

方法2完整说明

梯度提升回归模型说明 (Gradient Boosting Regression)

[基本原理](#)

[模型构建过程](#)

[损失函数](#)

[超参数及影响](#)

[优缺点](#)

[程序思路说明](#)

梯度提升回归模型说明 (Gradient Boosting Regression)

梯度提升回归 (Gradient Boosting Regression) 是机器学习领域里用于解决回归问题的一种极为强大的算法，它的核心构建在提升 (Boosting) 算法的基础之上。

基本原理

提升算法的核心思路是将多个能力相对较弱的学习器组合起来，形成一个预测能力强大的模型。这些弱学习器通常是比较简单的模型，它们自身的预测能力有限，但通过特定的迭代方式，后续的学习器能够专注于修正前面学习器所产生的错误。经过多轮这样的迭代和改进，最终将这些弱学习器组合起来，就可以获得预测精度很高的强学习器。

梯度提升回归的独特之处在于，它在每一次迭代过程中，会计算损失函数关于当前模型预测值的负梯度。简单来说，损失函数就是用来衡量模型预测值与真实值之间差距的一种工具。负梯度代表着模型预测误差最大的方向，算法让新的弱学习器朝着这个方向去进行拟合，也就是让新的弱学习器努力去减少当前模型的损失，通过这样一轮又一轮的迭代，不断提升整个模型的预测性能。

模型构建过程

- 1. 初始化模型：**一开始，我们通常会把初始模型设定为一个常数。这个常数一般可以是训练数据中目标变量的平均值。之所以这样设定，是因为在没有任何先验知识的情况下，用平均值作为初始预测是一个比较合理的起点，它为后续的迭代改进提供了一个基础。
- 2. 迭代训练弱学习器：**在每一轮的迭代中，首先要计算损失函数关于当前模型预测值的负梯度。这个负梯度实际上反映了当前模型预测结果与真实值之间的偏差情况。然后，我们利用这些负梯度值作

为新的目标值，去训练一个新的弱学习器，这里的弱学习器一般采用决策树。决策树会根据这些负梯度值来构建，目标是尽量缩小负梯度值与自身预测值之间的误差。接着，要确定这个新的弱学习器在整个模型中的权重。这个权重的确定方式，是通过寻找一个合适的值，使得加入这个弱学习器后，模型的整体损失能够最小化。最后，把这个带有合适权重的弱学习器添加到当前模型中，完成模型的更新。

3. **最终模型**：经过多轮这样的迭代之后，我们就得到了最终的梯度提升回归模型。这个模型是由最初的初始模型，以及在每一轮迭代中添加进来的一系列带有不同权重的弱学习器共同组成的。

损失函数

损失函数在梯度提升回归中起着至关重要的作用，它主要用于衡量模型预测值和真实值之间的差异。常见的损失函数有以下几种：

- **平方损失函数**：计算方式是用真实值减去预测值的差的平方。它的优点是计算起来比较简单，容易理解。它对误差的平方进行惩罚，这就意味着如果模型的预测误差较大，那么它受到的惩罚会更严重，所以模型会更加关注那些预测偏差较大的数据点，努力去减少这些较大的误差。
- **绝对损失函数**：计算方式是真实值与预测值差值的绝对值。与平方损失函数相比，绝对损失函数对于异常值更加鲁棒。因为它对误差的惩罚是线性的，不会像平方损失函数那样，对较大的误差进行过度的惩罚，所以在处理含有异常值的数据时，表现会相对稳定。
- **Huber损失函数**：它综合了平方损失函数和绝对损失函数的特点，是两者的一种折衷。当模型的预测误差比较小时，Huber损失函数的表现和平方损失函数类似；而当预测误差较大时，它又和绝对损失函数类似。这种特性使得它在不同的误差情况下，都能有比较好的表现，既能在误差较小时保证模型对数据的精确拟合，又能在误差较大时避免对异常值的过度敏感。

超参数及影响

- **n_estimators**：这个参数代表的是弱学习器（也就是决策树）的数量。当我们增加这个参数的值时，就相当于增加了模型中参与组合的弱学习器的数量。更多的弱学习器意味着模型有更强的拟合能力，能够更好地捕捉数据中复杂的模式和规律。但是，如果这个值设置得过大，模型就可能会对训练数据过度拟合。也就是说，模型虽然在训练数据上表现得非常好，能够准确地预测训练数据中的每一个样本，但在面对新的、未见过的测试数据时，就可能无法准确地进行预测，泛化能力变差。
- **learning_rate**：学习率决定了每次迭代中，弱学习器对模型更新的幅度大小。如果学习率比较小，那么每次迭代时模型的更新幅度就会比较小，这就需要更多次的迭代（也就是更大的 `n_estimators`），模型才能达到比较好的效果。不过，较小的学习率也有好处，它可以使模型的训练过程更加稳定，降低模型过拟合的风险。相反，如果学习率比较大，模型在训练过程中收敛的速度会比较快，能够更快地完成训练。但是，较大的学习率也可能导致模型在训练过程中跳过最优解，使得模

型无法收敛到一个较好的结果，甚至可能出现训练结果发散的情况，也就是模型的性能变得越来越差。

- `max_depth`: 这个参数指的是决策树的最大深度。如果决策树的深度比较大，它就能够学习到数据中非常复杂的模式。然而，这也容易导致模型过拟合，因为深度较大的决策树可能会过于关注训练数据中的一些细节和噪声，而忽略了数据的整体规律。当面对新的数据时，模型就无法准确地进行预测。相反，如果决策树的深度较小，它的拟合能力就会受到限制，可能无法捕捉到数据中一些重要的特征和关系，从而导致模型欠拟合，也就是模型的预测能力不足，无法很好地对数据进行建模和预测。

优缺点

- **优点:**
 - **预测精度高:** 通过将多个弱学习器按照特定的方式组合起来，梯度提升回归模型能够捕捉到数据中非常复杂的非线性关系。这使得它在很多回归问题上，都能够展现出非常高的预测精度，能够准确地预测出目标变量的值。
 - **灵活性强:** 它可以使用不同类型的弱学习器，比如决策树、线性模型等，而且对于数据的分布没有严格的要求。无论是数据呈现正态分布，还是其他各种不规则的分布，梯度提升回归模型都有可能适用，这使得它在处理各种类型的数据时都具有很大的优势。
 - **鲁棒性较好:** 在一定程度上，它对噪声数据具有一定的鲁棒性。因为它是通过一轮又一轮的迭代来逐步改进模型的，所以个别异常的数据点不会对整个模型的训练结果产生过大的影响。即使数据中存在一些噪声，模型也能够通过迭代的方式，逐渐调整和优化，从而保持相对稳定的性能。
- **缺点:**
 - **计算成本高:** 由于梯度提升回归模型需要迭代训练多个弱学习器，每一次迭代都需要进行大量的计算，所以它的训练过程计算量比较大。特别是当数据集规模较大，或者 `n_estimators` 这个参数设置得比较大的时候，模型的训练时间会明显增加，对计算资源的要求也会更高。
 - **容易过拟合:** 如果在设置超参数的时候不合理，比如 `n_estimators` 设置得过大，`max_depth` 设置得过深等，模型就很容易出现过拟合的情况。过拟合的模型虽然在训练数据上表现出色，但在面对新的数据时，往往无法准确地进行预测，泛化能力较差。
 - **对异常值敏感:** 虽然前面提到它在一定程度上对噪声有鲁棒性，但如果数据中存在大量的异常值，还是会对模型的性能产生较大的影响。因为在模型训练的早期阶段，弱学习器可能会过度关注这些异常值，从而影响整个模型的学习和预测效果。

程序思路说明

该程序旨在实现一个基于梯度提升回归的深度预测系统，整体思路如下：

1. 数据读取：

- 通过指定路径从Excel文件读取数据。若文件不存在或读取出错，给出相应提示。

2. 数据预处理：

- 筛选出深度在特定范围内的数据。
- 将数据转换为数值类型，对缺失值使用均值填充。

3. 确定不同深度范围参数：

- 定义不同深度区间及其对应的各类参数。这些参数后续用于与预测结果关联展示。

4. 模型训练与调优：

- 划分训练集和测试集。
- 对每个预测目标（y的各列），使用网格搜索交叉验证来寻找梯度提升回归模型的最优参数。
- 记录每个预测目标对应的最优模型。

5. 单位转换：

- 定义预测结果单位转换函数，当前仅返回预测值，可按需修改。

6. 预测：

- 用户输入深度值，程序判断其有效性。
- 找到对应深度的最优模型进行预测，并转换预测结果单位。
- 根据深度所在范围，从预定义参数中筛选相关参数，展示预测结果。

7. 异步预测：

- 使用线程实现异步预测，防止界面卡顿。
- 实时更新进度条，显示预测进度。

8. 创建GUI：

- 创建图形用户界面，包含输入框、按钮、进度条和输出文本框。
- 布局各组件，实现用户与程序的交互。

9. 主程序：

- 执行上述各步骤，读取数据、训练模型、评估模型性能，最后启动GUI。

WPS Office 找稿壳模板 zongbiao.xlsx

开始 插入 页面 公式 数据 审阅 视图 工具 会员专享 效率 WPS AI

深度 钻井参数 补偿中子 声波时差 自然伽马 光电吸收界面指数

	A	B	C	D	E	F	G	H	I	J	K	L	M	N
1	深度	井径	补偿中子	声波时差	自然伽马	光电吸收界面指数	深度预测系统							
2	4399.7	8.7408	32566	67.7396	101.301	14.8437	输入深度值	5555	预测结果:	自然电位	密度校正	岩性密度校正	岩石电阻率1	岩石电阻率
3	4399.8	8.7206	12636	67.1166	94.9837	15.1055				测井深度校正				
4	4399.9	8.7472	97956	66.4173	95.2985	15.4008				-20.0114	844200	2.4812	-	-
5	4400.0	8.7690	90100	65.6899	99.4178	15.6246				-19.9367	786700	2.4765	-	-
6	4400.1	8.7866	86749	65.1785	106.12	15.6807				-19.8505	714900	2.4787	-	-
7	4400.2	8.7949	0.081818	64.88	112.8741	15.4887	井径: 6.708816	补偿中子: 18.167895		-19.7563	646000	2.4775	-	-
8	4400.3	8.7949	0.084592	65.0126	116.6975	15.1471	声波时差: 55.100623			-19.6729	592400	2.4738	-	-
9	4400.4	8.7871	0.093271	65.3243	117.9616	14.772	自然伽马: 294.108721			-19.6561	-0.5631	2.4648	-	-
10	4400.5	8.7415	0.083647	65.7035	113.4521	14.5123	光电吸收界面指数: 3.732071			-19.7387	-0.5513	2.457	-	-
11	4400.6	8.7116	0.080307	66.0999	109.1712	14.4275	深侧向电阻率: 192.413413			-19.9032	-0.5512	2.4504	-	-
12	4400.7	8.6974	0.08325	66.4022	105.1187	14.4956	微球型聚焦电阻率: 24.908133			-20.0757	-0.5577	2.4565	-	-
13	4400.8	8.7195	0.086999	66.6145	103.9531	14.6529	浅侧向电阻率: 168.098740			-20.1653	-0.5655	2.4561	-	-
14	4400.9	8.7398	0.087654	66.59	101.0284	14.8164	密度校正: 1604.505166			-20.147	-0.5719	2.4495	-	-
15	4401.0	8.7568	0.084532	66.495	95.8849	14.9523	岩性密度校正: 2.560020			-20.0398	-0.575	2.4519	-	-
16	4401.1	8.7341	0.084508	66.3682	93.8753	15.0615				-19.8771	-0.5742	2.4543	-	-
17	4401.2	8.7214	0.084964	66.292	93.4536	15.1508				-19.688	-0.5704	2.4563	-	-
18	4401.3	8.7278	0.085777	66.3045	95.0158	15.2074				-19.4924	-0.5647	2.4525	-	-
19	4401.4	8.7666	0.091958	66.3873	99.2267	15.1977				-19.3243	-0.5591	2.4489	-	-
										-19.2114	-0.5552	2.4461	-	-
										-19.1629	-0.5534	2.4559	-	-

Sheet1 +

平均值=4991.7612413443 计数=1万5742 求和=7857万5313.7