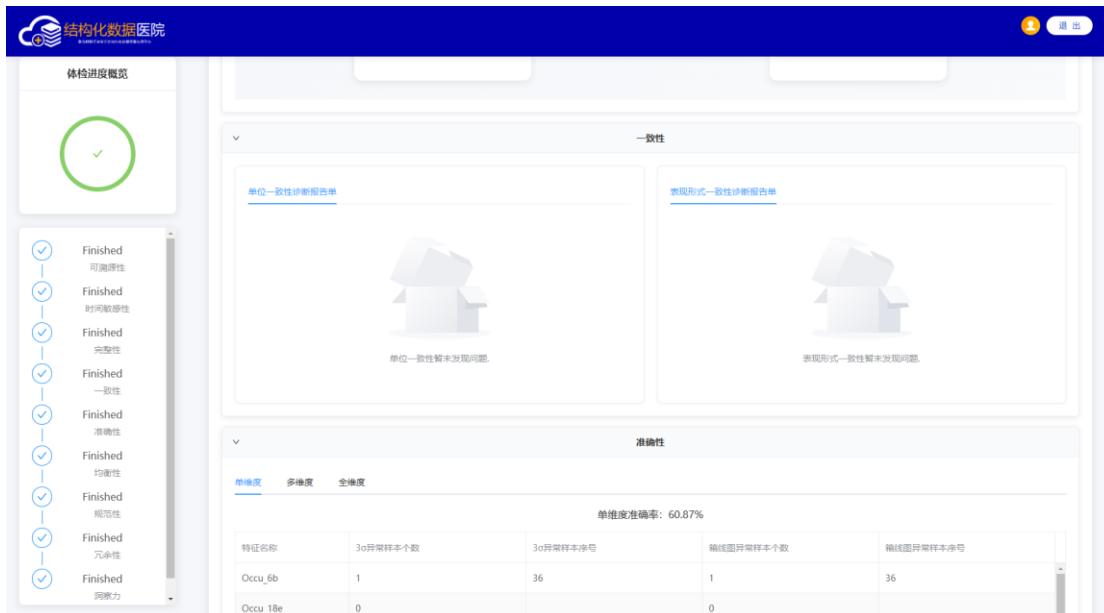


图 3-4 数据体检一致性

### 3.1.5 准确性

平台在准确性维度展示了单维度准确率、多维度准确率和全维度准确率。此处以 nasicon 数据集为例进行测试，单维度运行结果如下图 3-5 所示，可知单维度准确率为 60.87%；多维度运行结果如下图 3-6 所示，可知多维度准确率为 89.08；全维度运行结果如下图 3-7 所示，

可知全维度准确率为 86.67%。

特征名称	3σ异常样本个数	3σ异常样本序号	箱线图异常样本个数	箱线图异常样本序号
Occu_6b	1	36	1	36
Occu_18e	0		0	
Occu_36f	0		21	27, 28, 29, 37, 38, 39, 40, 41, 42, 43, 6, 2, 63, 64, 65, 68, 69, 70, 71, 77, 78, 89
C_Na	0		0	
Occu_M1	0		0	
Occu_M2	0		4	3, 4, 29, 50
EN_M1	1	30	1	30
EN_M2	0		0	
avg_EN_M	2	30, 56	2	30, 56

图 3-5 数据体检准确性单维度

特征	Occu_6b	Occu_18e	Occu_36f	C_Na	Occu_M1	Occu_M2	EN_M1	EN_M2	avg_EN_M	Radius_M1	Radius_M2	avg_Raduis_M	Valence_M1	Valence_M2	avg_Valence_M
Occu_6b	1	-0.3434	-0.6476	-0.2982	-0.1913	0.2105	0.2381	0.2311	0.2164	-0.3392	0.2192	-0.3248	0.0983	0.1774	0.1244
Occu_18e	-0.3434	1	0.2589	0.9806	0.0358	-0.0237	-0.0615	0.0148	-0.0458	0.1918	0.0167	0.1671	-0.453	-0.0802	-0.653
Occu_36f	-0.6476	0.2589	1	0.3755	0.1393	-0.1366	-0.344	-0.2252	-0.3463	0.4326	-0.1969	0.4537	-0.0524	-0.1881	-0.095
C_Na	-0.2982	0.9806	0.3755	1	0.0315	-0.015	-0.0989	0.0048	-0.0884	0.229	0.0114	0.213	-0.4395	-0.0904	-0.642
Occu_M1	-0.1913	0.0358	0.1393	0.0315	1	-0.9881	-0.3215	-0.8483	-0.3214	0.4972	-0.8793	0.3961	0.5014	-0.8889	0.1317
Occu_M2	0.2105	-0.0237	-0.1366	-0.015	-0.9881	1	0.3037	0.8369	0.2672	-0.5196	0.8806	-0.3991	-0.482	0.8674	-0.1651
EN_M1	0.2381	-0.0615	-0.344	-0.0989	-0.3215	0.3037	1	0.3379	0.9114	-0.6335	0.3211	-0.6342	-0.3318	0.3195	-0.201
EN_M2	0.2311	0.0148	-0.2252	0.0048	-0.8483	0.8369	0.3379	1	0.4165	-0.4513	0.9594	-0.4115	-0.4413	0.9309	-0.224
avg_EN_M	0.2164	-0.0458	-0.3463	-0.0884	-0.3214	0.2672	0.9114	0.4165	1	-0.5202	0.3311	-0.6261	-0.3167	0.3476	-0.177

图 3-6 数据体检准确性多维度

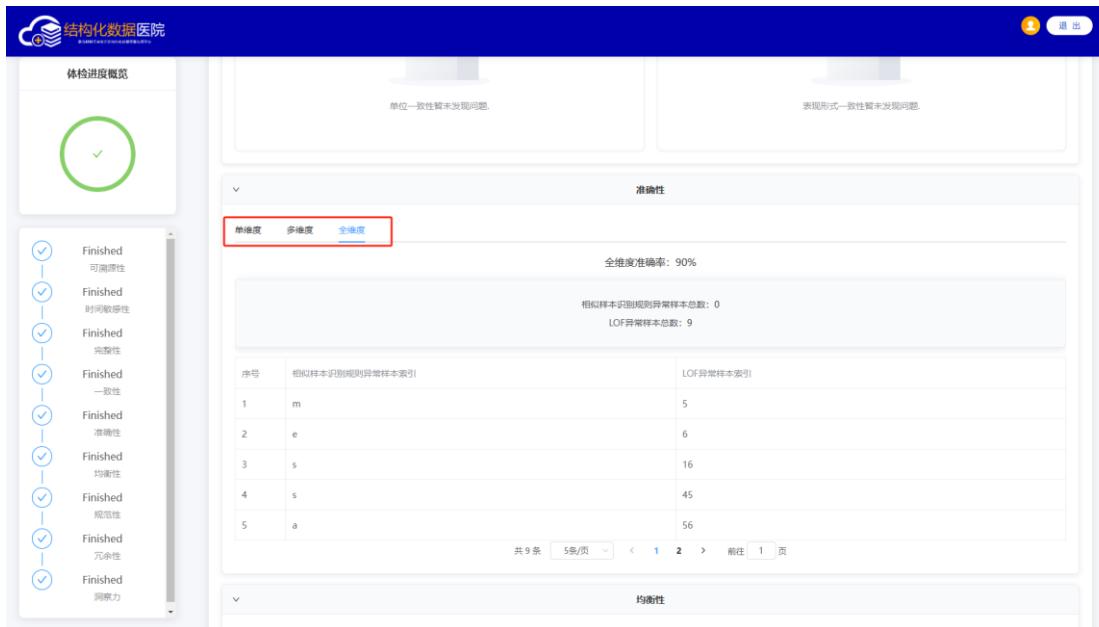


图 3-7 数据体检准确性全维度

### 3.1.6 均衡性

平台在均衡性维度展示了数据均衡性指标。此处以 nasicon 数据集为例进行测试，运行结果如下图 3-8 所示，可知数据均衡性指标为 7.41%。

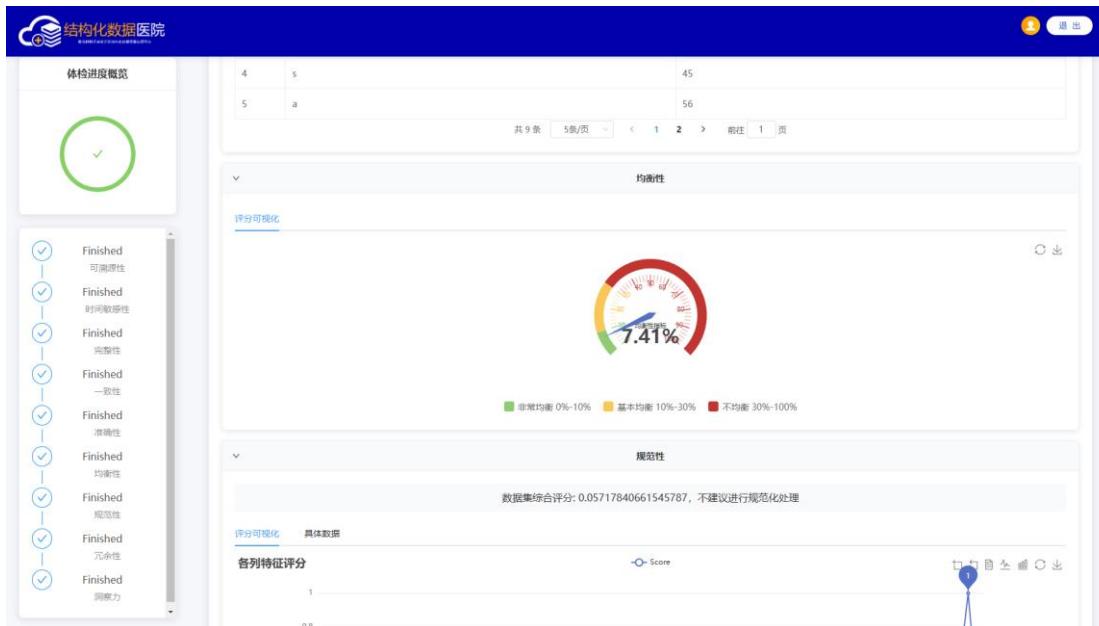


图 3-8 数据体检均衡性

### 3.1.7 规范性

平台在规范性维度展示了数据集综合评分以及各列特征评分可视化和具体数据。此处以

nasicon 数据集为例进行测试，运行结果如下图 3-9 所示，可知数据集综合评分为 0.057。



图 3-9 数据体检规范性

### 3.1.8 冗余性

平台在冗余性维度展示了冗余样本绝对数、冗余样本相对数、冗余特征绝对数、冗余特征相对数及其评分可视化和具体数据。此处以 nasicon 数据集为例进行测试，运行结果如下图 3-10 所示，可知数据集冗余样本绝对数为 29，冗余特征绝对数为 22。

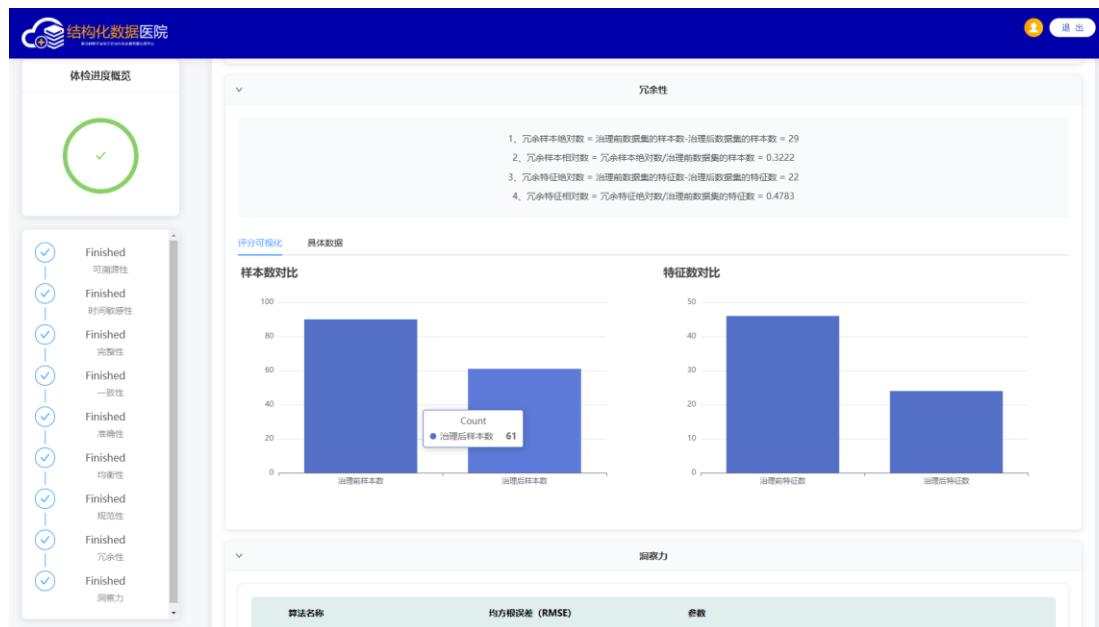


图 3-10 数据体检冗余性

### 3.1.9 洞察力

平台在洞察力维度展示了该数据集在回归模型上 RMSE 最小的前三个模型及其对应的 RMSE 和参数。此处以 nasicon 数据集为例进行测试，运行结果如下图 3-11 所示。



图 3-11 数据体检洞察力

### 3.2 元特征报告

平台在数据体检界面提供了下载元特征报告功能，点击“下载元特征报告”即可生成元特征报告。此处以 nasicon 数据集为例进行测试，生成元特征报告如下图 3-12 所示。

**数据基本情况及主诉**

数据名:	Nasicon	所有者:	shu
获取方式:	实验测量	样本量:	90
特征量:	45	目标属性:	Activation energy (Ea)
材料类型:	[“复合材料”, “基体材料”, “金属基复合材料”]		
下潜任务类型:	机器学习的回归任务		

**面向机器学习建模的数据画像**

基于基础信息的元特征			
实例数量	90.000	实例数量对数值	4.500
特征数量	45.000	特征数量对数值	3.807
特征数量/实例数量	0.500	特征数量/实例数量对数值	-0.693
实例数量/特征数量	2.000	实例数量/特征数量对数值	0.693
基于统计信息的元特征			
峰度最大值	-1.871	偏度最大值	9.197
峰度最小值	0.370	偏度最小值	2.139
峰度平均值	-2.008	偏度平均值	2.843
峰度标准差	0.434	偏度标准差	1.118
基于主成分分析的元特征			
解释95%方差的成分比例	0.267		
第一个主成分峰度	-0.979	第一个主成分偏度	0.449
基于回归基准模型的元特征			
Lasso 回归指标	0.743	KNN 回归指标	0.857
SVM 回归指标	0.514	MLP 回归指标	0.888
DT 回归指标	0.553		
基于决策属性统计信息的元特征			
决策属性峰度	10.780	决策属性偏度	3.234
基于不确定性的元特征			
(均值, 峰, 超峰)	1.284, 0.595, 0.133		

本次检测结果由融合材料领域知识的结构化数据质量治理平台提供

图 3-12 数据体检的元特征报告

### 3.3 体检报告

平台在数据体检界面提供了下载体检报告功能，点击“下载体检报告”即可生成体检报告。此处以 nasicon 数据集为例进行测试，生成体检报告如下图 3-13 所示。

**结构化数据医院**

**结构化数据医院数据质量体检报告**

体检时间: 2025年08月09日 12:13 体检科室: 可溯源性等0个科室

**数据基本情况及主诉**

数据名:	Nasicon	所有者:	shu
获取方式:	["实验测量"]	样本量:	90
特征量:	45	目标属性:	Activation energy (Ea)
材料类型:	["复合材料", "基体材料", "金属基复合材料"]	下游任务类型:	机器学习的回归任务

**面向机器学习建模的体检项目及体检结果**

项目名称	结果	优秀指标	提示	项目名称	结果	优秀指标	提示
<b>可溯源性</b>				<b>时间敏感性</b>			
元数据完整性	75.00%	=100%	有问题	时序属性名称	a	--	--
				最小嵌入维数	4	--	--
<b>完整性</b>				<b>一致性</b>			
决策属性完整性	100%	=100%	无问题	单位一致性	100%	=100%	无问题
条件属性完整性	100%	=100%	无问题	表现形式一致性	100%	=100%	无问题
<b>准确性</b>				<b>均衡性</b>			
单维度准确率	60.87%	=100%	有问题	变异系数值	7.41%	0%-10%	高度均衡
多维度准确率	100%	=100%	无问题				
全维度准确率	90%	=100%	有问题				
<b>规范性</b>				<b>冗余性</b>			
特征规范率	0.06	=0	有问题	冗余样本量	29	=0	有问题
				冗余特征量	22	=0	有问题
<b>洞察力</b>							
推荐前三种算法		KRR	0.564	模型1及其RMSE		模型2及其RMSE	
				SGD	0.57	MLP	0.595
<b>诊断结论与治理建议</b>							
<b>诊断结论:</b> Nasicon 数据的 完整性、一致性、洞察力较好，但是 可溯源性、准确性、均衡性、规范性、冗余性需要进行治理。				<b>治理建议:</b> 首先需要针对数据可溯源性、准确性进行治理，并关注均衡性、规范性；在此基础上，进行数据冗余性的治理。			

本次检测结果由融合材料领域知识的结构化数据质量治理平台提供

图 3-13 数据体检的体检报告

# 第 4 章 数据诊治模块

## 本章概述

本章主要介绍了结构化数据医院——融合材料领域知识的结构化数据质量治理平台的数据诊治部分的内容。用户可通过此模块对数据体检中数据存在的质量问题从可溯源性、时间敏感性、完整性、一致性、准确性、均衡性、规范性、冗余性、洞察力九个治理维度分别进行质量治理。

## 本章内容

主题	页码
可溯源性诊治	18
时间敏感性诊治	19
完整性诊治	20
一致性诊治	24
准确性诊治	26
均衡性诊治	32
规范性诊治	33
冗余性诊治	36
洞察力诊治	41

## 4.1 可溯源性诊治

点击平台页面左侧导航栏中的“可溯源性科室” - “检查报告”，进入到可溯源性页面。平台在可溯源性维度展示了元信息完整率和数据溯源表格展示。此处以 nasicon 数据集为例进行测试，系统可通过在数据挂号上传的信息来进行元信息完整率的计算以及数据溯源中各个字段的信息显示，运行结果如下图 4-1 所示。

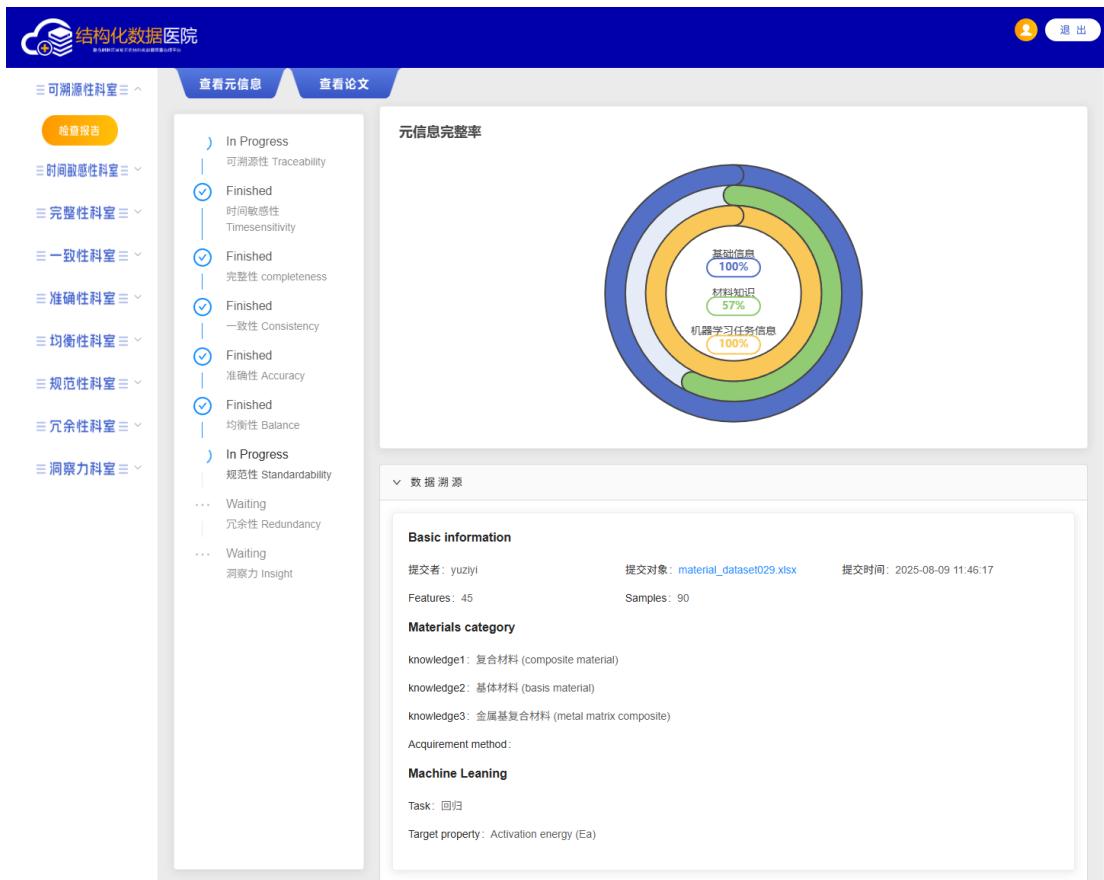


图 4-1 数据诊治可溯源性

## 4.2 时间敏感性诊治

(1) 点击平台页面左侧导航栏中的“时间敏感性科室” - “导医台”，进入到时间敏感性文件选择页面，运行结果如下图 4-2 所示。

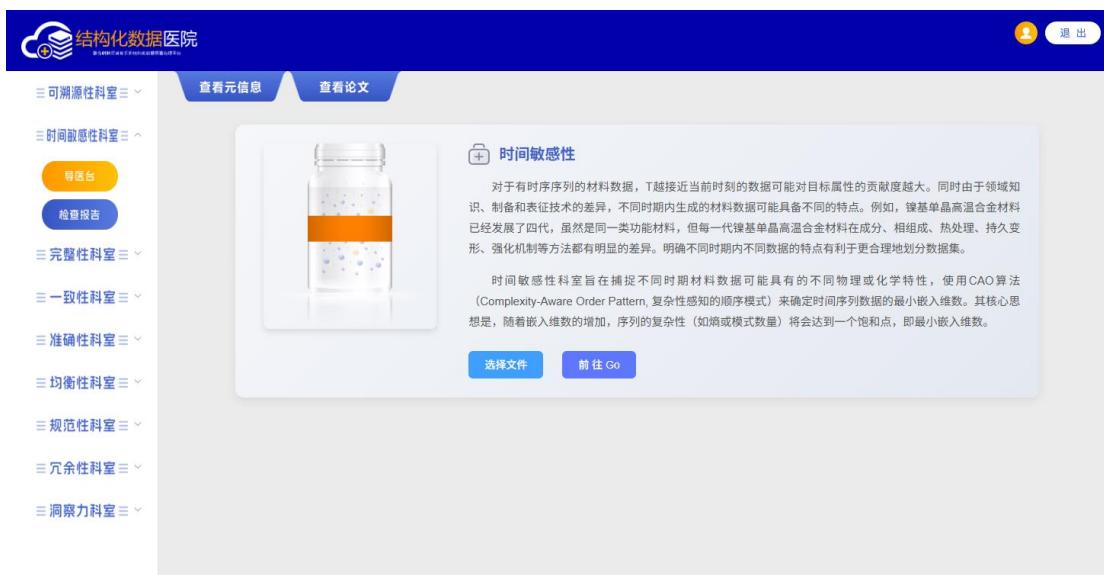


图 4-2 数据诊治时间敏感性导医台

(2) 点击“选择文件”，选择数据挂号上传的文件，点击“前往 go”，进入到时间敏

感性最小嵌入维计算页面。平台会根据在数据上传时选择的时序特征，通过 CAO 算法计算其最小嵌入维数，并在页面上进行显示，运行结果如下图 4-3 所示。

**时间敏感性**

**检查报告说明**

本页面展示的检查报告是基于您提供的时间序列数据生成的特征样本，每个表格对应一个特定的特征，展示了该特征的时间序列数据经过处理后的结果。以下是对表格内容及其生成过程的详细说明：

- 特征与时间序列：页面上的每个标签页对应一个特定的特征，点击即可查看对应特征的样本数据。
- 最小嵌入维数：表格上方显示的“最小嵌入维数”通过 CAO 算法计算得到，表示预测数据所需要的最小嵌入维数，最大值为 10。
- 表格结构：
  - (a) Pre列：表示利用嵌入维数从时间序列中提取的相距的数据点。例如，如果某特征的最小嵌入维数为3，则将显示三个“Pre”列进行数据提取。
  - (b) Next列：表示紧随“Pre”数据段后的一个数据点。例如，如果某特征的最小嵌入维数为3，则表示至少需要前三个数据（“Pre”列）来预测“Next”列的内容。

本页面展示的检查报告旨在提供一种直观的方式来查看和分析时间序列数据的内在结构和特征，为您的研究或项目决策提供参考和支持。

Pre1	Pre2	Pre3	Pre4	Next
8.46	8.729	8.628	8.475	8.766
8.729	8.628	8.475	8.766	9.186
8.628	8.475	8.766	9.186	9.186
8.475	8.766	9.186	9.186	9.186
8.766	9.186	9.186	9.186	8.729

图 4-3 数据诊治时间敏感性体检报告

### 4.3 完整性诊治

(1) 点击平台页面左侧导航栏中的“完整性科室” - “导医台”，进入到完整性文件选择页面，运行结果如下图 4-4 所示。

**完整性**

元数据的完整性治理可以通过预先定义好元数据组织结构和具体内容，并在数据的获取、分析和处理的过程中实时监督元数据的生成和记录来实现。主数据的完整性治理则可以利用更灵活的工具实现。目前，部分机器学习模型可学习含有空缺值的数据，如人工神经网络、决策树和支持向量机。空缺值插补是另一种主数据完整性的治理途径，如 KNN、模糊 C 均值法和自组织神经网络等等。

完整性科室旨在评估关于材料数据本身特征取值的完整性，判断数据集是否存在缺失值，如果存在，则首先计算有缺失的数据缺失率，之后用多重插补方法填补缺失数据，得到完整的数据集，之后分别对缺失数据集和完整数据集进行测试集和训练集划分，训练KNN，MLR模型进行回归检测，并使用RMSE和R<sup>2</sup>两个标准，评估两个模型的精度（在测试集上）。

图 4-4 数据诊治完整性导医台

(2) 点击“选择文件”，进入到选择文件页面，运行结果如下图 4-5 所示。在待分析的数据集中选择数据挂号上传的数据集，点击“分析”，等待平台分析完成之后会将此数据集

加入到已分析的数据集中，运行结果如下图 4-6 所示。在已分析的数据集中点击“报告”，会跳转到数据详情页面，展示说明文档、数据完整性概览、特征相关性矩阵、特征缺失率、Top5 重要特征、筛选条件/决策属性，运行结果如下图 4-7 所示。

待分析的数据集		
Nasicon nasicon 45 aaa 32	test001 12e2e 42 aa 57	test32 sfewfwe data32 qqq 1 40
<a href="#">预览</a> <a href="#">删除</a> <a href="#">分析</a>	<a href="#">预览</a> <a href="#">删除</a> <a href="#">分析</a>	<a href="#">预览</a> <a href="#">删除</a> <a href="#">分析</a>

图 4-5 数据诊治完整性分析数据集

已分析的数据集		
Nasicon nasicon 58 1 41	nasicon nasicon数据 1 1 40	55 1 22
<a href="#">报告</a> <a href="#">分析</a>	<a href="#">报告</a> <a href="#">分析</a>	<a href="#">报告</a>

图 4-6 数据诊治完整性查看已分析数据集



图 4-7 数据诊治完整性分析结果

(3) 点击平台页面左侧导航栏中的“完整性科室”-“诊室台”，进入到算法填补页面，运行结果如下图 4-8 所示。平台提供双模式填补策略，可以选择推荐算法和手动填补两种方式对数据进行填补，此处以推荐算法为例。勾选待填补数据集前面的方框，点击“生成推荐算法”，平台会根据领域知识自动生成适合该数据集的最佳填补算法。再次勾选待填补数据集前面的方框，点击“一键填补”，平台对待填补数据集进行填补，运行结果如下图 4-9 所示。

数据集名称	数据集规模	特征属性规模	数据集描述	数据集状态
Nasicon	90	45	nasicon	完整性良好无需处理
nasicon	90	45	nasicon数据	完整性良好无需处理
58	413	7	1	完整性良好无需处理
57	1016	81	1	完整性良好无需处理
56	38	20	1	完整性良好无需处理

图 4-8 数据诊治完整性诊室台

数据集名称	数据集规模	特征属性规模	数据集描述	数据集状态
Nasicon	90	45	nasicon	完整性良好无需处理
nasicon	90	45	nasicon数据	完整性良好无需处理
58	413	7	1	完整性良好无需处理
57	1016	81	1	完整性良好无需处理
56	36	20	1	完整性良好无需处理

推荐算法: Nasicon, MissForest  
生成推荐算法  
手动选择填补算法: 统计学算法, 机器学习算法  
一键填补  
手动填补

图 4-9 数据诊治完整性诊治结果

(4) 点击平台页面左侧导航栏中的“完整性科室” - “检查报告”，进入到填补效果评估页面，运行结果如下图 4-10 所示。在“选择数据集”中选择刚刚填补后的数据集，在“选择回归模型”中根据需求选择相应的模型，在“选择评估指标”中根据需求选择相应的指标，点击“开始评估”，平台会根据用户选择进行填补效果评估，并在“评估结果”区域可视化呈现填补效果，运行结果如下图 4-11 所示。

数据集	回归模型	RMSE	MAPE	R2	操作
4	SVR	23.6275	13.1375%	0.5561	<span>立即评估</span> <span>标记可量化</span>
3	LASSO	0.2879	14.6152%	0.1689	<span>立即评估</span> <span>标记可量化</span>
2	LASSO	89.3316	7.3569%	0.8955	<span>立即评估</span> <span>标记可量化</span>
2	LASSO	89.3316	7.3569%	0.8955	<span>立即评估</span> <span>标记可量化</span>
1	Ridge	2.0644	22.5824%	0.5624	<span>立即评估</span> <span>标记可量化</span>

图 4-10 数据诊治完整性检测报告页面

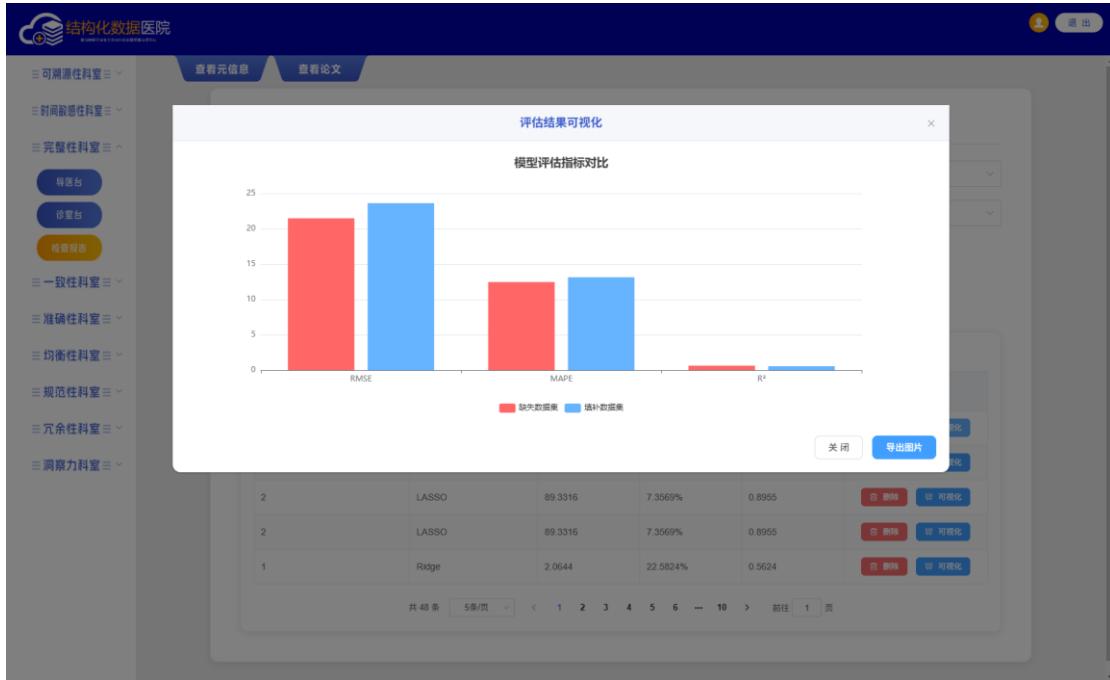


图 4-11 数据诊治完整性评估结果

## 4.4 一致性诊治

(1) 点击平台页面左侧导航栏中的“一致性科室” - “导医台”，进入到一致性文件选择页面，在一致性治理中平台实现了对于数据量纲不一致和表现形式不一致的分别治理，运行结果如下图 4-12 所示。



图 4-12 数据诊治一致性导医台

(2) 点击“量纲不一致”区域的“选择文件”，选择数据挂号上传的文件，点击“前往”，进入到量纲不一致治理页面，运行结果如下图 4-13 所示。在“单位不一致判断”中根据自身是否需要治理的需求选择“是”或者“否”，此处以需要治理为例，选择“是”，平

台会在“单位一致性诊断报告单”区域展示存在单位问题的数据，并在“单位一致化结果”区域展示治理后的数据集，运行结果如下图 4-14 所示。

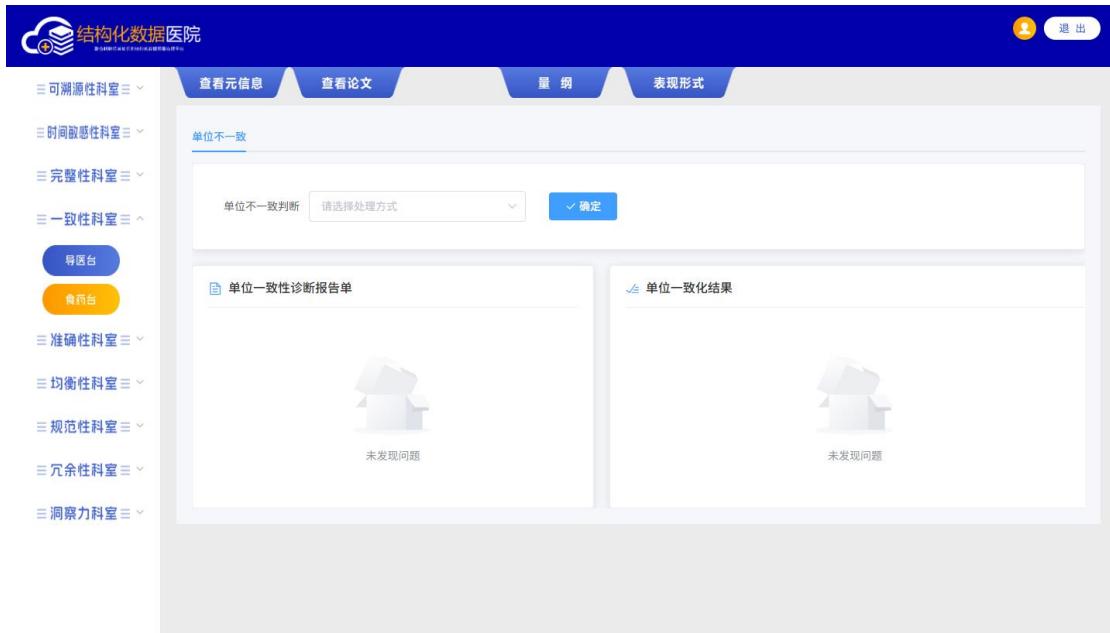


图 4-13 数据诊治量纲一致性食药台

Occu_6b	Occu_18e	Occu_36f	C_Na	Occu_M1	Occu_18
0.5	0.5	0	12	0.75	0
1	0.83	0	20.94	0.75	0
0.78	0.407	0	12.006	0.5	0
1	0.087	0	7.566	0.152	0
1	0.031	0	6.558	0.046	0
1	1	0	24	1	0
共 90 条 10条/页 < 1 2 3 4 5 6 ... 9 > 前往 1					

图 4-14 数据诊治量纲一致性诊治结果

(3) 点击平台页面上方的“表现形式”，进入到表现形式不一致治理页面，运行结果如下图 4-15 所示。在“表现形式不一致判断”中根据自身是否需要治理的需求选择“是”或者“否”，此处以需要治理为例，选择“是”，平台会在“表现形式一致性诊断报告单”区域展示存在表现形式问题的数据，并在“表现形式一致化结果”区域展示治理后的数据集，运行结果如下图 4-16 所示。

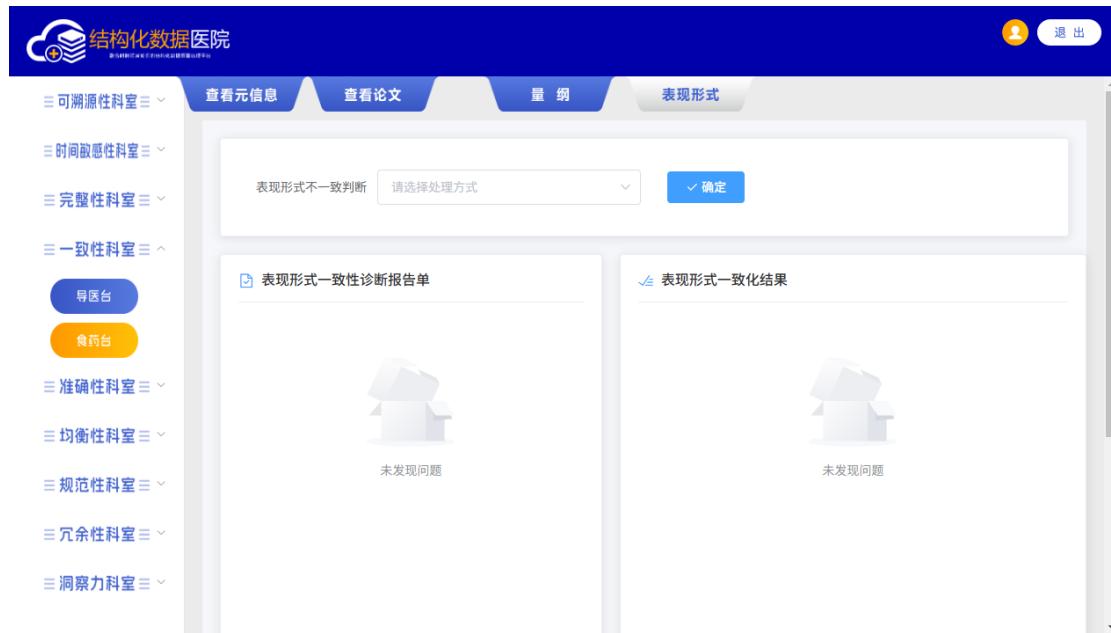


图 4-15 数据诊治表现形式一致性食药台

图 4-16 数据诊治量纲一致性诊治结果

## 4.5 准确性诊治

(1) 点击平台页面左侧导航栏中的“准确性科室” - “导医台”，进入到准确性文件选择页面，运行结果如下图 4-17 所示。



图 4-17 数据诊治准确性导医台

(2) 点击“选择文件”，选择数据挂号上传的文件，页面会展示当前数据的单描述符规则和相似性识别规则（如果没有立即显示，请等待几秒钟后刷新页面），用户可根据领域知识对其进行修改，运行结果如下图 4-18 所示。点击“前往 Go”，进入到数据准确性分维度检测页面，运行结果如下图 4-19 所示。

	A	B	C	D	E	F
1	1	Occu_6b	Float	[0, 6]		
2	2	Occu_18e	Float	[0, 6]		
3	3	Occu_36f	Float	[0, 5.98]		
4	4	C_Na	Float	[0, 29]		
5	5	Occu_M1	Float	[0, 6]		
6	6	Occu_M2	Float	[0, 5.95]		
7	7	EN_M1	Float	[0, 7.16]		
8	8	EN_M2	Float	[0, 7.55]		
9	9	avg_EN_M	Float	[0, 7.35]		
10	10	Radius_M1	Float	[0, 5.88]		
11	11	Radius_M2	Float	[0, 6.01]		
12	12	avg_Radius_M	Float	[0, 5.88]		
13	13	Valence_M1	Float	[0, 9.25]		
14	14	Valence_M2	Float	[0, 11]		
15	15	avg_Valence_M	Float	[0, 9.36]		
16	16	Occu_X1	Float	[0, 6]		
17	17	Occu_X2	Float	[0, 5.97]		
18	18	EN_X1	Float	[0, 7.16]		
19						

	A	B	C	D	E	F
1	Occu_6b	Occu_18e	Occu_36f	C_Na	Occu_M1	
2	0	0	0	0	0	
3	0	0	0	0	0	
4	0	0	0	0	0	
5	0	0	0	0	0	
6	0	0	0	0	0	
7	0	0	0	0	0	
8	0	0	0	0	0	
9	0	0	0	0	0	
10	0	0	0	0	0	
11	0	0	0	0	0	
12	0	0	0	0	0	
13	0	0	0	0	0	
14	0	0	0	0	0	
15	0	0	0	0	0	
16	0	0	0	0	0	
17	0	0	0	0	0	
18	0	0	0	0	0	
19						

图 4-18 数据诊治准确性导医台规则



图 4-19 数据诊治准确性诊室台

(3) 点击“单维度”区域的“检测”，平台会根据异常检测算法和所修改的领域知识对异常值进行检测并展示在页面上，运行结果如下图 4-20 所示。点击“多维度”区域的“检测”，平台会根据异常检测算法和所修改的领域知识对异常特征对进行检测并展示在页面上，运行结果如下图 4-21 所示。点击“全维度”区域的“检测”，平台会根据异常检测算法和所修改的领域知识对异常样本进行检测并展示在页面上，运行结果如下图 4-22 所示。

最终异常值检测结果														
箱线图异常值检测结果														
3sigema异常值检测结果														
Z-score异常值检测结果														
Occup_6b	Occup_16	Occup_36f	C_Na	Occup_M1	Occup_M2	EN_M1	EN_M2	avg_EN_M	Radius_M1	Radius_M2	avg_Radi_us_M	Valence_M1	Valence_M2	avg_Val_e_M
36	29	3	30	30								37	37	
	38	4		56								38	38	
	39	29										39	39	
	40	50										40	40	
	41											41	41	
	42											42	42	
	43											68	68	
	77											69	69	
	78											70	70	
												71	71	

图 4-20 数据诊治准确性诊室台单维度

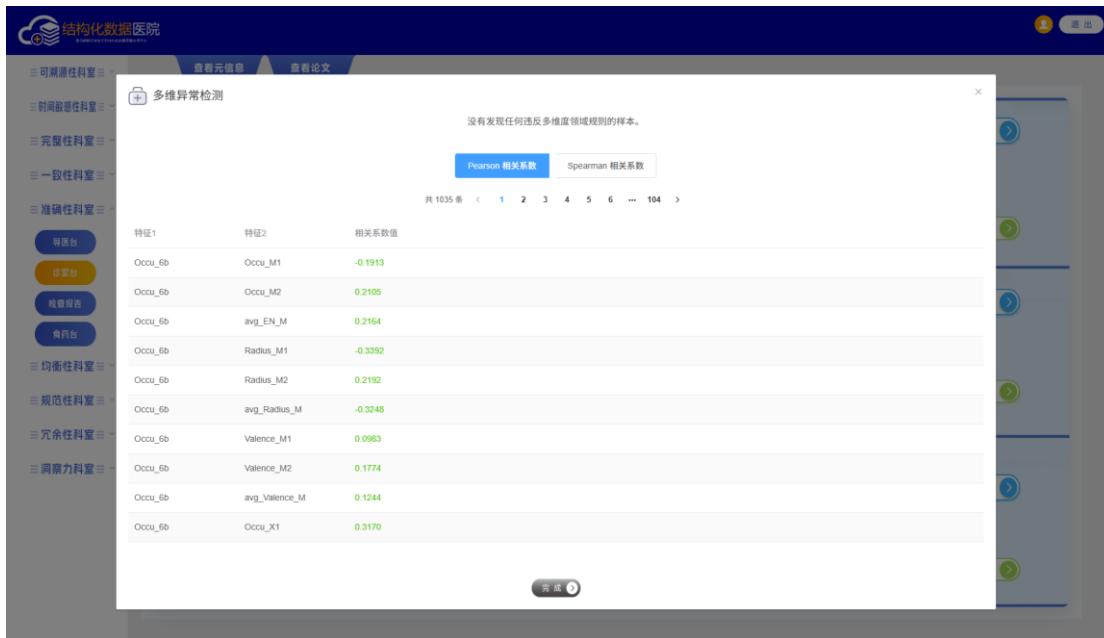


图 4-21 数据诊治准确性诊室台多维度

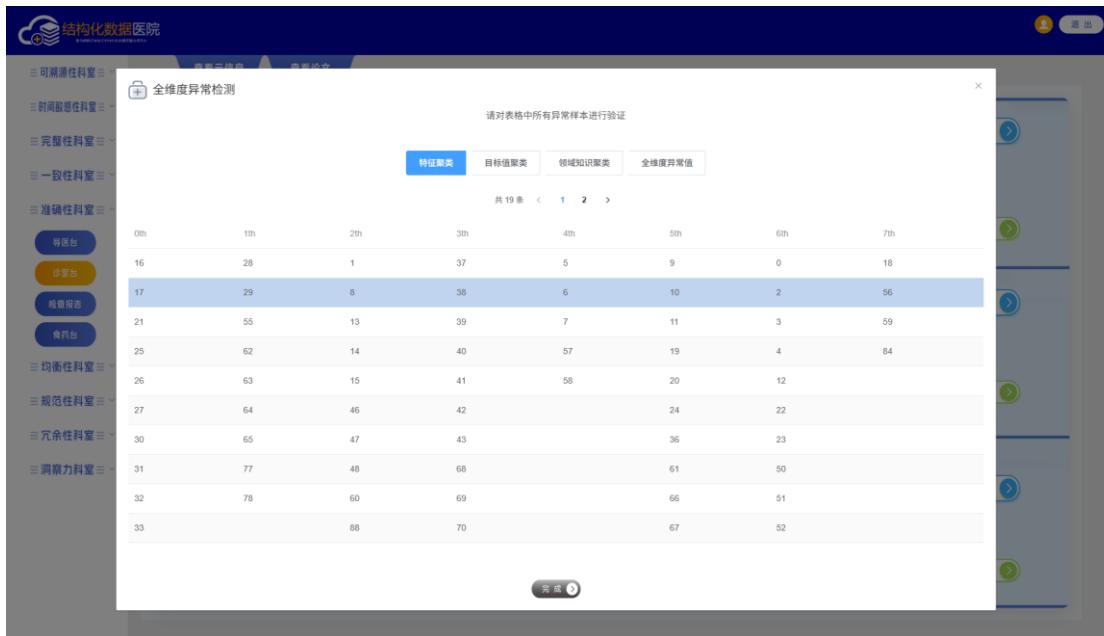


图 4-22 数据诊治准确性诊室台全维度

(4) 点击“单维度”区域的“查看详情”，进入到单维度准确性治理页面，页面会呈现字段准确性以及各特征的异常率等相关单维度准确性指标，并根据检测结果提供治疗策略，运行结果如下图 4-23 所示。点击页面上方的“多维度”，进入到多维度准确性治理页面，页面会呈现特征准确性等相关多维度准确性指标，并根据检测结果提供治疗策略，运行结果如下图 4-24 所示。点击页面上方的“全维度”，进入到全维度准确性治理页面，页面会呈现样本准确性等相关全维度准确性指标，并根据检测结果提供治疗策略，运行结果如下图 4-25 所示。

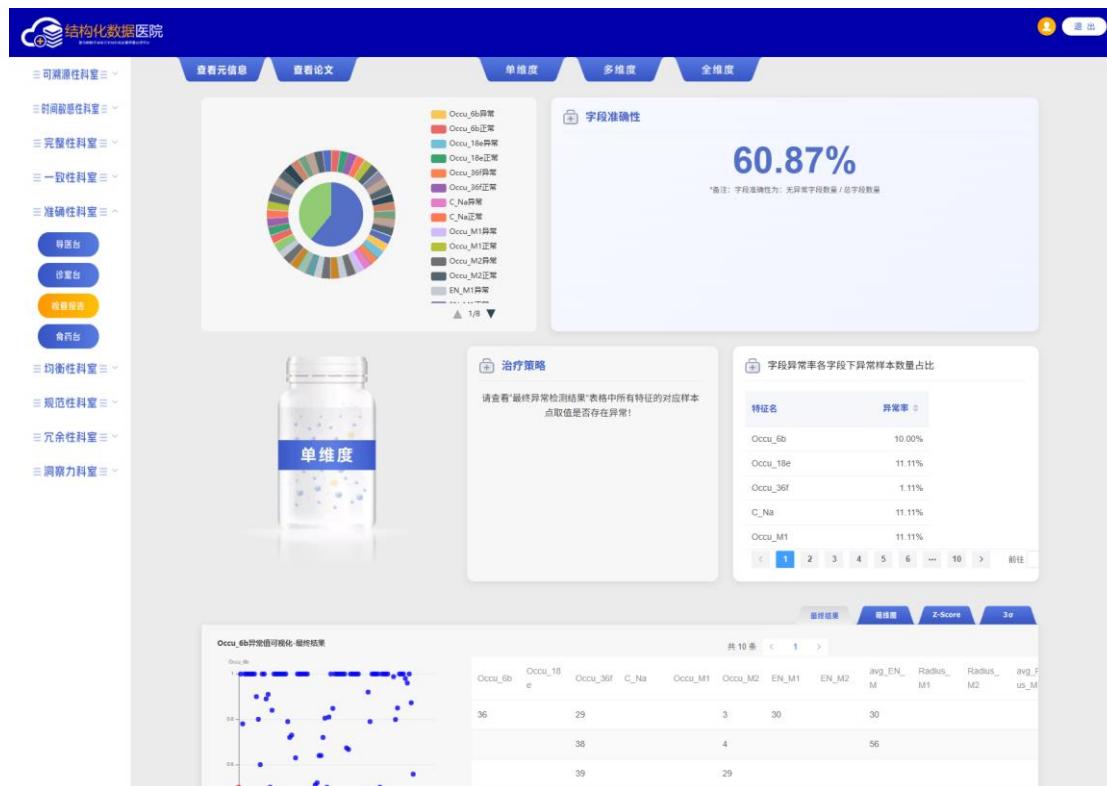


图 4-23 数据诊治准确性检查报告单维度

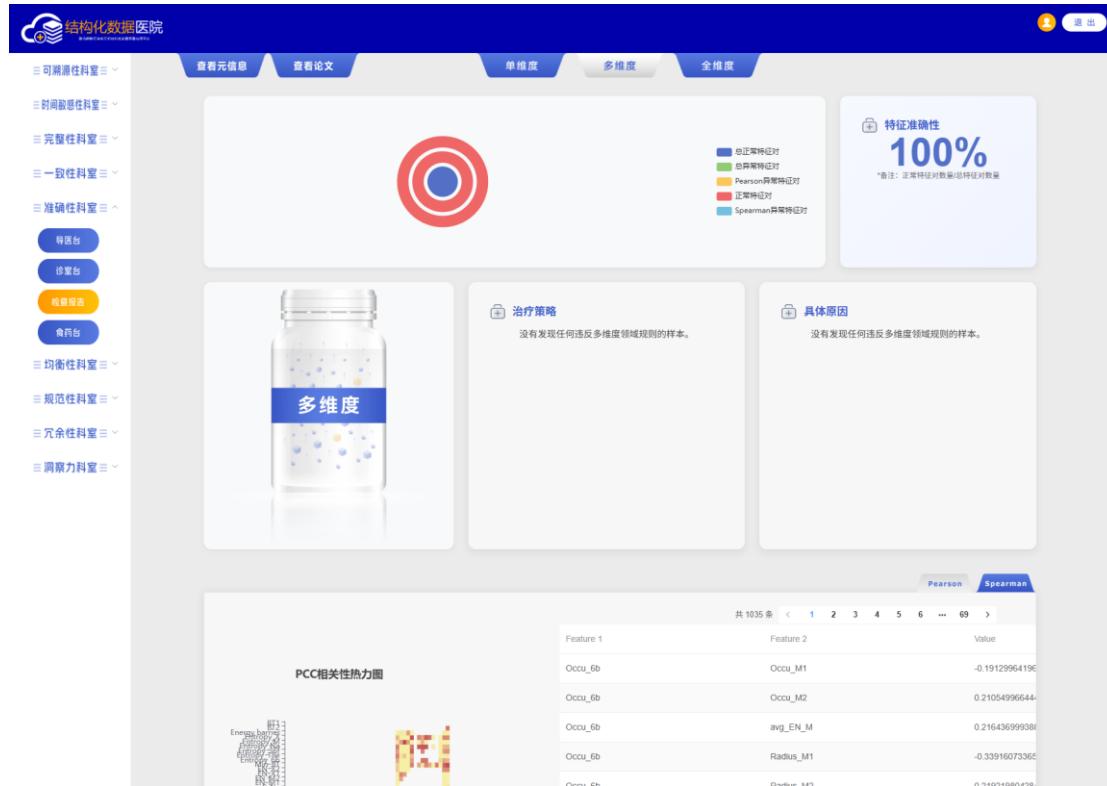


图 4-24 数据诊治准确性检查报告多维度

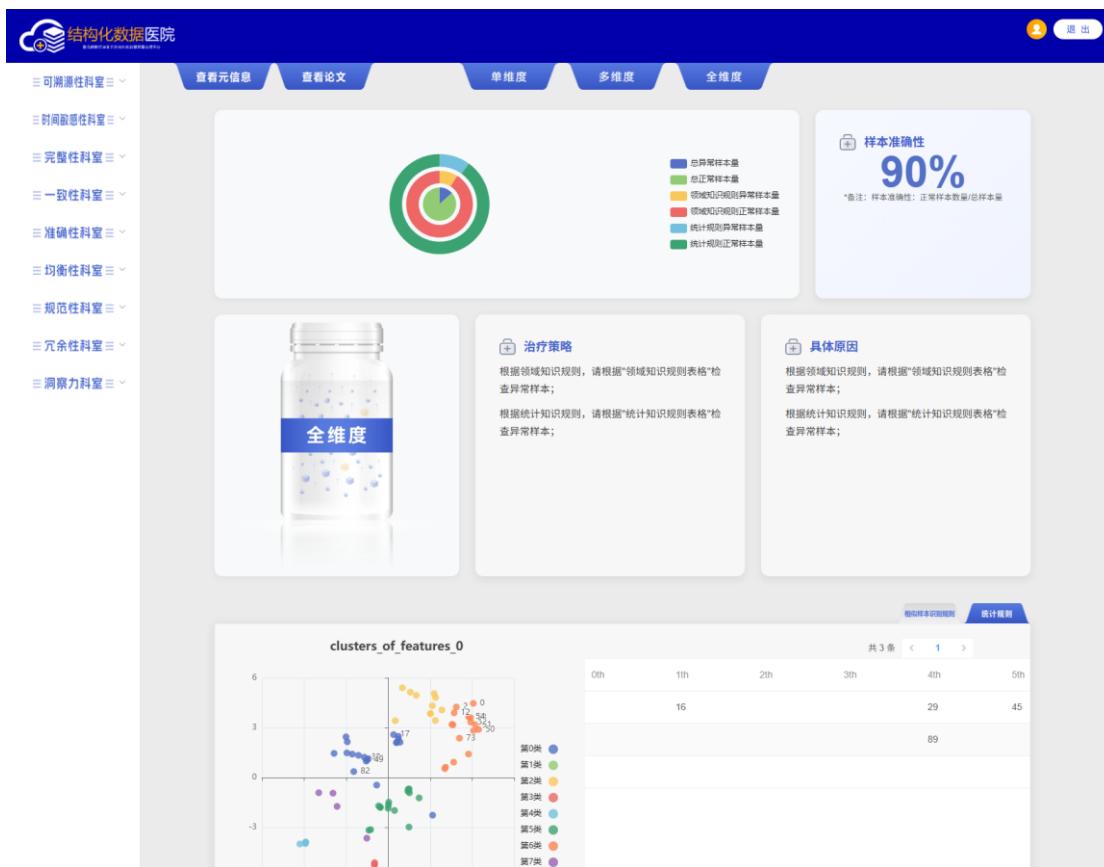


图 4-25 数据诊治准确性检查报告全维度

(5) 点击平台页面左侧导航栏中的“准确性科室” - “食药台”，进入到数据准确性治理页面，页面鼠标各个维度的食药瓶上会自动呈现需要治理的样本。点击“一键修正”平台会对有问题的数据进行治理，运行结果如下图 4-26 所示。

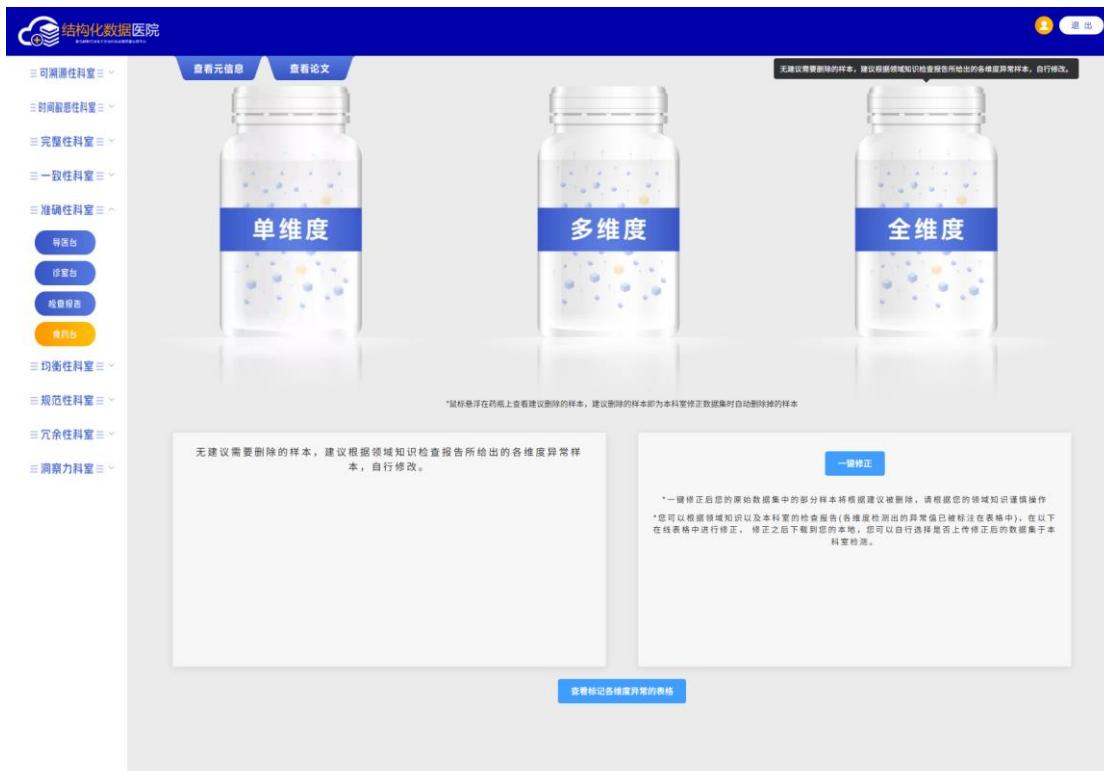


图 4-26 数据诊治准确性食药台

## 4.6 均衡性诊治

(1) 点击平台页面左侧导航栏中的“均衡性科室” - “导医台”，进入到均衡性文件选择页面。点击“选择文件”，选择数据挂号上传的文件，点击“均衡性检测”，平台会计算选择数据的均衡性指标并进行显示，运行结果如下图 4-26 所示。



图 4-26 数据诊治均衡性导医台

(3) 点击“前往治理”，平台进入均衡性治理页面，可根据是否需要对样本进行均衡性治理的需求选择“不对样本做修改”或者“采样”，此处以“采样”为例进行测试。点击“采样” - “确定”，平台根据均衡性指标判断如何进行治理并进行治理，随后点击“选择文件”，选择刚刚治理后生成的文件，点击“再次均衡性检测”，平台会展示治理后的均衡性指标，运行结果如下图 4-27 所示。

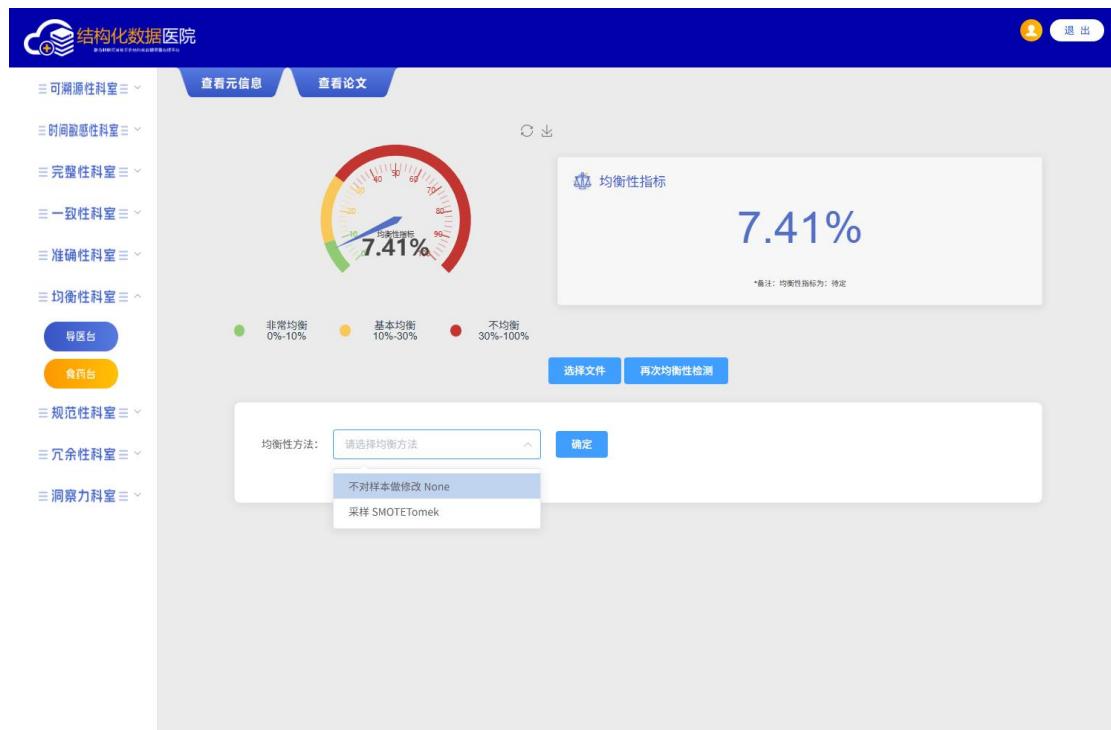


图 4-26 数据诊治均衡性食药台

## 4.7 规范性诊治

(1) 点击平台页面左侧导航栏中的“规范性科室” - “导医台”，进入到规范性文件选择页面，在规范性治理中平台实现了对于取值规范化和划分规范化的分别治理，运行结果如下图 4-28 所示。



图 4-26 数据诊治规范性导医台

(2) 点击“取值规范化”区域的“选择文件”，选择数据挂号上传的文件，点击“前往 Go”，进入到取值规范化治理页面，运行结果如下图 4-27 所示。在“规范化方法”中根据自身需求选择不同的规范化方法，此处以最小-最大规范化为例。选择“最小-最大规范化”，点击“确定”，平台会展示治理后的数据集，运行结果如下图 4-28 所示。

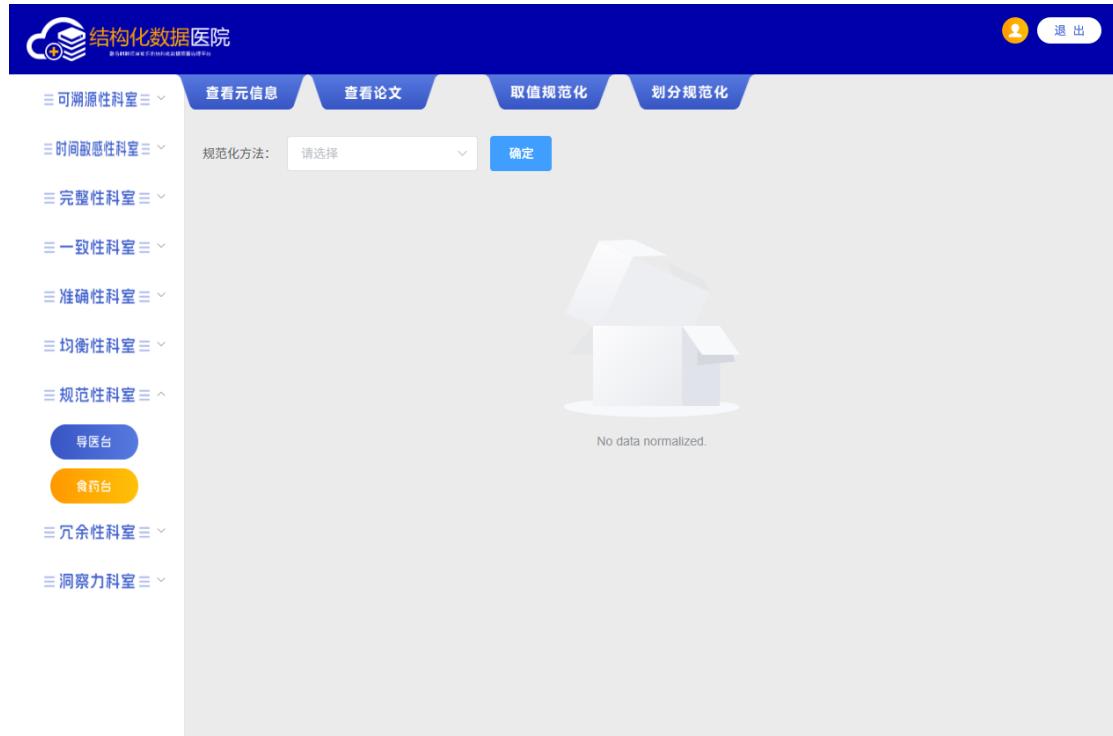


图 4-27 数据诊治规范性取值规范化食药台

Occu_6b	Occu_18e	Occu_36f	C_Na	Occu_M1	Occu_M2	EN_M1	EN_M2	avg_EN_M	Radius_M1	Radius_M2	avg_Radius_M	VN
0.5	0.5	0	0.4285714286	0.7379454927	0.2620545073	0.3604651163	0.5882352941	0.2698564593	0	0.7722772277	0	0.4076086957
1	0.83	0	0.8542857143	0.7379454927	0.2620545073	0.4186046512	0.737254902	0.3942583732	0.380952381	0.8292079208	0.4076086957	0.4076086957
0.78	0.407	0	0.4288571429	0.4758909653	0.5241090147	0.6162790698	0.768627451	0.5684210526	0.2857142857	0.8217821782	0.4456521739	0.4456521739
1	0.087	0	0.2174285714	0.1111111111	0.8888888889	0.3604651163	0.6039215686	0.2392344498	0	0.7376237624	0.1793478261	0.1793478261
1	0.031	0	0.1694285714	0	1	0.3604651163	0.5098039216	0.0124401914	0	0.8415841584	0.7663043478	0.7663043478
1	1	0	1	1	0	0.0348837209	0	0.028708134	0.880952381	0	0.8641304348	0.8641304348

共 90 条 10条/页 < 1 2 3 4 5 6 ... 9 > 前往 1 页

图 4-28 数据诊治规范性取值规范化治理结果

(3) 点击平台页面上方的“划分规范化”，进入到划分规范化治理页面，运行结果如下图 4-29 所示。在“规范化方法”中根据自身需求选择对应的划分方法，此处以留一法为例进行测试。选择“留一法”，点击“确定”，平台会在“训练集”、“测试集”和“验证集”区域展示划分后的数据，运行结果如下图 4-30 所示。

方法	说明	方法推荐
留出法 (Hold-Out)	适用规则：适合于各种规模的数据集，特别是中到大规模数据集。 描述：将数据集分成两部分，通常一部分用于训练模型，另一部分用作测试集评估模型性能。有时还会划分出第三部分作为验证集，用于模型的选择和调参。 优点：实现简单，计算成本较低。 缺点：结果可能受到数据划分方式的影响较大，尤其是在数据集较小的情况下。	数据集规模定义： 小规模数据集：少于100个样本 中规模数据集：100到1000个样本 大规模数据集：超过1000个样本 方法建议： 小规模数据集：使用留一法 中规模数据集：使用k折交叉验证 大规模数据集：使用留出法  数据集样本数：90 属于：小规模数据集 <b>推荐使用留一法。</b>
k折交叉验证 (k-Fold Cross-Validation)	适用规则：适用于小到中等规模的数据集。 描述：将数据集分成k个大小相似的互斥子集，每个子集尽量保持数据分布的一致性。每次用k-1个子集的合集作为训练集，剩下的一个子集作为测试集，这个过程重复进行k次，每次选择不同的子集作为测试集，最后用这k次的平均测试结果作为模型性能的评估。 优点：评估结果更加稳健和可靠，可以有效避免模型过拟合。 缺点：计算成本高，尤其是当k值较大或数据集很大时。	
留一法 (Leave-One-Out, LOO)	适用规则：适用于非常小的数据集。 描述：是k折交叉验证的一种特例，其中k等于样本总数。这意味着如果有N个样本，就将每个样本单独作为测试集，其余N-1个样本作为训练集。因此，这个过程会重复N次，每次都用一个不同的测试样本。 优点：可以充分利用有限的数据，评估结果非常详尽。 缺点：计算成本极高，对于稍大的数据集几乎不可行。	

**Data Input**

规范化方法：请选择

图 4-29 数据诊治规范性划分规范化食药台

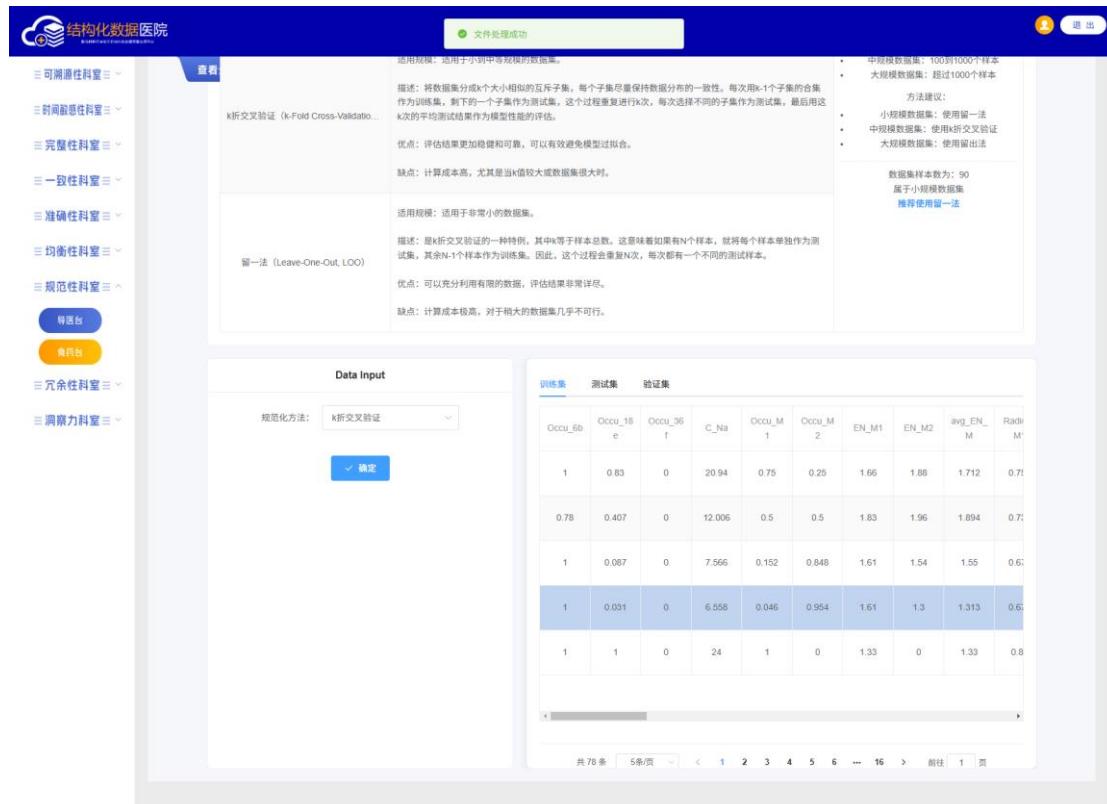


图 4-30 数据诊治规范性划分规范化治理结果

## 4.8 冗余性诊治

(1) 点击平台页面左侧导航栏中的“冗余性科室”-“检查报告”，进入到冗余性数据管理界面，平台实现了对于所有上传数据集的查询、查看、下载等功能，运行结果如下图 4-31 所示。

	Data Set Summary	Data Submitter	Data Submission Unit	Data Validator	File Name	Keyword	Sample Size	Dimension	Operation
<input type="checkbox"/>	Nasicon	qs	shu	qs	material_dataset029.xlsx		90	45	<span>View</span> <span>Download</span>
<input type="checkbox"/>	nasicon	qs	shu	qs	material_dataset029.xlsx		90	45	<span>View</span> <span>Download</span>
<input type="checkbox"/>	test001	qs	shu	qs	material_dataset001.xlsx		111	11	<span>View</span> <span>Download</span>
<input type="checkbox"/>	test32	qs	shu	qs	material_dataset032.xlsx		111	11	<span>View</span> <span>Download</span>
<input type="checkbox"/>	45	qs	shu	qs	material_dataset045.xlsx		511	498	<span>View</span> <span>Download</span>

共 516 条 5条/页 < 1 2 3 4 5 ... 104 > 前往 1 页

图 4-30 数据诊治冗余性检查报告页面

(2) 点击平台页面左侧导航栏中的“冗余性科室” - “导医台”，进入到冗余性文件选择页面，运行结果如下图 4-31 所示。

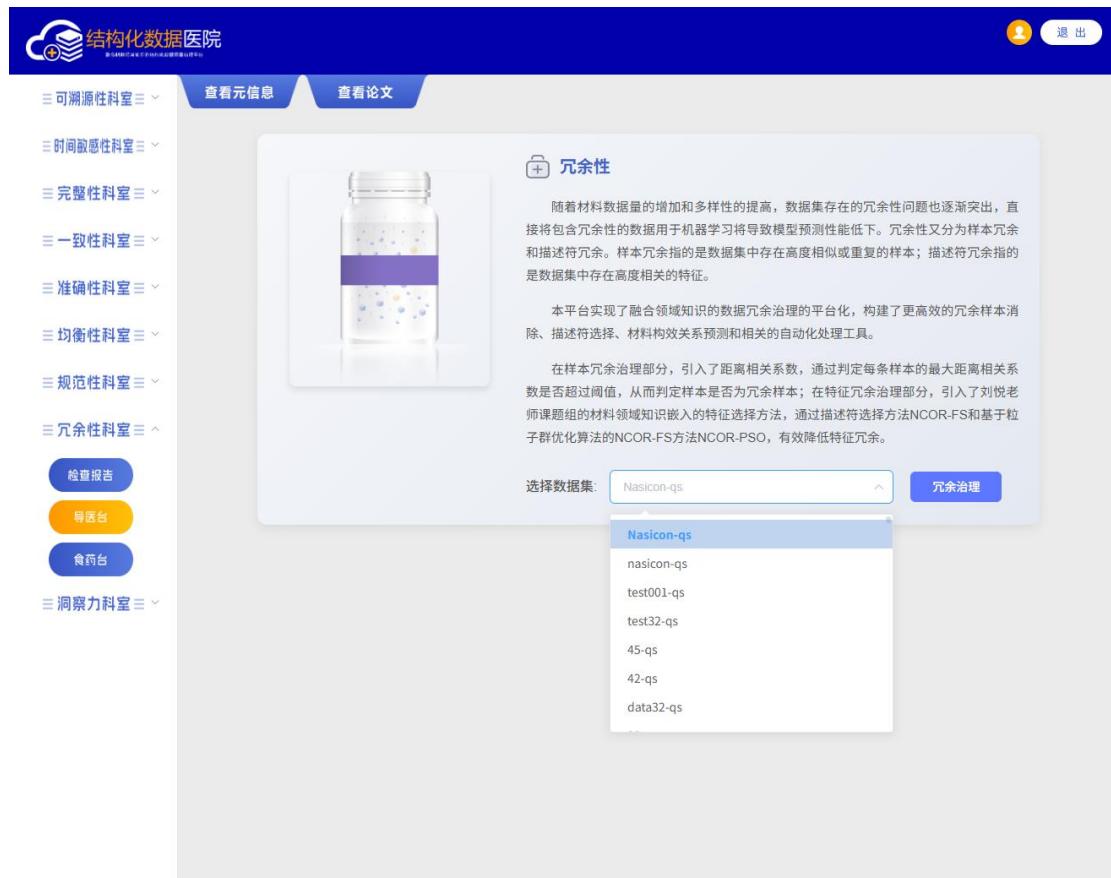


图 4-31 数据诊治冗余性导医台

(3) 在“选择数据集”区域选择需要进行冗余性治理的数据，点击“冗余治理”，进入到数据集评估页面，运行结果如下图 57 所示。点击“下一步”，可在“距离相关系数阈值”区域根据需求输入对应的值，点击“下一步”平台开始计算冗余样本并进行呈现，运行结果如下图 58 所示。可在冗余样本概况区域根据自身需求对冗余样本进行勾选，并点击“删除选中项”对其进行删除，删除之后点击“前往可视化分析”，进入到冗余样本治理后的详情展示页面，运行结果如下图 4-32 所示。

样本冗余治理																
数据集评估		样本冗余治理		数据概况		可视化分析										
数据集概况		距离相关系数		冗余样本数		样本特征数										
Occu_6k	Occu_18k	Occu_24k	O_Na	Occu_M1	Occu_M2	EN_M1	EN_M2	avg_EN_M	Radius_M1	Radius_M2	avg_Radius_M	Valence_M1	Valence_M2	avg_Valence_M	Occu_X1	
0.5	0.5	0	12	0.75	0.25	1.61	1.5	1.9200000000000001	0.6750000000000001	0.78	0.7050000000000001	999999	3	5	3.5	1
1	0.83	0	20.94	0.75	0.25	1.66	1.68	1.712	0.715	0.8375000000000001	0.7760000000000001	000002	3	2	2.75	1
0.78	0.4000000000000001	0	12.006	0.5	0.5	1.83	1.96	1.8300000000000001	0.7349999999999999	0.83	0.7800000000000001	999999	3	4	3.5	1
1	6.0000000000000005	0	7.8000000000000005	0.152	0.8479999999999999	1.61	1.54	1.55	0.6750000000000001	0.7310000000000001	0.745	0.7349999999999999	3	4	3.8479999999999995	1
1	3.1E-2	0	6.5579999999999995	4.0399999999999995	0.36399999999999995	1.61	1.53	1.3129999999999999	0.6750000000000001	0.85	0.8419999999999999	999997	3	4	3.9540000000000002	1
1	1	0	24	1	0	1.33	0	1.33	0.86	0	0.86	4	0	4	4	1
1	1	0	24	1	0	1.33	0	1.33	0.86	0	0.86	4	0	4	4	1
1	0.83	0	20.94	0.75	0.25	1.66	1.68	1.712	0.715	0.8375000000000001	0.7760000000000001	00002	3	2	2.75	1

图 4-32 数据诊治冗余性食药台样本冗余

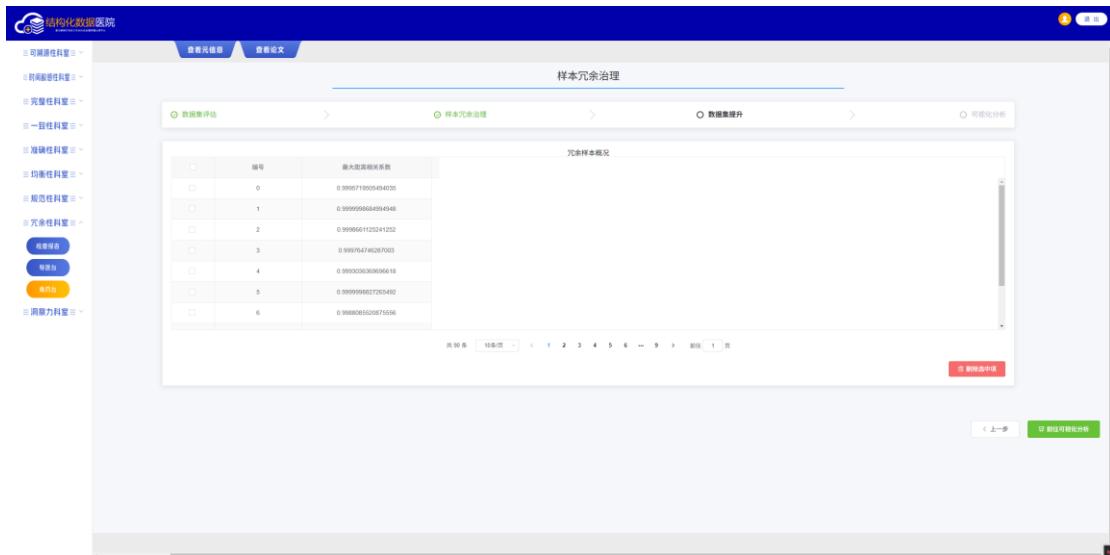


图 4-33 数据诊治冗余性食药台样本冗余

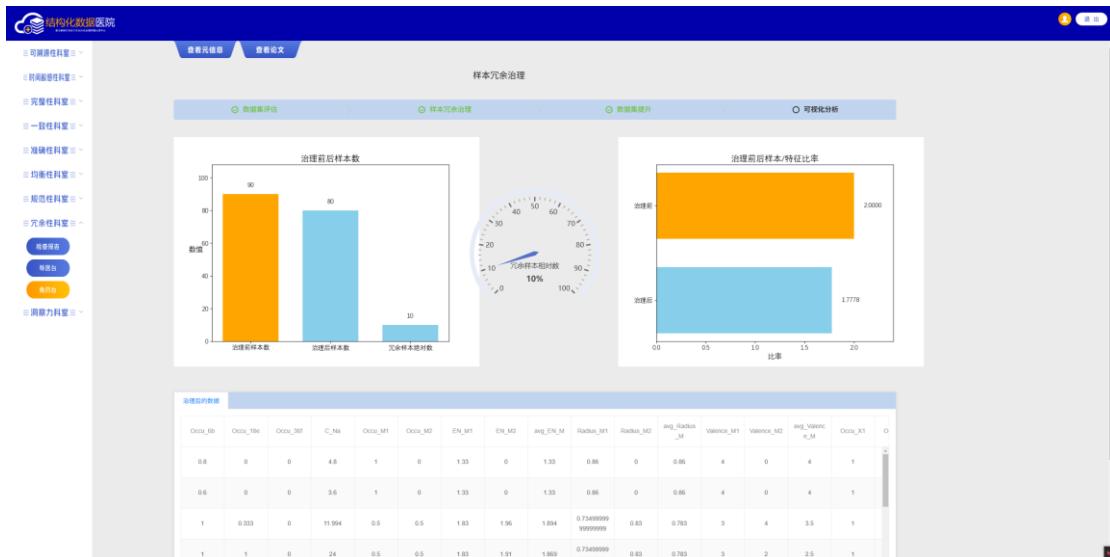


图 4-34 数据诊治冗余性食药台样本冗余

(4) 点击“进入特征冗余治理”，进入到数据集评估页面，运行结果如下图 60 所示。点击“下一步”进入到 KNCOR 上传页面，可根据自身需求上传 KNCOR 文件，点击“选择文件”后选取对应文件即可，运行结果如下图 4-35 所示。点击“下一步”，进入到特征子集查看和添加页面，运行结果如下图 4-36 所示。点击“下一步”，进入到计算最优特征子集页面，点击“开始计算”，平台自动计算该数据集最有特征子集并进行呈现，运行结果如下图 4-37 所示。点击“下一步”，进入到冗余特征治理后的详情展示页面，运行结果如下图 4-38 所示。

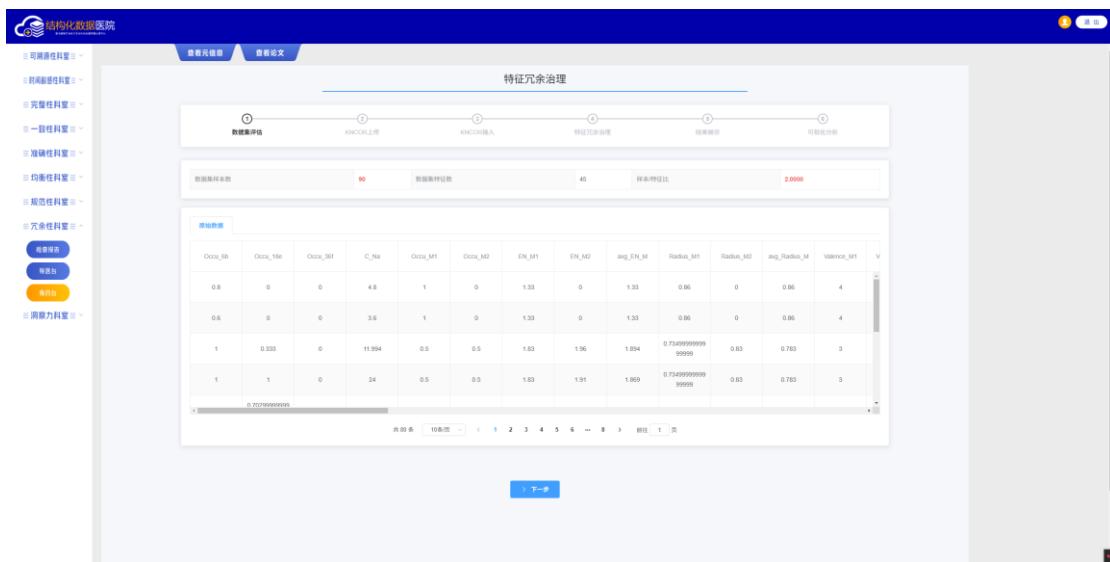


图 4-35 数据诊治冗余性食药台样本冗余



图 4-36 数据诊治冗余性食药台特征冗余

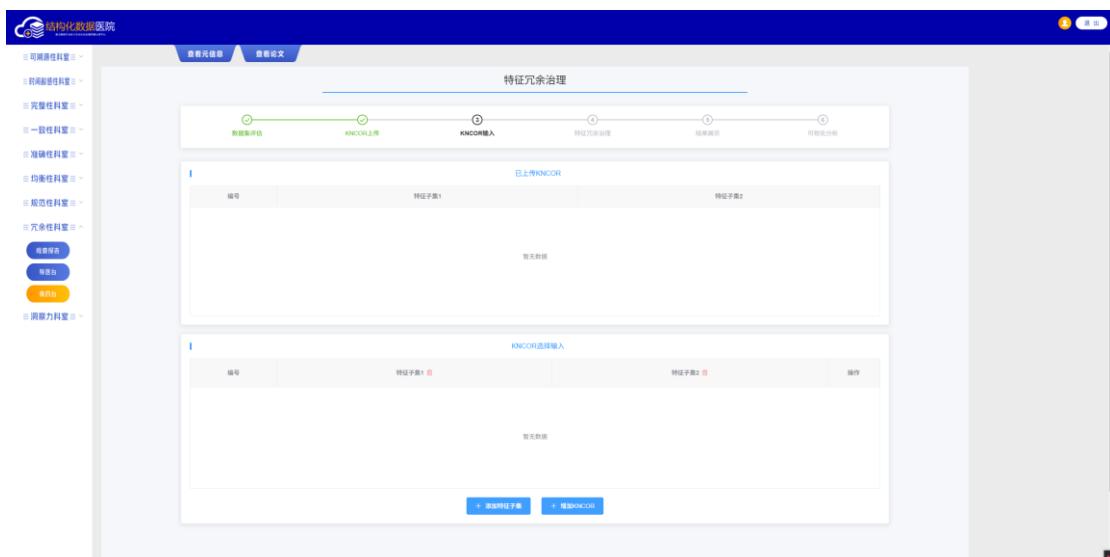


图 4-37 数据诊治冗余性食药台特征冗余

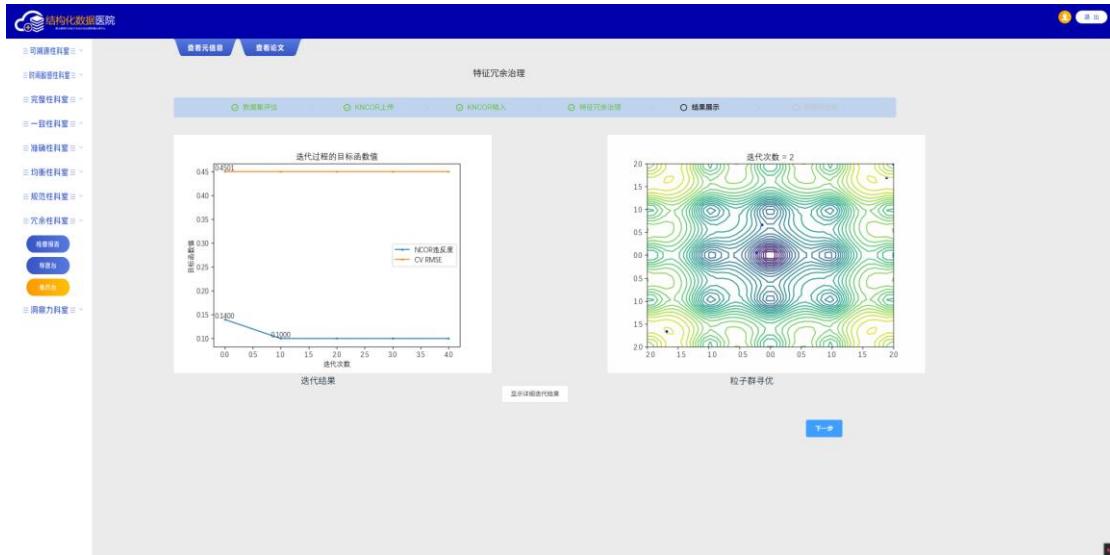


图 4-38 数据诊治冗余性食药台特征冗余

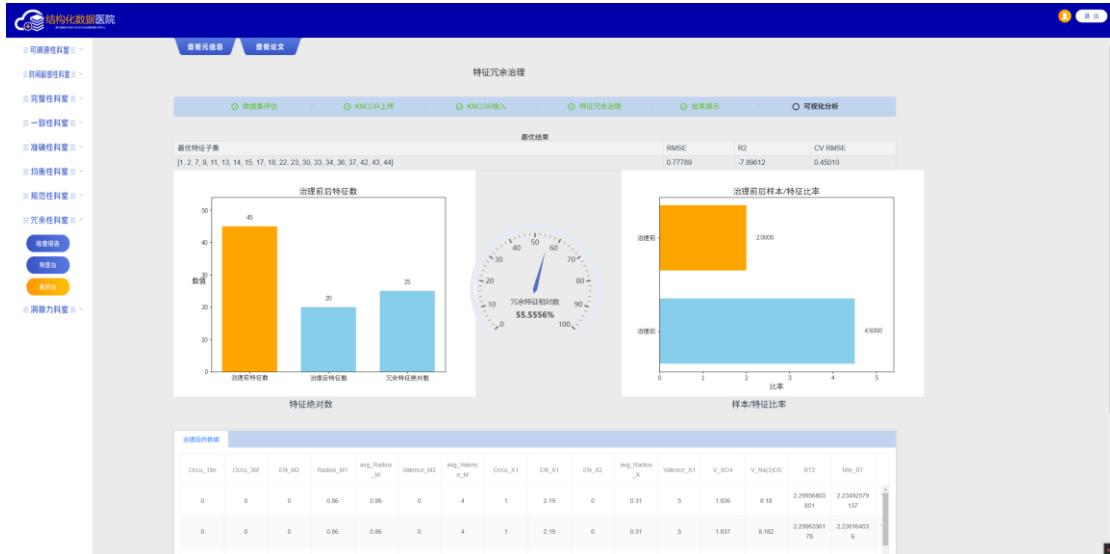


图 4-39 数据诊治冗余性食药台特征冗余

## 4.8 洞察力诊治

(1) 点击平台页面左侧导航栏中的“洞察力科室” - “导医台”，进入到文件和功能选择页面，运行结果如下图 4-40 所示。分别选择划分规范性所生成的训练集和验证集，并根据下游任务选择“自动回归”或“自动聚类”，进入到洞察力检测页面，运行结果如下图 4-41 所示。



图 4-40 数据诊治洞察力导医台

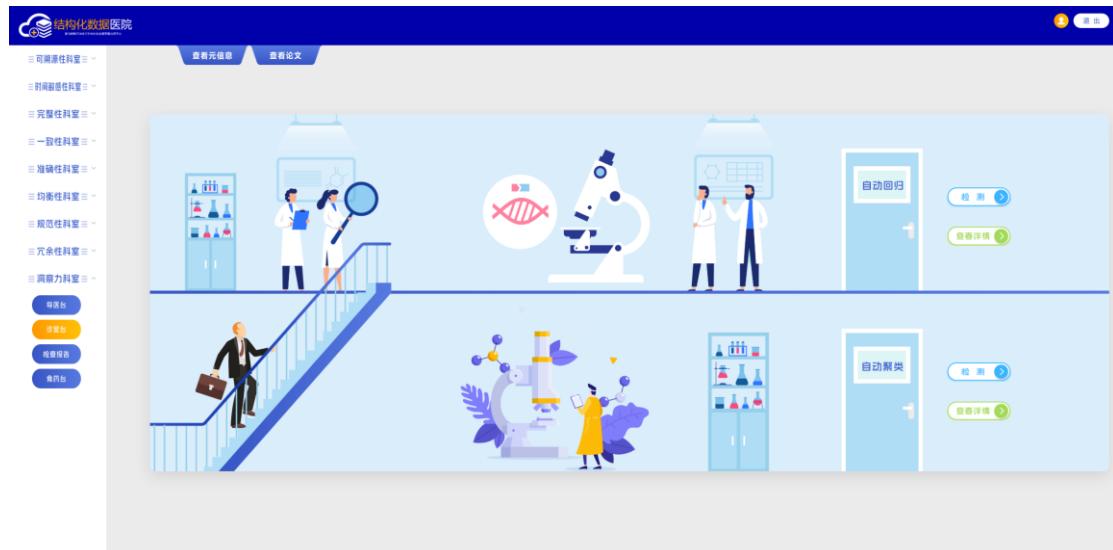


图 4-41 数据诊治诊室台

(2) 点击“自动回归”区域的“检测”，进入到元特征计算页面，运行结果如下图 4-42 所示。点击“元特征计算”，平台开始计算数据的元特征并进行呈现，运行结果如下图 4-43 所示。



图 4-42 数据诊治诊室台元特征计算前

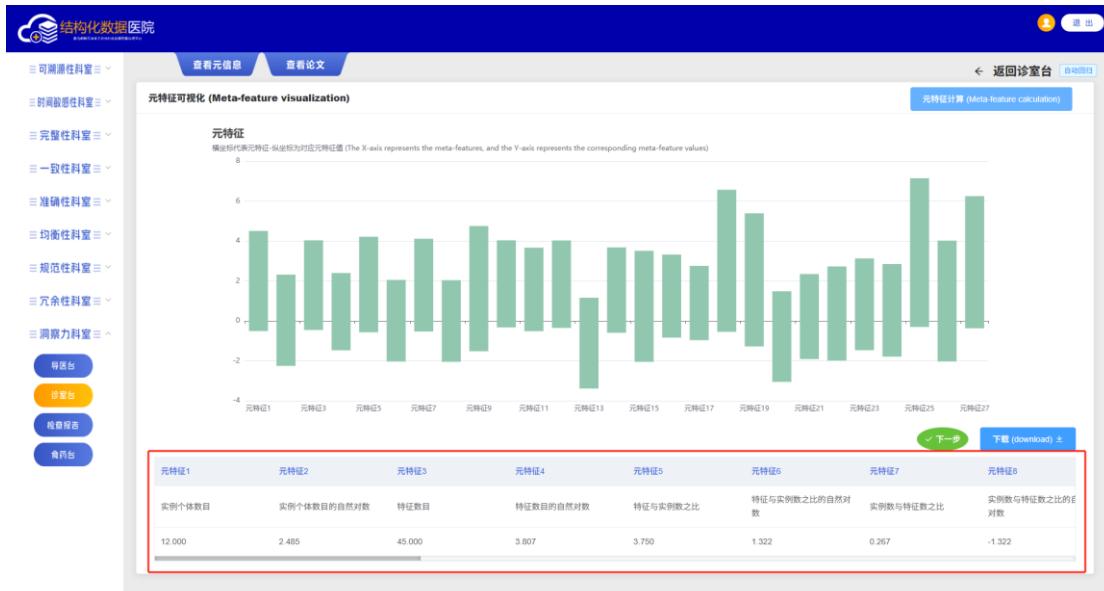


图 4-43 数据诊治诊室台元特征计算后

(3) 点击平台左侧导航栏中的“洞察力科室” - “检查报告”，进入到算法推荐页面，运行结果如下图 4-44 所示。点击“确认推荐”，平台开始计算每个算法的评分进行推荐并展示在页面上，运行结果如下图 4-45 所示。



(4) 点击“下一步”，进入到模型对比效果页面，运行结果如下图 4-46 所示。点击“开始预测”，平台计算在前三种模型上该数据集对应的均方根误差等结果并呈现在页面上，运行结果如下图 4-47 所示。



图 4-46 数据诊治食药台



图 4-47 数据诊治食药台预测后