

T/SSSC

中国土壤学会团体标准

T/SSSC 0—2026

# 红壤长期植被覆盖区土壤有机质退化遥感 监测技术规程

Technical Code of Practice for Remote Sensing Monitoring of Soil Organic Matter  
Degradation in Long-term Vegetation-covered Red Soil Regions

(征求意见稿)

2026-0-0 发布

2026-0-0 实施

中国土壤学会 发布

# 目 次

|                        |    |
|------------------------|----|
| 前 言 .....              | II |
| 1 范围 .....             | 1  |
| 2 规范性引用文件 .....        | 1  |
| 3 术语和定义 .....          | 1  |
| 4 监测流程 .....           | 1  |
| 5 数据获取 .....           | 2  |
| 6 遥感指标提取与筛选 .....      | 3  |
| 7 有机质遥感监测模型构建及验证 ..... | 3  |
| 8 有机质退化评价 .....        | 4  |
| 9 监测报告 .....           | 4  |
| 附录 A .....             | 1  |
| 附录 B .....             | 4  |
| 附录 C .....             | 5  |
| 附录 D .....             | 6  |
| 参考文献 .....             | 7  |

## 前 言

本文件按照GB/T 1.1—2020《标准化工作导则 第1部分：标准化文件的结构和起草规则》的规定起草。

请注意本文件的某些内容可能涉及专利。本文件的发布机构不承担识别专利的责任。

本文件由中国科学院南京土壤研究所提出。

本文件由中国土壤学会归口。

本文件起草单位：中国科学院南京土壤研究所、南京工业大学、南京大学

本文件主要起草人：马利霞、于东升、解宪丽、刘杰、刘明、张乾、关廷宇、郑光

# 红壤长期植被覆盖区土壤有机质退化遥感监测技术规程

## 1 范围

本文件确立了针对红壤长期植被覆盖区土壤有机质退化遥感监测的技术规程，包括术语和定义、监测流程、数据获取、遥感指标提取与筛选、有机质遥感监测模型构建及验证、有机质退化评价、监测报告。

本文件适用于红壤长期植被覆盖区土壤有机质退化监测。

## 2 规范性引用文件

下列文件中的内容通过文中的规范性引用而构成本文件必不可少的条款。其中，注日期的引用文件，仅该日期对应的版本适用于本文件；不注日期的引用文件，其最新版本（包括所有的修改单）适用于本文件。

GB/T 14950 摄影测量与遥感术语

DB37/T 3242 农作物种植面积遥感监测技术规程 马铃薯

NY/T 3527 农作物种植面积遥感监测规范

NY/T 4151 农业遥感监测无人机影像预处理技术规范

NY/T 1121.1 土壤检测第1部分：土壤样品的采集、处理和贮存

NY/T 1121.6 土壤检测第6部分：土壤有机质的测定

HJ 1231 土壤环境 词汇

HJ 1068 土壤粒度的测定吸液管法和比重计法

## 3 术语和定义

下列术语和定义适用于本文件。

### 3.1 黏粒含量 Clay content

土壤中粒径 $<2\ \mu\text{m}$ 的矿物质颗粒所占的比例。

### 3.2 砂粒含量 Sand content

土壤中粒径小于2mm大于等于0.05mm的矿物质颗粒所占的比例。

### 3.3 粉粒含量 Silt content

土壤中粒径小于0.05mm大于等于0.002mm的矿物质颗粒所占的比例。

### 3.4 训练样本 Training sample

可由实地调查或图像解释方法选取确定的已知地物属性或特征的图像像元，用于进行分类的学习和训练，以建立分类模型或分类函数的样本。

### 3.5 验证样本 Validation sample

可由实地调查或图像解释方法选取确定的已知地物属性或特征的图像像元，用于验证分类结果精度的样本数。

### 3.6 长期植被覆盖区 Long-term vegetation-covered regions

由于常绿林占比高以及集约化种植导致植被多年（如 $\geq 5$ 年）覆盖，加之云雾频发导致裸土时期遥感缺失的区域。

## 4 监测流程

红壤长期植被覆盖区土壤有机质退化遥感监测主要包括数据获取、遥感指标提取与筛选、有机质遥感监测模型构建及验证、有机质退化评价4个步骤（图1）。

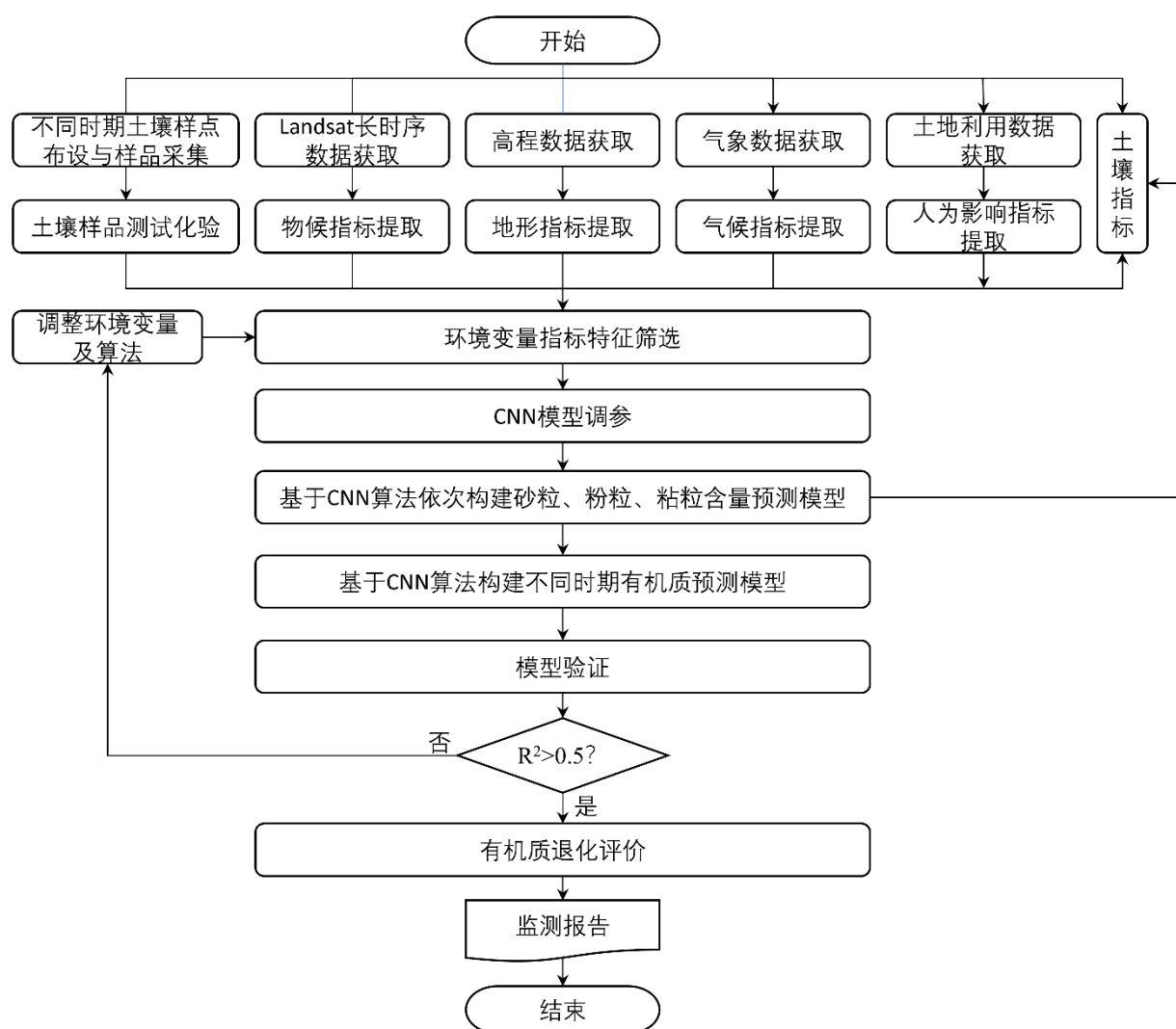


图1 流程图

## 5 数据获取

### 5.1 土壤样品采集与化验

#### 5.1.1 土壤样点布设

地面样点布设主要依据土地利用（森林、农田、草地等）和土壤类型占比，同时考虑样点在空间分布的均匀性。

#### 5.1.2 土壤样品采集

不同时期通过多点混合采样方法采集0 - 20 cm深度的土壤样品，质量约500 g。样品密封于塑料自封袋并标注点位信息，同时使用GPS记录采样坐标。

#### 5.1.3 土壤样品测试化验

按照NY/T 1121.1—2006规定的方法进行土壤样品的处理和贮存，按照HJ 1068—2019规定的方法化验土壤砂粒、粉粒、黏粒含量，按照NY/T1121.6—2006规定的方法化验土壤有机质含量。

### 5.2 遥感影像数据获取

#### 5.2.1 Landsat长时序数据

基于经过辐射定标、大气校正、几何校正等预处理后的Landsat数据集[1]，选取目标年份及前三年的卫星影像，并按照监测区范围、行政区划图进行裁剪和掩膜处理。

#### 5.2.2 高程数据

选取12.5米的ALOS高程数据用于提取相关地形因子。

### 5.2.3 气象数据

通过中国国家地球系统科学数据中心下载目标年份逐月温度和降雨数据集。

### 5.2.4 土地利用数据

为提取不同土地利用空间分布数据。可通过中国国家地球系统科学数据中心下载基于Landsat数据通过目视解译得到的土地利用空间分布。也可自行通过目标年份的Landsat数据集进行相关目视解译。提取不同土地利用（林地、草地、水田、旱地）的空间分布数据。

## 6 遥感指标提取与筛选

### 6.1 特征提取

基于遥感数据提取相关环境变量，包括气候、地形、人为影响、物候及土壤指标。

#### 6.1.1 气候指标

针对目标年份，提取温度和降雨的月均数据作为气候指标。

#### 6.1.2 地形指标

基于ALSO高程数据提取14个基础地形指标，包括高程值、坡度、坡向、地形湿度指数、坡长坡度因子、分析性山体阴影、河道网络基准面、河道网络距离、闭合洼地、汇流指数、平面曲率、剖面曲率、相对坡度位置、山谷深度。具体参见附录A。

#### 6.1.3 人为影响指标

基于样点土壤属性变量大小对土地利用进行编号（例如旱地1、草地2、林地3、水田4），作为人为影响指标。

#### 6.1.4 物候指标

基于Landsat时序反射率数据提取两类物候指标。即年内物候指标和年际物候指标。具体参见附录A。

#### 6.1.5 土壤特征

预测有机质的土壤特征包括土壤颗粒组成（模型预测的砂、粉、黏粒量）。为保证颗粒组成之间的一致性，砂、粉、黏粒量采用逐步建模策略：首先，以土壤类型为土壤特征预测粉粒含量；其次，以土壤类型和预测粉粒为土壤特征预测砂粒含量；随后，计算初始黏粒含量（即1-粉粒-砂粒），并结合土壤类型、预测的粉粒和砂粒含量为主要土壤特征进一步预测黏粒含量。最终，将土壤类型、预测的粉粒、砂粒、以及两种黏粒含量作为土壤特征共同用于有机质预测模型。

### 6.2 特征筛选

为了减少冗余以及多重共线性，通过两步来对特征进行选择。首先，通过对所有变量进行皮尔逊相关分析，对于两个变量相关性较高（相关性超过0.98），选择删除其中一个。对于剩余的所有变量通过随机森林算法进行重要性排序，具体参见附录B。进而计算重要性累计曲线，设置阈值选择用于土壤属性建模的特征。其中阈值作为一个参数通过与模型参数共同调参确定。

## 7 有机质遥感监测模型构建及验证

### 7.1 模型调参

通过一维卷积神经网络（CNN）算法构建基于遥感指标预测目标土壤属性的模型。CNN模型架构包括两个卷积层和丢弃层，后接一个全局平均池化层和全连接输出层。每个卷积层都需设置滤波器、卷积核，丢弃层需设定丢弃率参数。所有层均采用ReLU激活函数。模型通过Adam优化器进行编译，需设置学习率参数。此外，为平衡训练速度和稳定性，需设定批量处理大小参数。模型所有参数通过贝叶斯优化进行调优，具体参数和调优区间见附录表。参数的性能可通过交叉验证进行评估，重点关注较低的训练集平均绝对误差（MAE）。

### 7.2 模型构建

先对土壤砂粒、粉粒、黏粒含量进行模型构建。将筛选特征作为输入变量，采用CNN算法依次建立土壤砂粒、粉粒、黏粒含量的模型。将其模型预测值作为特征变量加入到环境变量中，再次采用CNN算法分别建立T1和T2不同时期土壤有机质含量预测模型。CNN算法具体参见附录C。

### 7.2 模型验证

通过基于 X - Y 距离的样本集划分方法 (SPXY) 抽取70%的样本用于模型训练, 30%的样本用于模型验证, 具体参见附录D。基于验证样本对各属性建模结果进行精度评价, 选择决定系数 ( $R^2$ )、均方根误差 (RMSE)、相对分析误差 (RPD) 作为精度评价指标。按照公式 (1) 计算  $R^2$ , 按照公式 (2) 计算 RMSE, 按照公式 (3) 计算 RPD。原则上土壤有机质的决定系数  $R^2$  应大于 0.5, 精度验证不合格的, 需要调整输入变量或算法, 直至满足精度要求。

$$R^2 = 1 - \frac{\sum_{i=1}^n (\hat{y}_i - y_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \quad (1)$$

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2} \quad (2)$$

$$RPD = \frac{SD}{RMSE} \quad (3)$$

式中,

$R^2$ —决定系数;

RMSE—均方根误差;

SD—标准差;

RPD—相对分析误差;

n—采样点数量;

$\hat{y}_i$ —土壤样点 i 的预测值;

$y_i$ —土壤样点 i 的实测值;

$\bar{y}$ —土壤样点验证样本实测值的平均值。

## 8 有机质退化评价

根据不同时期样点构建有机质模型, 进而分别生成各个时期有机质空间分布图, 通过后一时期结果减去前一时期结果得到有机质在目标期间的变化结果。考虑两个时期制图误差 RMSE 设定分级间隔, 间隔等于两个时期制图较大的 RMSE 值, 将有机质变化分为六个级别, 包括“显著降低”、“降低”、“轻微降低”、“轻微升高”、“增加”、“显著增加”。

## 9 监测报告

有机质退化监测报告包括采用的遥感数据和样点数据、监测技术流程、有机质制图结果和精度评价、有机质退化评价等信息。监测结果形式宜采用统计表格和图表等, 统计表格包括样点特征等信息。图片包括土壤样点分布和有机质退化分级分布图片等。

## 附录 A

(资料性)  
遥感指标计算

表A.1 地形指标

| 指标         | 含义   |
|------------|--|
| 高程值        | 高程数据的值 $Z$ ，直接由DEM提取   |
| 坡度         | $Slope = \arctan\left(\sqrt{\left(\frac{\partial Z}{\partial x}\right)^2 + \left(\frac{\partial Z}{\partial y}\right)^2}\right)$   |
| 坡向         | $Aspect = \arctan 2\left(\frac{\partial Z}{\partial y}, \frac{\partial Z}{\partial x}\right)$  |
| 地形湿度指数     | $TWI = \ln\left(\frac{SCA}{\tan(\beta)}\right)$ SCA为汇水面积， $\beta$ 为局部坡度  |
| 坡长坡度因子     | $LS = \left(\frac{SCA}{22.13}\right)^{0.4} \times \left(\frac{\sin \beta}{0.0896}\right)^{1.3}$  |
| 分析性山体阴影    | $Hillshade = 255 \times (\cos(slope) \times \cos(90^\circ - altitude) + \sin(slope) \times \sin(90^\circ - altitude) \times \cos(azimuth - aspect))$   |
| 河道网络基准面[2] | 基于河道网络栅格，采用张力样条插值对河道单元格的高程值进行插值，生成整个区域的基准面高程。  |
| 河道网络距离     | $D_{channel} = Z - Z_{baselevel}$ $Z_{baselevel}$ 表示插值的河道基准面高程。  |
| 闭合洼地[3]    | 通过洼地填充算法识别，原始DEM与填充后DEM的差值即为洼地深度。  |
| 汇流指数[4]    | 采用D8算法计算每个栅格单元的汇流累积量，表示上游汇流区域内的栅格数量。   |
| 平面曲率       | $Plan\_curvature = \frac{\partial^2 Z}{\partial x^2} \left(\frac{\partial Z}{\partial y}\right)^2 - 2 \frac{\partial^2 Z}{\partial x \partial y} \frac{\partial Z}{\partial x} \frac{\partial Z}{\partial y} + \frac{\partial^2 Z}{\partial y^2} \left(\frac{\partial Z}{\partial x}\right)^2$ 正值表示汇聚流，负值表示发散流。  |
| 剖面曲率       | $Plan\_curvature = \frac{\frac{\partial^2 Z}{\partial x^2} \left(\frac{\partial Z}{\partial x}\right)^2 + 2 \frac{\partial^2 Z}{\partial x \partial y} \frac{\partial Z}{\partial x} \frac{\partial Z}{\partial y} + \frac{\partial^2 Z}{\partial y^2} \left(\frac{\partial Z}{\partial y}\right)^2}{\left(\frac{\partial Z}{\partial x}\right)^2 + \left(\frac{\partial Z}{\partial y}\right)^2}$ 沿坡度方向的地表曲率，正值表示凸坡，负值表示凹坡。 |
| 相对坡度位置     | $RSP = \frac{HO - HU}{2}$ HO为坡高，HU为谷深  |
| 山谷深度       | $Valley\_Depth = Z_{ridge} - Z$  |

表A.2 年内物候指标

| 类型               | 指标   | 含义   |
|------------------|--|--|
| 统计值              | 蓝光波段反射率 Blue   | <b>统计值:</b><br>最小值<br>最大值<br>第二小值<br>第二大值<br>平均值<br>最小值与第一四分位数之间的平均值<br>第三四分位数与最大值之间的平均值<br>第一四分位数与第三四分位数之间的平均值<br>除去最大值和最小值后的平均值<br>中值<br>标准差<br>各数据区间的差异值<br>最后一次观测值   |
|                  | 绿光波段反射率 Green  |  |
|                  | 红光波段反射率 Red  |  |
|                  | 近红外波段反射率 Nir   |  |
|                  | 短波红外波段反射率 1 Swir1  |  |
|                  | 短波红外波段反射率 2 Swir2  |  |
|                  | 归一化植被指数 $NDVI=(Nir-Red)/(Nir+Red)$   |  |
|                  | 归一化指数 $S1N=(Nir-Swir1)/(Nir+Swir1)$  |  |
|                  | 归一化指数 $S2N=(Nir-Swir2)/(Nir+Swir2)$  |  |
|                  | 归一化指 $GN=(Green-Red)/(Green+Red)$  |  |
|                  | 归一化指 $S1S2=(Swir1-Swir2)/(Swir1+Swir2)$  |  |
|                  | 光谱变异指数 $SVVI=SD(Blue, Green, Red, Nir, Swir1, Swir2) - SD(Nir, Swir1, Swir2)$  |  |
|                  | 缨帽变换绿度<br>$TCG=-0.1603 \times Blue + -0.4934 \times Red + 0.2819 \times Green + 0.7940 \times Nir + -0.0002 \times Swir1 + -0.1446 \times Swir2$ |  |
| 特定指标统计值对应日期的光谱数据 | 蓝光波段反射率 Blue   | <b>由以下指标决定的日期:</b><br>最小 NDVI<br>最大 NDVI<br>第二小 NDVI<br>第二大 NDVI<br>NDVI 中位数<br>最小值与第一四分位数之间的 NDVI 均值<br>第三四分位数与最大值之间的 NDVI 均值<br>最小 LST<br>最大 LST<br>第二小 LST<br>第二大 LST<br>LST 中值<br>最小值与第一四分位数之间的 LST 均值<br>第三四分位数与最大值之间的 LST 均值<br>最小 S2N<br>最大 S2N<br>最小值与第一四分位数之间的 S2N 均值<br>第三四分位数与最大值之间的 S2N 均值 |
|                  | 绿光波段反射率 Green  |  |
|                  | 红光波段反射率 Red  |  |
|                  | 近红外波段反射率 Nir   |  |
|                  | 短波红外波段反射率 1 Swir1  |  |
|                  | 短波红外波段反射率 2 Swir2  |  |
| NDVI 相关的季节特征数据   | RNph_sos   | 生长季开始的 NDVI  |
|                  | RNph_eos   | 生长季结束的 NDVI  |
|                  | RNph_sos_slope   | 从生长季开始到最大值的 NDVI 增长斜率  |
|                  | RNph_eos_slope   | 从最大值到生长季结束的 NDVI 下降斜率  |
|                  | RNph_sos_amp   | 从生长季开始到最大值的 NDVI 幅度  |
|                  | RNph_eos_amp   | 从最大值到生长季结束的 NDVI 幅度  |
|                  | RNph_ave   | 生长季开始与结束之间的 NDVI 平均值   |
| RNph_sum         | 生长季开始与结束之间的 NDVI 总和  |  |

表A.3 年际物候指标

| 类型                       | 指标                                  | 含义  |
|--------------------------|-------------------------------------|---|
| 目标年份前三年不同波段和指标指数的统计值     | 蓝光波段反射率 Blue                        | 最小值   |
|                          | 绿光波段反射率 Green                       | 最大值   |
|                          | 红光波段反射率 Red                         | 第二小值  |
|                          | 近红外波段反射率 Nir                        | 第二大值  |
|                          | 短波红外波段反射率 1 Swir1                   | 中位数   |
|                          | 短波红外波段反射率 2 Swir2                   | 标准差   |
|                          | 归一化植被指数 $NDVI=(Nir-Red)/(Nir+Red)$  | 全部数值的均值<br>除去最小值和最大值后的均值                          |
|                          | 归一化指数 $S1N=(Nir-Swir1)/(Nir+Swir1)$ | 最后一次观测值   |
| 目标年以及前三年对应指标的线性拟合坡度和标准差  | 蓝光波段反射率 Blue                        | 数值与观测日期线性回归的斜率<br>标准差                             |
|                          | 绿光波段反射率 Green                       |   |
|                          | 红光波段反射率 Red                         |   |
|                          | 近红外波段反射率 Nir                        |   |
|                          | 短波红外波段反射率 1 Swir1                   |   |
|                          | 短波红外波段反射率 2 Swir2                   |   |
|                          | 归一化植被指数 $NDVI=(Nir-Red)/(Nir+Red)$  |   |
|                          | 归一化指数 $S1N=(Nir-Swir1)/(Nir+Swir1)$ |   |
| 目标年和前三年的差异差值统计值          | 蓝光波段反射率 Blue                        | 最小值   |
|                          | 绿光波段反射率 Green                       | 最大值   |
|                          | 红光波段反射率 Red                         | 第二小值  |
|                          | 近红外波段反射率 Nir                        | 第二大值  |
|                          | 短波红外波段反射率 1 Swir1                   | 最大值之后的数值  |
|                          | 短波红外波段反射率 2 Swir2                   | 最小值之后的数值  |
|                          | 归一化植被指数 $NDVI=(Nir-Red)/(Nir+Red)$  | 全部数值的均值<br>去除最小值和最大值后的均值                          |
|                          | 归一化指数 $S1N=(Nir-Swir1)/(Nir+Swir1)$ |   |
| 特定指标（例如NDVI）统计值对应日期的光谱指标 | 蓝光波段反射率 Blue                        | 最小 NDVI   |
|                          | 绿光波段反射率 Green                       | 最大 NDVI   |
|                          | 红光波段反射率 Red                         | 第二小 NDVI  |
|                          | 近红外波段反射率 Nir                        | 第二大 NDVI  |
|                          | 短波红外波段反射率 1 Swir1                   | NDVI 中位数  |
|                          | 短波红外波段反射率 2 Swir2                   | 最小 LST<br>最大 LST<br>第二小 LST<br>第二大 LST<br>LST 中位数 |

## 附录 B

(资料性)

## 基于随机森林进行特征重要性排序算法

随机森林是一种集成学习方法，由大量决策树组成。每棵决策树在训练过程中，随机抽取样本和特征进行节点分裂，以减少预测误差（如均方误差MSE）。集成多棵树后，随机森林能够有效提高预测精度并增强模型的鲁棒性。通过分析每个特征在树中分裂节点所带来的目标变量误差减小量来评估其重要性。具体步骤为：

(1) 节点分裂误差减小量

对第  $t$  棵树的第  $j$  个节点，如果特征  $X_i$  被用于分裂，则该节点的误差减小量为：

$$\Delta MSE_{i,j}^{(t)} = MSE_{parent} - \left( \frac{N_{left}}{N_{parent}} MSE_{left} + \frac{N_{right}}{N_{parent}} MSE_{right} \right) \quad (B.1)$$

$MSE_{parent}$  为分裂前节点的均方误差， $MSE_{left}$  和  $MSE_{right}$  为分裂后左右子节点的均方误差， $N_{parent}$ ， $N_{left}$  和  $N_{right}$  分别为父节点和左右子节点的样本数。

(2) 单棵树特征总贡献

对第  $t$  棵树中，特征  $X_i$  在所有节点的误差减小量进行累加：

$$imp_i^{(t)} = \sum_{j \in nodes_{X_i}} \Delta MSE_{i,j}^{(t)} \quad (B.2)$$

(3) 随机森林的总体重要性

对随机森林所有  $T$  棵树，将每棵树中特征贡献取平均，得到特征  $X_i$  的总体重要性：

$$imp(X_i) = \frac{1}{T} \sum_{t=1}^T imp_i^{(t)} \quad (B.3)$$

(4) 特征排序

对所有特征重要性归一化到0-1的范围，进而对特征从高到低进行排序。

## 附录 C

(资料性)

## 卷积神经网络算法

卷积神经网络算法 (CNN) 是一种常用于处理高维数据的深度学习模型。与传统神经网络相比, CNN 能够通过卷积操作自动提取局部特征, 并通过多层网络进行特征组合与抽象。其核心原理和操作流程如下:

## (1) 输入层

输入层包括每个样本的所有特征, 为提高训练稳定性和收敛速度, 需对输入数据进行标准化或归一化处理。

算法中, 需要指定批量大小, 即在一次参数更新过程中, 网络所使用的训练样本数量。样本梳理可以控制内存占用, 平衡训练稳定性与收敛速度。

## (2) 卷积层

基于数据自动提取局部特征。使用一组可训练的卷积核 (滤波器) 在输入数据上滑动, 通过卷积运算生成特征图。需指定滤波器和核的大小。这个过程中, 可引入激活函数, 使得网络可拟合复杂的函数关系。

## (3) 池化层

可降低特征图的空间维度, 减少参数数量和计算量, 同时保留主要特征。

## (4) 丢弃层

在训练过程中, 按照设定概率随机“丢弃”部分神经元, 以减少过拟合, 提高模型泛化能力, 防止某些神经元对预测过度依赖。

## (5) 多层卷积、池化和丢弃层组合

卷积层和池化层可重复堆叠, 逐层提取更抽象、更高层次的特征。

## (6) 全连接层

将最终卷积/池化得到的特征图展开为一维向量, 输入全连接层。通过矩阵乘法和激活函数, 将抽取到的特征映射到预测目标空间。

## (7) 输出层

输出模型预测值。

## (8) 损失函数与训练

定义损失函数衡量预测结果与真实值的差距, 回归模型常通过均方误差 (MSE) 表征。使用优化算法 (如Adam) 更新网络参数。其中, 学习率是其中一个重要参数。

## (9) 模型评估与应用

使用验证集或测试集评估CNN模型的性能, 常用指标包括 $R^2$ 、RMSE等。

针对一个两层网络的CNN, 其重要参数及调参区间见表C. 1。

表C. 1 卷积神经网络参数及调参区间

| 类型     | 参数    | 调参区间                 |
|--------|-------|----------------------|
| 第一层卷积层 | 滤波器个数 | [64, 128, 256]       |
|        | 核尺寸   | [3, 7]               |
| 第一层丢弃层 | 丢弃率   | [0, 0.5]             |
| 第二层卷积层 | 滤波器个数 | [64, 128, 256, 512]  |
|        | 核尺寸   | [3, 7]               |
| 第二层丢弃层 | 丢弃率   | [0, 0.5]             |
| 损失函数   | 学习率   | [0.01, 0.0001]       |
| 模型表现   | 批量大小  | [8, 16, 32, 64, 128] |

## 附录 D

(资料性)

## 基于X - Y距离的样本集划分方法 (SPXY)

基于X - Y距离的样本集划分方法 (SPXY) 旨在根据特征空间 (X) 和响应变量空间 (Y) 的分布合理选择训练集和测试集, 使模型训练和评估更加稳定和代表性。其核心思想是选择既能覆盖样本特征空间, 又能反映目标变量变化范围的样本用于训练。生成的训练集和测试集在空间分布上更均匀, 有利于建模稳定性和预测性能。其计算步骤为:

(1) 计算样本间X-Y距离

对每个样本, 构建包含输入特征向量 $X_i$ 和目标变量 $Y_i$ 的联合向量 $[X_i, Y_i]$ 。

计算样本之间的欧氏距离或加权距离:

$$D_{ij} = \sqrt{\|X_i - X_j\|^2 + \alpha \|Y_i - Y_j\|^2} \quad (\text{D. 1})$$

其中 $\alpha$ 为权重系数, 用于平衡特征空间和响应变量空间的重要性, 默认可取1。

(2) 选择初始样本

从样本集中选择一堆距离最远的样本作为初始训练集样本, 该步骤保证训练集覆盖特征和响应变量的边界。

(3) 逐步添加样本到训练集

对剩余样本, 计算其到已选训练样本的最小距离

$$d_i = \min_{j \in \text{训练集}} D_{ij} \quad (\text{D. 2})$$

选择距离最远的样本加入训练集, 重复该过程, 指导训练集达到预设比例或与其样本数量。

(4) 划分测试集

剩余样本作为测试集, 保证测试集覆盖特征空间和响应变量空间的未选区域, 增强模型评估的代表性。

参考文献

- [1] Potapov, P., Hansen, M.C., Kommareddy, I., Kommareddy, A., Turubanova, S., Pickens, A., Adusei, B., Tyukavina, A., & Ying, Q. (2020). Landsat Analysis Ready Data for Global Land Cover and Land Cover Change Mapping. *Remote Sensing*, 12, 426
- [2] Conrad, O. (2002). Vertical Distance to Channel Network. SAGA-GIS Tool Library Documentation.
- [3] Conrad, O. (2007). Sink Removal (Fill Sinks). SAGA-GIS Tool Library.
- [4] Tarboton, D.G. (1997). A new method for the determination of flow directions and upslope areas in grid digital elevation models. *Water Resources Research*, 33, 309-319