

分类号：
学号：20212012028

密级：公开
单位代码：10759

石河子大学

硕士学位论文



基于机器学习算法的新疆棉花品质 预测模型研究

学位申请人	陆永迪
指导教师	田景山 副教授 张旺锋 教授
申请学位门类级别	农学硕士
学科、专业名称	作物学
研究方向	数学建模与作物品质生态
所在学院	农学院

中国·新疆·石河子
2024年7月

分类号：
学号：20212012028

密级：公开
单位代码：10759

石河子大学

硕士学位论文



基于机器学习算法的新疆棉花品质 预测模型研究

学位申请人	陆永迪
指导教师	田景山 副教授
	张旺锋 教授
申请学位门类级别	农学硕士
学科、专业名称	作物学
研究方向	数学建模与作物品质生态
所在学院	农学院

中国·新疆·石河子

2024年7月

**Research on prediction model of cotton quality in Xinjiang based on
machine learning**

A Dissertation Submitted to

Shihezi University

In Partial Fulfillment of the Requirements

for the Degree of

Master of Agriculture

By

Lu Yongdi

Mathematical modeling and crop quality ecology

Dissertation Supervisor: Associate Prof. Tian Jing-shan

Prof. Zhang Wang-feng

July, 2024

石河子大学学位论文独创性声明及使用授权声明

学位论文独创性声明

本人所提交的学位论文是在我导师的指导下进行的研究工作及取得的研究成果。据我所知，除文中已经注明引用的内容外，本论文不包含其他个人已经发表或撰写过的研究成果。对本文的研究做出重要贡献的个人和集体，均已在文中作了明确的说明并表示谢意。

研究生签名： 陆永迪 时间： 2024 年 7 月 10 日

使用授权声明

本人完全了解石河子大学有关保留、使用学位论文的规定，学校有权保留学位论文并向国家主管部门或指定机构送交论文的电子版和纸质版。有权将学位论文在学校图书馆保存并允许被查阅。有权自行或许可他人将学位论文编入有关数据库提供检索服务。有权将学位论文的标题和摘要汇编出版。保密的学位论文在解密后适用本规定。

研究生签名： 陆永迪 时间： 2024 年 7 月 10 日

导师签名： 田京山 时间： 2024 年 7 月 10 日

摘 要

【目的】新疆棉区跨越幅度大，纤维品质具有区域性和生态多样性的特点，优良品种配合适宜的气候条件才能取得最佳的纤维品质。如何通过气象变化科学预测纤维品质变化趋势和区域分布，这对棉花产业可持续发展具有重要意义。因此，开展基于机器学习算法的新疆棉花纤维品质预测模型的研究，寻求适宜纤维品质各指标的最优预测模型，这将为棉花区域性育种和栽培提供依据。

【方法】本研究利用新疆历年原棉品质公证检验数据，通过大数据分析气象因素对纤维品质的相对贡献及其关系；采用随机森林(Random Forest)、支持向量回归算法(Support Vector Regression)和长短期记忆网络(Long Short-Term Memory)等机器学习算法构建气象因子与纤维品质的预测模型，筛选出预测效果最优的机器学习算法。

【结果】棉花生育期气象特征变量与原棉纤维品质存在多重共线性，在随机选择最优样本训练集和自变量多重共线性干扰方面，应优先选择随机森林算法。通过不同气象特征变量组合对原棉品质指标的方差解释率高低，可以确定数学模型的最优特征变量个数。不同气象特征变量组合对纤维长度的方差解释率在 10.43%~15.45%，其中，最高气温 T_{\max} 和最低气温 T_{\min} 特征变量组合对纤维品质指标的方差解释率最高、达 15.45%；日照时数 Sun 和降水量 $Prec$ 组合则对马克隆值的方差解释率最高，为 10.64%。输入变量为 T_{\max} 、 T_{\min} 、 Sun 、 $Prec$ 和有效积温 $GDD_{\geq 12^{\circ}\text{C}}$ 组合时，气象特征变量对断裂比强度的方差解释率最高可达 22.34%，所有组合特征变量的方差解释率在 20.06%~22.34%；当平均气温 T_{ave} 、 T_{\max} 、 T_{\min} 、 $GDD_{\geq 10^{\circ}\text{C}}$ 、 $GDD_{\geq 12^{\circ}\text{C}}$ 、 $GDD_{\geq 15^{\circ}\text{C}}$ 、 $GDD_{\geq 20^{\circ}\text{C}}$ 作为输入变量时，对整齐度指数的方差解释率最高为 12.92%。

选择方差解释率较高的气象特征变量组合作为模型输入变量，得到相应的纤维品质指标预测结果。与支持向量机模型相比，随机森林模型预测原棉纤维长度、断裂比强度、马克隆值和整齐度指数的效果更佳，其精度均在 88.59% 以上，均方根误差 RMSE 在 0.0826~0.3192。长短期记忆网络模型预测纤维品质各指标的平均绝对百分比误差 MAPE 在 0.5%~2.6%，RMSE 在 0.021~0.514，决定系数 R^2 范围在 0.020~0.130，其预测误差均较小，对原棉纤维指标的预测性能总体符合预期；但是相较于纤维长度、断裂比强度和整齐度指数的预测，长短期记忆网络模型对马克隆值的预测效果最好，RMSE 值最低为 0.021， R^2 最高为 0.130，可以优先选择长短期记忆网络模型。

【结论】随机森林算法可以更好的随机选择最优样本训练集，并在解决自变量多重共线性方面有较大的优势，通过随机森林算法对自变量进行特征选择会明显提高模型的准确性。随机森林模型预测原棉纤维长度、断裂比强度和整齐度指数的相对性能更好，长短期记忆网络模型则对马克隆值的预测效果最好。

关键词：纤维品质；预测模型；气象；温度；特征选择

Abstract

【Objective】 The cotton region in Xinjiang has a great span, and the fiber quality has the characteristics of regional and ecological diversity. Only by combining excellent varieties with suitable climatic conditions can the best fiber quality be obtained. How to scientifically predict the changing trend and regional distribution of fiber quality through meteorological changes is of great significance to the sustainable development of cotton industry. Therefore, the research of Xinjiang cotton fiber quality prediction model based on machine learning algorithm is carried out, and the optimal prediction model suitable for each index of fiber quality is sought, which will provide basis for cotton regional breeding and cultivation.

【Method】 This study used the data of raw cotton quality notarization inspection in Xinjiang over the years to analyze the relative contribution of meteorological factors to fiber quality and their relationship through big data. Machine learning algorithms such as Random Forest, Support Vector Regression and Long Short-Term Memory were used to build the prediction model of meteorological factors and fiber quality, and the machine learning algorithm with the best prediction effect is screened out.

【Result】 There is multicollinearity between meteorological characteristic variables in cotton growth period and raw cotton fiber quality, so random forest algorithm should be preferred in randomly selecting the optimal sample training set and the multicollinearity interference of independent variables. The optimal number of characteristic variables in the mathematical model can be determined by the variance interpretation rate of raw cotton quality indexes in different combinations of meteorological characteristic variables. The variance explanation rate of fiber length by different combinations of meteorological characteristic variables is 10.43% to 15.45%, and the variance explanation rate of fiber quality index by the combination of maximum temperature T_{max} and minimum temperature T_{min} is the highest, reaching 15.45%. The combination of Sunshine hours Sun and Precipitation $Prec$ has the highest explanation rate of variance of micronaire value, which is 10.64%; When the input variables are the combination of T_{max} , T_{min} , Sun , $Prec$ and effective accumulated temperature $GDD_{s \geq 12^{\circ}C}$, the maximum variance explanation rate of fracture specific strength is 22.34%, and the variance explanation rate of all combined characteristic variables is 20.06% to 22.34%. When the average temperatures T_{ave} , T_{max} , T_{min} , $GDD_{s \geq 10^{\circ}C}$, $GDD_{s \geq 12^{\circ}C}$, $GDD_{s \geq 15^{\circ}C}$ and $GDD_{s \geq 20^{\circ}C}$ are used as input variables, the maximum variance explanation rate of the uniformity index is 12.92%.

The combination of meteorological characteristic variables with high variance explanation rate is selected as the input variables of the model, and the corresponding fiber quality index prediction results are obtained. Compared with the support vector machine model, the random forest model has a better effect in

predicting the length, specific breaking strength, micronaire value and evenness index of raw cotton fiber, with the accuracy above 88.59% and the root mean square error RMSE ranging from 0.0826 to 0.3192. The average absolute percentage error of fiber quality predicted by long-term and short-term memory network model is 0.5% to 2.6%, RMSE is 0.021 to 0.514, and determinant coefficient R^2 is 0.020 to 0.130. The results show that the prediction errors of fiber length, micronaire value, specific strength at break and uniformity index by long-term and short-term memory network model are all small, and the MAPE values are all less than 10%. Therefore, the prediction performance of the long-term and short-term memory network model for each index is generally in line with expectations, but compared with the prediction of fiber length, specific breaking strength and uniformity index, the prediction value of the long-term and short-term memory network is the best relative to the actual value, with the lowest RMSE value of 0.021 and the highest R^2 of 0.130. Compared with the random forest model, the RMSE of the long-term and short-term memory network model for predicting the micronaire value is also lower, so the long-term and short-term memory network model can be preferred for predicting the micronaire value.

【Conclusion】 The random forest algorithm can better randomly select the optimal sample training set, and has great advantages in solving the multicollinearity of independent variables. The feature selection of independent variables through the random forest algorithm will obviously improve the accuracy of the model. The random forest model has better relative performance in predicting the length, specific breaking strength, and evenness index of raw cotton fiber, while the long-term and short-term memory network model has the best prediction effect on the micronaire value.

Key words: Fiber quality; Prediction model; Weather; Temperature; Feature selection

目 录

摘 要.....	I
Abstract.....	II
目 录.....	IV
缩略词表.....	VI
第 1 章 绪论.....	1
1.1 研究背景及意义.....	1
1.2 研究进展.....	2
1.2.1 环境因素对棉纤维品质的影响.....	2
1.2.2 机器学习算法对作物产量与品质的预测.....	3
1.2.3 预测模型构建的特征选择和特征提取.....	4
1.2.4 棉纤维品质预测模型的构建.....	4
1.3 研究内容.....	5
第 2 章 材料与方法.....	7
2.1 研究区域.....	7
2.2 数据来源与计算.....	8
2.2.1 数据来源.....	8
2.2.2 数据计算与特征归一化.....	8
2.3 研究方法.....	9
2.3.1 相关性分析与特征选择.....	9
2.3.2 模型评价.....	9
2.3.3 机器学习回归预测算法.....	10
第 3 章 结果与分析.....	14
3.1 气象因子及品质指标相关性分析和多元共线性诊断.....	14
3.2 影响原棉品质各气象因子的相对重要性排序.....	16
3.3 随机森林和支持向量机回归模型的构建.....	17
3.3.1 纤维长度.....	17
3.3.2 马克隆值.....	20
3.3.3 断裂比强度.....	22
3.3.4 整齐度指数.....	25

3.4 长短期记忆网络模型的构建	27
3.4.1 参数调优	27
3.4.2 模型比较	28
第 4 章 讨论	32
4.1 特征选择对模型精度的影响	32
4.2 参数寻优与模型精度的关系	33
4.3 纤维品质预测模型的精度表现及分析	34
第 5 章 研究结论与展望	36
5.1 研究结论	36
5.2 创新点	37
5.3 研究展望	37
参考文献	39
附 录	48
6.1 随机森林模型参数寻优及回归程序	48
6.2 支持向量机参数寻优及回归程序	54
6.3 长短期记忆网络超参数寻优及回归程序	60
致 谢	65
作者简介	66
石河子大学硕士研究生学位论文导师评阅表	67

缩略词表

英文缩写	英文全称	中文全称
<i>RF</i>	Random Forest	随机森林
<i>SVR</i>	Support Vector Machine Regression	支持向量机回归
<i>LSTM</i>	Long Short-Term Memory	长短期记忆网络
T_{\max}	Daily maximum temperature	最高气温
T_{\min}	Daily minimum temperature	最低气温
T_{ave}	Daily average temperature	平均气温
<i>Sun</i>	Sunshine duration	日照时数
<i>Prec</i>	Precipitation	降水量
<i>DTR</i>	Daily temperature difference	日温差
GDD_s	Growing degree days	有效积温
<i>Len</i>	Upper half mean length	纤维长度
<i>Mic</i>	Micronaire	马克隆值
<i>Str</i>	Strength	断裂比强度
<i>Uni</i>	Uniformity index	整齐度指数
MAPE	Mean Absolute Percentage Error	平均绝对百分比误差
RMSE	Root Mean Square Error	均方根误差

第1章 绪论

1.1 研究背景及意义

棉花纤维是重要的纺织工业原料，优良品质的纤维增强了原棉在国际市场的竞争力，纺织技术的持续发展也要求更好的品质性状的原棉。我国原棉综合品质虽已处于国际中上等水平，但生态区域间存在差异（唐淑荣等，2017）。长江流域棉区纤维长度和断裂比强度较优，黄河流域棉区的则马克隆值偏高，西北内陆棉区马克隆值和纺纱均匀性较好棉区，但断裂比强度偏低（许乃银等，2017）。纤维品质主要受品种遗传特性的影响，环境生态因素也有强烈影响（Darawshen *et al.*, 2010; Echer *et al.*, 2015）。新疆棉花种植区域的生态条件差异大，热量资源的变化形成了纤维品质类型的多样性（Zhu *et al.*, 2023）。可见，通过气象变化科学预测纤维品质变化趋势和区域分布，不仅满足纺织企业用棉的专业性和区域性，还对棉花产业可持续发展具有重要意义。

温度、降水和光照等环境因子对棉纤维品质有显著影响（Li *et al.*, 2020a; Han *et al.*, 2022a）。气候变暖延长了棉花生育期和增加了光合产物积累（张丽霞等，2016；赵彦茜等，2019），日照时数越多则利于纤维品质发育（熊宗伟等，2014），较多的光辐射量增加了光合产物供给和促进优质棉铃形成（Pettigrew *et al.*, 2001）。作物产量与环境因子间的关系通常是非线性的（Jeong *et al.*, 2016; Li *et al.*, 2019），基于机器学习的作物产量预测模型能够构建出多种变量组合影响产量形成的数学模型，较为准确地对进行产量预测（Peng *et al.*, 2020; Dilli *et al.*, 2021）；灰色关联度 GM(1,1) 预测模型可综合评价棉花纤维品质区域分布和预测发展趋势（唐淑荣等，2017）。机器学习预测模型虽已广泛应用于产量预测，但在作物品质预测方面的研究较少，且主要集中于品质时空分布规律的研究。我国棉花加工企业生产的原棉实施包包检测，基于新疆历年原棉公证检验数据，通过大数据分析气象因素对纤维品质的相对贡献及其关系，使用机器学习算法构建气象因子影响纤维品质的预测模型，明确纤维品质预测模型建立的关键步骤及参数的优化调整，筛选出适宜纤维品质各指标的最优预测模型，为棉花区域性育种和生产提供理论支撑。

1.2 研究进展

1.2.1 环境因素对棉纤维品质的影响

棉花优质高产可持续是棉花栽培和育种工作者的长期目标,纤维品质形成受品种遗传特性、环境生态因素、栽培管理措施等诸多因素的综合影响,具有很强的地域性(赵文青等, 2008; Raper *et al.*, 2019),调控纤维品质往往很难综合考虑多个因素的交互作用,且准确预测品质变化规律的难度亦较大(Mohammed *et al.*, 2022)。纤维发育与生态因子的关系已有大量研究,纤维品质性状受多种因素制约,某一性状的变化势必引起其它品质性状的改变(陈兵林等, 2006; Chen *et al.*, 2019)。断裂比强度和马克隆值广义遗传率高,受环境影响相对较小,整齐度指数广义遗传率小,受环境影响大(袁有禄等, 2002)。

棉花花铃期的温度是非常重要的因素,直接影响棉铃的发育状况和产量品质的形成(McClendon *et al.*, 1981)。单铃重和纤维品质均与花铃期温度密切相关,随着有效积温的减少,铃重和纤维品质呈明显降低趋势(熊宗伟等, 2014)。夜间温度对纤维伸长的影响大于日均温,温度升高使纤维长度变短,也有研究认为温度身高对纤维长度无影响(Reddy *et al.*, 1999; 单世华等, 2000; Pettigrew, 2008)。铃期 $\geq 15^{\circ}\text{C}$ 有效积温、日均温和最低温与纤维比强度和马克隆值呈明显正相关关系,平均温度和降水量有利于断裂比强度增加(马富裕等, 2005; Conaty *et al.*, 2015)。纤维长度与品种遗传特性和环境因素的关系备受关注(Reddy *et al.*, 1993; Raper *et al.*, 2019),特殊气候环境会使相同棉花品种在本地区种植后长度减短 2.0~3.0 mm(Qin *et al.*, 2011)。低温不仅限制了纤维伸长,也阻碍了次生壁发育和纤维素沉积。前人研究表明,当日均温小于 21.1°C 时,断裂比强度的降幅较大,但对纤维长度则无显著影响;小于 15°C 的日均温则纤维发育受阻(赵新华等, 2010; 卞云海等, 2009; Conaty *et al.*, 2015)。棉铃铃龄 13~19 d 遭遇高温胁迫会使棉花光合能力急剧下降,干物质积累能力和光合输出能力显著减少,导致纤维长度变短、马克隆值偏离最适范围(贺新颖等, 2013; 郭林涛等, 2015; Conaty *et al.*, 2015; Xu *et al.*, 2017)。

光照是影响棉花纤维品质的另一个关键因子,光照不足影响了棉叶光合产物的形成,减少了碳水化合物向棉铃的输出,降低了棉纤维品质(王庆材等, 2005)。光照也通过影响纤维素合成而影响纤维品质(Kasperbauer, 2000)。多数研究认为弱光会增加纤维长度,这可能与蔗糖转化和 β -1,3-葡聚糖的沉积和降解有关(Lv *et al.*, 2013; Pettigrew *et al.*, 2001)。也有研究发现弱光使纤维长度变短或对纤维长度和伸长率影响较小(王庆材等, 2005; Zhao *et al.*, 2000)。花铃期遮阴条件下,纤维素最大累积速率降低,抑制了纤维素的合成与累积,使得纤维比强度、成熟度和马克隆值在其形成过程中的增加速率以及

纤维细度在加厚期内的降低速率降低,因此弱光导致纤维比强度和马克隆值降低,对纤维比强度和长度整齐度的影响与弱光时期及持续时间有关(Lv *et al.*, 2013; Zhao *et al.*, 2000; 张新新等, 2015)。纤维长度、比强度、马克隆值均与花后日照时数呈显著或极显著正相关关系(周治国等, 1999)。棉纤维伸长发育期的累积辐热积(PTP)可综合温光复合因子的效应,PTP与棉纤维伸长特征值、比强度快速和稳定增加期特征值以及最终长度、比强度均存在极显著相关(赵文青等, 2011)。

1.2.2 机器学习算法对作物产量与品质的预测

作物产量与气候及多环境因素间关系通常是非线性的(Jeong *et al.*, 2016; Li *et al.*, 2019),基于机器学习的作物产量预测模型也有很多,如人工神经网络(Chlingaryan *et al.*, 2018)、最小绝对收缩和选择算子回归(Cao *et al.*, 2020)、支持向量机(Liakos *et al.*, 2018)和随机森林(Norouzi *et al.*, 2010; Peng *et al.*, 2020., 钟仁海, 2022)等。采用随机森林模型解析了冬小麦实际单产、气象产量和相对气象产量的关系,构建了多种变量组合模型并对产量进行回归预测,结合袋外数据重要性结果提出了突破冬小麦产量限制因子的关键技术途径(刘峻明等, 2019; Gaso *et al.*, 2019)。利用支持向量机、随机森林和反向神经网络构建大豆产量估算模型发现,大豆关键生育时期光谱指数与产量相关性较好(唐子竣等, 2023),全莢期光谱指数与大豆产量的相关系数高达0.72,引入一阶微分光谱指数构建出产量最优估算模型(唐子竣等, 2023)。灰色系统理论GM(1, 1)模型与时间序列算法的组合算法模型能够较好且稳定的预测吉林省玉米产量变化趋势(季宇等, 2018),决策树、随机森林、回归支持向量机、人工神经网络、堆叠稀疏自动编码器、卷积神经网络和长短期记忆等多种机器学习方法可以预测水稻、玉米和大豆的产量,且支持向量机预测结果的更为准确(Sungha *et al.*, 2021)。Dilli等提出了一种基于机器学习的多空间区域作物产量预测方法,较为精准地预测了9个国家的6种作物产量趋势(Dilli *et al.*, 2021)。

区域级别的机器学习模型具有更低的归一化均方根误差(NRMSE)和不确定性。回归模型和机器学习模型都很好再现观察到的产量模式,但基于过程作物模型的使用协调参数存在很大偏差(Feng *et al.*, 2019);在收益率概率分布方面,机器学习模型的预测表现最好,其次是回归模型和基于过程模型(Leng *et al.*, 2020)。从宏观尺度来看,机器学习模型预测作物产量的性能优于回归模型和基于过程的模型,其对产量变异性的解释度为93%,回归模型和基于过程模型的仅有42%~51%(Leng *et al.*, 2020);与LASSO回归方法相比,基于机器学习算法构建的作物产量预测经验模型更优(Cai *et al.*, 2019)。从区域尺度来看,卷积神经网络、高斯过程回归和长短期记忆网络能更好的预测县级冬小麦产量,得到了相对较好的越策结果(Wang *et al.*, 2020; Han *et al.*, 2020)。利用随

机森林、长短期记忆网络和深度神经网络结合多源环境变量也预测了县级小麦产量趋势, 预测结果的精度也较高 (Cao *et al.*, 2021)。基于决策树的随机森林模型已被广泛应用 (Rehfeldt *et al.*, 2012; Singh *et al.*, 2017; Zhao *et al.*, 2019), 并在产量预测方面有良好的表现 (Han *et al.*, 2020; Maya Gopal and Bhargavi., 2019)。

1.2.3 预测模型构建的特征选择和特征提取

高维数据解析具有较高的挑战性, 选取有效的特征建立模型就会运用特征选择和特征提取。特征提取是从原始特征集中提取出新特征, 所提取的新特征可能与原来特征含义不同, 目的是通过新的少量特征来表示原始特征所含的大量信息; 特征选择是从原始特征序列中依据特定的方式选择出对建模有效的特征, 以提高建模效率和准确率, 所选特征与原始特征含义一致。特征选择算法一般包括过滤法、嵌入法和包装法 (Venkatesh *et al.*, 2019; Augusto *et al.*, 2018; Zhou *et al.*, 2019), 基于过滤特征选择的算法往往都使用了互信息特征 (Maldonado *et al.*, 2014), 这不仅能解决分类问题, 也能用于回归问题的特征处理 (王李娟等., 2020)。互信息特征选择算法和 Pearson 系数改进了冗余特征的筛选过程, 创建了较为精准的交互特征筛选系统 (包芳, 2021)。此外, 采用并行遗传算法进行特征选择提高了预测的精度, 与多元线性回归和支持向量机相比, 基于并行遗传算法的机器学习算法构建的预测模型更科学准确, 土壤团聚体几何平均直径实测值与预测值的决定系数提高了 12.2%~20.0% (Besalatpour *et al.*, 2014)。最大特征数对随机森林的影响较大, 运用随机森林回归模型, 对所选的 10 个影响土壤侵蚀的因素进行重要性排序得出影响土壤侵蚀发生显著的因素有 5 个, 并对不同年份土壤侵蚀预测值和实际值进行比较, 模型的 R^2 为 0.89, 模型拟合效果明显提高 (赵琦等, 2024; Genuer *et al.*, 2010); 特征优选的随机森林算法提取土地利用信息的效果最佳, 总体精度高达 88.2% (王李娟等, 2020)。通过相关性分析和特征选择算法结合优选变量, 能够较单独通过相关分析明显提升 LSTM 和 PSO-BPNN 模型的建模精度, 但对 RF 模型则无法优化变量 (马宇欣等, 2024)。建立不同输入特征变量下的覆膜冬小麦 LAI 反演模型会发现采用适宜的特征降维方法结合机器学习算法能够提高覆膜冬小麦 LAI 的反演精度和稳定性 (谷晓博等, 2023)。在选取相同特征指标参数的情况下, 建立小麦冠层叶绿素含量估测模型, SVR 的预测能力优于多元线性回归和岭回归 (苑迎春等, 2022)。

1.2.4 棉纤维品质预测模型的构建

新疆棉花种植区域分布广泛, 影响纤维品质的气象因子不尽相同。空间分析和时间序列分析方法解析了气候时空分布规律和演变特征 (李谢辉等, 2015), 通过灰色关联度综合评价了棉花区试品种的纤维品质, 揭示了纤维品质指标演变规律和发展趋势 (陈荣

江等, 2007)。通过建立棉纤维主要品质性状的气象生态模型, 指出棉花铃期日均最低温、日均最高温和相对湿度、夜均温和日均降水量分别是影响棉纤维长度、比强度、马克隆值的关键气象因子(赵文青等, 2008)。目前, 高光谱指数构建的作物产量估算预测模型会因种植区域、作物种类和建模方法的不同, 其预测精度也有较大差异(李严明, 2019; 费帅鹏等; 2021), 提高作物产量和品质模型的估算精度仍是亟需解决的关键问题。机器学习预测模型虽已广泛应用于产量预测, 但在作物品质预测方面的研究较少, 且主要集中于品质时空分布规律的研究。目前用于作物品质预测的方法有时间序列法、传统回归预测模型、多元线性回归模型和神经网络模型等, 但这些方法存在较大的局限性。GM(1, 1)预测模型是基于—阶常微分方程, 通过小样本数据建立数学预测模型, 对于短期预测有较高的精度, 更适用于小样本数据和贫信息序列的预测(梁毅等, 2012)。多元线性回归模型主要用于中长期预测, 也能直观的显示解释变量对被解释变量的影响, 但会因多重共线性问题造成大量信息损失和较大误差(陶惠林等, 2020)。神经网络模型则由于数据收集条件限制和模型选取数据量较小, 其输出结果精度较低(张家瑜等, 2024)。可见, 机器学习算法在非线形拟合上具有较好拟合能力, 传统机器学习的随机森林算法可以从大量数据和参数中筛选出影响纤维品质的关键参数, 支持向量机算法在处理高维数据时有很强的学习能力, 神经网络模型的长短期记忆模型对于预测时序性数据方面具有更好的效果。

1.3 研究内容

本研究利用新疆历年原棉纤维品质公证检验数据, 开展基于机器学习算法的新疆棉花纤维品质预测模型的研究, 通过大数据分析气象因素对纤维品质的相对贡献及其关系; 采用机器学习算法构建气象因子影响原棉品质的预测模型, 寻求适宜纤维品质各指标的最优预测模型, 以期为新疆棉花区域性育种和栽培提供依据。技术路线如图 1-1 所示。

(1) 构建棉花纤维品质预测模型。

统计 2015~2022 年新疆各区域的原棉纤维品质数据, 分析气象因子对棉花纤维品质的影响及关键气象指标的相对重要性, 采用随机森林(Random Forest)、支持向量回归算法(Support Vector Regression)和长短期记忆网络(Long Short-Term Memory)等机器学习语言构建纤维品质预测模型, 明确纤维品质预测模型建立的关键步骤及参数的优化调整。

(2) 棉花纤维品质指标预测模型检验。

采用稳定性分析和误差分析等方法利用验证数据来检验模型的合理性和适用性, 使得模型符合预测要求; 比较检验模型预测实际效果, 筛选出纤维品质各指标的最优机器

学习算法预测模型，科学预测未来新疆各产棉区纤维品质变化趋势和区域分布。

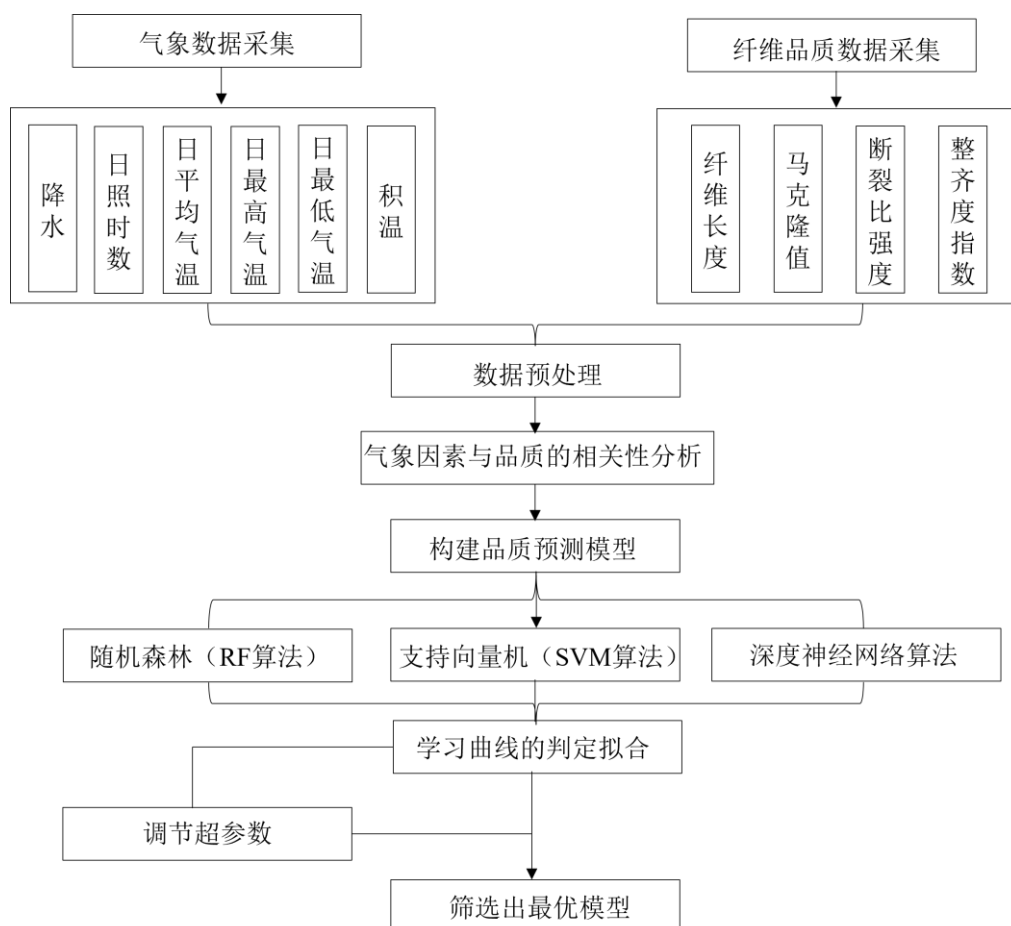


图 1-1 技术路线

Fig. 1-1 Structure of this research paper