机器学习实战:从 Scikit-learn 到 XGBoost (卷一)

前言

欢迎踏上机器学习的精彩旅程!在数据驱动的时代,机器学习已成为解锁数据价值、驱动智能决策的核心技术。无论你是希望系统学习机器学习的学生、渴望技能升级的开发者,还是寻求解决复杂业务问题的数据分析师,掌握机器学习的基本原理和实用工具都至关重要。

本教程(或称之为"小书")旨在为你提供一份全面而深入的指南,重点聚焦于 Python 生态中最流行、最强大的机器学习库: Scikit-learn 和 XGBoost。我们的目标不仅仅是罗列 API 或展示零散的代码片段,而是希望通过系统性的讲解、结合实际应用场景(例如你提供的天气预测相关代码),帮助你建立扎实的理论基础,掌握核心的编程技能,最终能够独立地分析问题、处理数据、构建模型并评估结果。

我们将从最基础的概念讲起,逐步深入到模型的选择、训练、评估、优化,以及集成学习(如 Stacking)、模型微调(Fine-tuning)等进阶主题。教程中会穿插大量的概念解释、关键参数分析以及代码示例(力求精练实用),并特别关注在实际项目中可能遇到的问题,如处理不平衡数据、应对大数据挑战等。

这不仅仅是一本关于"如何使用工具"的手册,更希望成为你理解"为何如此"以及"如何做得更好"的向导。通过学习本教程,你将有望达到能够理解并编写出类似提供的 8 个 Python 脚本(从数据准备、模型训练到集成预测)的水平。

这是一段需要耐心和实践的旅程。准备好你的 Python 环境,让我们一起开始探索机器学习的奥秘吧!

本书 (教程) 特点:

- 系统性: 从基础概念到进阶技术,结构清晰,循序渐进。
- **实用性**: 聚焦 Scikit-learn 和 XGBoost 两大核心库,紧密结合实际应用场景。
- 深度性: 不止于 API 调用,深入讲解关键参数、模型原理和评估指标。
- 中文语境: 使用流畅自然的中文表达, 力求易于理解。
- 目标导向: 旨在帮助读者掌握足以应对真实世界机器学习任务的知识和技能。

阅读建议:

- 建议按章节顺序阅读,确保概念的连贯性。
- 动手实践是关键,请务必运行代码示例,并尝试修改参数进行实验。
- 遇到不理解的概念或代码,可以查阅官方文档或通过搜索引擎深入学习。
- 结合你自己的项目或感兴趣的数据集进行练习,是巩固知识的最佳方式。

版本说明:

本教程基于编写时较为稳定和流行的库版本,如 Python 3.x, Scikit-learn 1.x, XGBoost 1.x 等。虽然库会不断更新,但核心概念和 API 设计通常具有较好的向后兼容性。

目录 (初步规划)

卷一

前言

- 目录
- 第一章: 机器学习启程
 - 1.1 什么是机器学习?
 - 1.2 机器学习的种类
 - 1.2.1 监督学习 (Supervised Learning)
 - 1.2.2 无监督学习 (Unsupervised Learning)
 - 1.2.3 强化学习 (Reinforcement Learning)
 - 1.2.4 本教程的重点: 监督学习
 - 1.3 核心概念: 特征、标签、模型
 - 1.4 典型机器学习工作流
 - 1.5 为何选择 Python、Scikit-learn 和 XGBoost?
 - 1.6 本书 (教程) 目标与结构
- 第二章: 环境搭建与工具准备
 - 2.1 Python 的选择与安装: 拥抱 Anaconda
 - 2.2 理解虚拟环境: 隔离的智慧
 - 2.3 创建与管理 Conda 环境
 - 2.4 安装机器学习核心库
 - 2.5 Jupyter Notebook/Lab: 交互式编程的首选
 - 2.6 集成开发环境 (IDE): VS Code
 - 2.7 代码版本控制: Git 与 GitHub

卷二 (后续生成)

第三章:数据处理基石: NumPy第四章:数据分析利器: Pandas

• 第五章: Scikit-learn 核心 API 与通用功能

卷三 (后续生成)

• 第六章: 监督学习: 分类模型详解(逻辑回归、决策树、随机森林等)

• 第七章:模型评估:度量你的模型 (混淆矩阵、准确率、召回率、F1、AUC、POD/FAR/CSI等)

卷四 (后续生成)

• 第八章:梯度提升王者: XGBoost 深入解析

• 第九章:特征工程:化数据为神奇

卷五 (后续生成)

• 第十章:集成学习进阶: Stacking 与 Blending

• 第十一章:模型调优与选择(网格搜索、随机搜索、交叉验证)

卷六 (后续生成)

• 第十二章: 处理不平衡数据

• 第十三章: 处理大型数据集策略 (内存映射、分块处理)

• 第十四章:模型微调 (Fine-tuning) 实战

• 第十五章: 案例分析: 天气预测模型代码解读

卷七 (后续生成)

• 第十六章: 模型部署与上线(简介)

• 第十七章: 总结与展望

(注:目录结构和内容可能会在后续生成过程中进行微调)

第一章: 机器学习启程

欢迎来到机器学习的世界!这一章将为你揭开机器学习的神秘面纱,介绍其基本概念、分类、典型工作流程,并解释为何 Python 及其相关库成为该领域的首选工具。

1.1 什么是机器学习?

机器学习(Machine Learning, ML)是人工智能(Artificial Intelligence, AI)的一个重要分支。与传统编程中开发者明确编写所有规则来处理数据不同,机器学习的核心思想是**让计算机系统能够从数据中"学习"并改进其性能,而无需进行显式编程**。

想象一下识别垃圾邮件的任务。传统方法可能需要程序员编写大量复杂的规则(例如,包含"免费"、"中奖"等词语,或者来自特定发件人)。但这种方法难以维护,且容易漏掉新型垃圾邮件。

机器学习的方法则不同:我们给计算机提供大量的邮件样本,并告诉它哪些是垃圾邮件(标签),哪些不是。 机器学习算法会分析这些样本,自动寻找区分垃圾邮件和非垃圾邮件的模式(特征)。学习完成后,当收到一 封新邮件时,模型就能根据它学到的模式来判断这封邮件是否是垃圾邮件。

核心要素:

- 数据 (Data): 机器学习的燃料。数据的质量和数量直接影响模型的性能。
- 算法 (Algorithm): 学习数据中模式的方法。不同的算法适用于不同的问题和数据类型。
- 模型 (Model): 算法从数据中学习到的结果,通常是一个可以进行预测的数学函数或一组规则。
- 学习 (Learning): 通过数据调整模型参数,使其在特定任务上表现更好的过程。

简单来说, 机器学习就是利用数据构建能够进行预测或决策的模型的过程。

1.2 机器学习的种类

机器学习根据学习方式和数据类型的不同,主要可以分为以下几类:

1.2.1 监督学习 (Supervised Learning)

这是最常见、应用最广泛的机器学习类型。在监督学习中,我们提供给算法的数据集包含了**输入特征** (Features) 和对应的**正确输出标签 (Labels)**。算法的目标是学习一个从输入特征到输出标签的映射函数。

- 分类 (Classification): 当输出标签是离散的类别时, 称为分类问题。
 - 示例:
 - 判断邮件是否为垃圾邮件(类别:是/否)
 - 识别图像中的动物(类别:猫/狗/鸟)
 - 预测客户是否会流失(类别:流失/不流失)
 - 本教程重点关联: 预测是否会下雨(类别:下雨/不下雨),根据特征区分 TP/TN/FP/FN 样本也属于分类任务。
- 回归 (Regression): 当输出标签是连续的数值时,称为回归问题。
 - 示例:

- 预测房价(数值:价格)
- 预测股票价格(数值:价格)
- 预测温度(数值:摄氏度)
- 本教程重点关联: 如果预测的是具体的降雨量大小(例如 5mm, 10.2mm),则属于回归问题。但教程中的例子是基于阈值判断是否下雨,将其转化为了分类问题。

1.2.2 无监督学习 (Unsupervised Learning)

与监督学习不同,无监督学习的数据集**没有提供输出标签**。算法需要自己发现数据中隐藏的结构或模式。

• **聚类 (Clustering):** 将数据点分成不同的组(簇),使得同一组内的数据点相似度高,不同组间的数据点相似度低。

○ 示例:

- 客户细分(将客户分成不同的购买群体)
- 社交网络分析 (发现社区结构)
- 图像分割
- **降维 (Dimensionality Reduction):** 在保留数据主要信息的前提下,减少数据的特征数量。这有助于数据可视化、提高后续模型训练效率、降低噪声。
 - 示例:
 - 主成分分析 (Principal Component Analysis, PCA)
 - t-分布随机邻域嵌入 (t-Distributed Stochastic Neighbor Embedding, t-SNE),常用于高维数据可视化。
- 关联规则学习 (Association Rule Learning): 发现数据项之间的有趣关系。
 - 示例:
 - 购物篮分析 ("购买了啤酒的顾客也倾向于购买尿布")

1.2.3 强化学习 (Reinforcement Learning)

强化学习关注的是智能体 (Agent) 如何在一个环境 (Environment) 中采取行动 (Action),以最大化累积奖励 (Reward)。智能体通过试错来学习最佳策略 (Policy)。

• 示例:

- 。 训练机器人走路
- 训练 AI 下棋 (AlphaGo)
- 。 自动驾驶策略优化
- 。 推荐系统动态调整

1.2.4 本教程的重点: 监督学习

本教程将主要聚焦于**监督学习**,特别是**分类问题**。因为你提供的代码示例(天气预测、FN/FP 专家模型)都是围绕着二元分类(下雨/不下雨,是 FN/否,是 FP/否)展开的。我们将深入探讨 Scikit-learn 和 XGBoost 在解决这类问题上的强大能力。当然,很多概念和工具(如数据处理、模型评估)同样适用于回归问题,我们也会在适当的时候提及。

1.3 核心概念:特征、标签、模型

理解以下核心概念对于学习机器学习至关重要:

• **样本 (Sample) / 实例 (Instance) / 数据点 (Data Point):** 数据集中的一条记录。例如,一次天气观测记录、一封邮件、一个客户的信息。

- 特征 (Feature) / 属性 (Attribute) / 输入变量 (Input Variable): 描述样本的特性或属性。这些是模型的输入。
 - **示例 (天气预测):** 温度、湿度、气压、风速、云量、历史降雨信息等。在你的 X_flat_features.npy 文件中,每一列代表一个特征。脚本中提到的 N_FEATURES = 100 或 101 指的就是特征的数量。
 - 。 特征可以是数值型(如温度)、类别型(如天气状况"晴天"、"阴天")或布尔型(如是否刮风)。
- 标签 (Label) / 目标变量 (Target Variable) / 输出变量 (Output Variable): 我们希望模型预测或学习的结果。这是监督学习中模型的输出。
 - **示例 (天气预测):** 是否下雨 (0 或 1),或者具体的降雨量。在你的 $Y_{flat_target_npy}$ 文件中存储的就是标签。RAIN_THRESHOLD = 0.1 用于将连续的降雨量标签转换为二元的分类标签 (大于 0.1 为 1,否则为 0)。
- 特征矩阵 (Feature Matrix): 通常用大写字母 X 表示。一个二维数组(或数据框),每一行代表一个样本,每一列代表一个特征。形状通常是 (n_samples, n_features)。
- **标签向量** (Label Vector): 通常用小写字母 y 表示。一个一维数组,包含每个样本对应的标签。形状通常是 (n_samples,)。
- 模型 (Model): 机器学习算法通过学习数据 (X, y) 后得到的产物。它可以接受新的、未见过的特征数据 X_new, 并预测其对应的标签 y_pred。
 - **示例:** 训练好的 XGBoost 分类器 (.joblib 文件存储的就是训练好的模型对象)、训练好的逻辑回归模型。

1.4 典型机器学习工作流

一个完整的机器学习项目通常遵循以下步骤,尽管实际过程往往是迭代和循环的:

1. 定义问题 (Define Problem):

- · 明确业务目标:希望解决什么问题? (例如,提高降雨预测准确率,减少漏报或误报)
- 。 确定问题类型: 是分类、回归还是其他?
- 选择评估指标:如何衡量模型的成功? (例如,准确率、召回率、CSI、FAR)

2. 收集数据 (Collect Data):

- o 获取与问题相关的数据。数据来源可能包括数据库、日志文件、传感器、API、公开数据集等。
- 3. 准备数据 (Prepare Data): 这是机器学习流程中最耗时但也至关重要的环节。
 - o 数据清洗 (Data Cleaning): 处理缺失值、异常值、重复值。
 - 数据探索 (Exploratory Data Analysis, EDA): 理解数据分布、特征之间的关系、发现潜在问题。
 可视化是 EDA 的重要手段。
 - 特征工程 (Feature Engineering):
 - 特征选择 (Feature Selection): 选择最相关的特征。
 - 特征提取 (Feature Extraction): 从原始数据中提取更有用的特征(例如,从日期中提取星期几)。
 - 特征转换 (Feature Transformation): 对特征进行缩放(标准化、归一化)、编码(处理类别特征)、组合等。
 - o **数据划分 (Data Splitting):** 将数据划分为训练集 (Training Set)、验证集 (Validation Set) 和测试集 (Test Set)。

- 训练集: 用于训练模型。
- **验证集**: 用于调整模型超参数(Hyperparameters)和初步评估模型性能,避免在测试集上过拟合。
- **测试集**: 用于最终评估模型在未见过数据上的泛化能力。**测试集绝对不能用于模型训练或调 优过程。** 你的代码中 TEST_SIZE_RATIO = 0.2 就是在进行训练集和测试集的划分。

4. 选择模型 (Choose Model):

根据问题类型、数据特点、计算资源等因素选择合适的机器学习算法。可能需要尝试多种模型。

5. 训练模型 (Train Model):

○ 使用**训练集**数据 (X_train, y_train) 和所选算法来学习模型参数。调用模型的 fit() 方法。

6. 评估模型 (Evaluate Model):

- 使用**验证集**或通过**交叉验证 (Cross-Validation)** 来评估模型性能,根据选定的评估指标进行度量。调用 predict() 或 predict_proba() 后计算指标。
- o 分析错误类型 (例如,查看混淆矩阵,识别哪些样本被错误分类)。你的脚本 1_generate_base_predictions.py 中识别 TP/TN/FP/FN 就是评估和错误分析的一部分。

7. 调优模型 (Tune Model):

- 根据评估结果调整模型的超参数(例如, XGBoost 的 max_depth, learning_rate)。可以使用网格搜索(Grid Search)、随机搜索(Random Search)等方法自动化调优过程。
- 。 可能需要回到特征工程步骤, 改进特征。
- 。 可能需要尝试不同的模型算法。
- 模型微调 (6_finetune_base_model.py) 也是一种调优策略。

8. 最终评估 (Final Evaluation):

○ 使用**测试集**对最终选定和调优好的模型进行一次性评估,得到模型泛化能力的客观估计。

9. 部署模型 (Deploy Model):

将训练好的模型集成到实际应用中,对外提供预测服务。例如,构建 API 接口。

10. 监控与维护 (Monitor & Maintain):

- 模型上线后需要持续监控其性能,因为数据分布可能随时间变化(概念漂移)。
- 。 根据监控结果定期重新训练或更新模型。

这个流程并非严格线性,经常需要在不同步骤之间来回迭代。例如,评估结果不理想可能需要重新进行特征工程或选择不同模型。

1.5 为何选择 Python、Scikit-learn 和 XGBoost?

Python 已经成为数据科学和机器学习领域的首选语言,主要原因包括:

- 简洁易学: Python 语法清晰, 上手相对容易。
- 强大的生态系统: 拥有海量用于数据处理、科学计算、机器学习和可视化的库。
- 社区活跃: 庞大的开发者社区意味着丰富的学习资源、活跃的讨论和持续的库更新。
- 胶水语言: 能够轻松集成其他语言(如 C/C++)编写的高性能模块。

Scikit-learn:

• 机器学习瑞士军刀: 提供了非常全面且一致的机器学习工具,涵盖了数据预处理、特征选择、模型训练(分类、回归、聚类等多种算法)、模型评估和选择等各个环节。

- API 设计优雅: 核心 API (fit, predict, transform) 设计一致, 易于学习和组合使用。
- 文档完善: 拥有高质量的官方文档和丰富的示例。
- 广泛使用: 是工业界和学术界最流行的通用机器学习库之一。

XGBoost:

- 梯度提升大杀器: 是梯度提升树 (Gradient Boosting Tree) 算法的高效、灵活和可移植实现。
- **性能卓越**: 在许多结构化数据(表格数据)的分类和回归任务中表现顶尖,是 Kaggle 等数据科学竞赛中的常胜将军。
- 速度快、内存占用优化: 实现了多种优化技术,包括并行处理、缓存优化、稀疏数据感知和近似分割算法 (tree_method='hist')。
- 内置正则化: 有助于防止过拟合。
- 灵活性高: 提供丰富的可调参数,允许用户精细控制模型行为。
- 原生支持缺失值: 能够自动处理数据中的缺失值。

结合这三者,我们可以构建一个强大、高效且灵活的机器学习开发环境,足以应对从简单到复杂的各种任务。你的代码示例充分体现了这一点,以 Python 为基础,利用 Scikit-learn 进行数据划分、模型评估(虽然评估代码在脚本中重复出现,但原理来自 Scikit-learn),并深度使用了 XGBoost 作为核心模型进行训练、预测,甚至基于 XGBoost 构建了更复杂的 Stacking 结构。

1.6 本书 (教程) 目标与结构

目标:

- 让你深入理解监督学习(特别是分类)的核心概念和流程。
- 熟练掌握 NumPy, Pandas 进行数据处理的基础操作。
- 熟练掌握 Scikit-learn 的核心 API,包括数据划分、预处理、模型训练(常用分类器)、评估指标计算和模型持久化。
- 深入理解 XGBoost 的工作原理、关键参数调优、早停、特征重要性等。
- 理解并能够实现 Stacking 集成学习策略,特别是理解示例代码中 FN/FP 专家模型的设计思想。
- 掌握模型微调 (Fine-tuning) 的概念和在 XGBoost 中的应用。
- 了解处理大型数据集的基本策略,如内存映射和分块处理。
- 最终具备独立分析、构建和评估机器学习模型(尤其是分类模型)的能力,能够读懂并编写类似提供的 天气预测脚本。

结构:

如前 目录 所示, 我们将遵循一个从基础到进阶的逻辑结构:

- 1. 入门与准备:介绍机器学习基础,搭建开发环境。
- 2. 数据处理基础: 学习 NumPy 和 Pandas, 为后续操作打下基础。
- 3. Scikit-learn 核心: 掌握 Scikit-learn 的通用功能和常用分类模型。
- 4. 模型评估: 深入理解各种评估指标及其选择。
- 5. XGBoost 深入: 详细解析 XGBoost 的原理、参数和使用技巧。
- 6. 特征工程: 探讨如何创造更好的特征。
- 7. **集成学习与调优:** 学习 Stacking、模型调优方法。

- 8. 实战技巧: 处理不平衡数据、大数据、模型微调。
- 9. 案例与总结:结合天气预测案例进行代码解读,并对整个教程进行总结。

每一章都会包含理论讲解、代码示例和必要的解释,力求将知识点讲透彻、讲明白。

第二章:环境搭建与工具准备

工欲善其事,必先利其器。在正式开始机器学习编程之前,我们需要搭建一个稳定、高效且易于管理的开发环境。本章将指导你完成必要的软件安装和环境配置。

2.1 Python 的选择与安装: 拥抱 Anaconda

Python 是我们进行机器学习的主要编程语言。虽然你可以从 Python 官网 (python.org) 下载安装标准 Python, 但对于数据科学和机器学习领域,强烈推荐使用 **Anaconda** 发行版。

什么是 Anaconda?

Anaconda 是一个包含了 Python 解释器、众多常用科学计算和数据科学库(如 NumPy, Pandas, SciPy, Matplotlib, Jupyter 等)以及强大的包管理和环境管理工具 (conda) 的**免费开源发行版。**

为什么推荐 Anaconda?

- 1. **一站式安装:** 安装 Anaconda 会自动为你安装好 Python 和数百个常用的科学计算库,省去了手动逐一安装的麻烦和可能遇到的依赖问题。
- 2. **强大的环境管理:** conda 工具可以轻松创建、切换和管理多个独立的 Python 环境。这对于处理不同项目可能依赖不同库版本的情况至关重要。
- 3. **便捷的包管理:** conda 不仅能管理 Python 包,还能管理非 Python 的依赖库(如 C/C++ 库),解决了pip (Python 官方包管理器)在某些情况下难以处理的依赖问题。
- 4. **跨平台:** 支持 Windows, macOS 和 Linux。

如何安装 Anaconda?

- 1. 访问 Anaconda 官方网站: https://www.anaconda.com/products/distribution
- 2. 根据你的操作系统 (Windows/macOS/Linux) 下载对应的图形化安装包 (推荐)。
- 3. 运行安装包,按照提示进行安装。关键步骤:
 - 。 同意许可协议。
 - 选择安装路径: 建议选择一个空间充足、路径不含中文或空格的目录。
 - 高级选项 (Windows):
 - "Add Anaconda to my PATH environment variable" (不推荐勾选): 官方不推荐勾选此项,以避免与其他 Python 安装冲突。我们将通过 Anaconda Navigator 或 Anaconda Prompt (命令行) 来启动工具和环境。
 - "Register Anaconda as my default Python [version]" (推荐勾选): 勾选此项可以让你在支持的 IDE (如 VS Code) 中更容易识别和使用 Anaconda 的 Python 解释器。
 - 完成安装。

安装完成后,你可以在开始菜单(Windows)或应用程序文件夹(macOS)中找到 "Anaconda Navigator"(图形化界面)和 "Anaconda Prompt"(命令行终端)。

2.2 理解虚拟环境:隔离的智慧

想象一下你同时在开发两个项目:

- 项目 A 需要使用库 library_x 的 1.0 版本。
- 项目 B 需要使用同一个库 library_x 的 2.0 版本。

如果你将所有库都安装在全局(基础) Python 环境中,这两个项目就会产生冲突。虚拟环境就是为了解决这个问题而生的。

虚拟环境 (Virtual Environment) 是一个独立的、隔离的 Python 运行环境。每个虚拟环境都有自己独立的 Python 解释器、库和脚本。在一个环境中安装、卸载或更新库,不会影响到其他环境。

使用虚拟环境的好处:

- 依赖隔离: 避免不同项目之间的库版本冲突。
- 环境复现: 可以轻松地将项目的依赖列表导出 (environment.yml 或 requirements.txt), 方便他人在不同机器上复现完全相同的环境。
- 保持基础环境清洁: 避免将各种库随意安装到基础环境中, 保持其整洁和稳定。

conda 是管理虚拟环境的利器。

2.3 创建与管理 Conda 环境

我们将使用 Anaconda Prompt (Windows) 或终端 (macOS/Linux) 来执行 conda 命令。

- 1. 打开 Anaconda Prompt 或终端。
- 2. 查看已有环境:

```
conda env list
# 或者
conda info --envs
```

你会看到一个名为 base 的环境,这是 Anaconda 安装时自带的基础环境。星号*表示当前激活的环境。

3. **创建新的虚拟环境**: 假设我们要为本教程创建一个名为 ml_book 的环境,并指定使用 Python 3.9 (你可以选择其他版本,如 3.8, 3.10 等)。

```
conda create --name ml_book python=3.9
```

- o conda create: 创建环境的命令。
- --name ml book: 指定环境的名称。
- o python=3.9: 指定环境中要安装的 Python 版本。

conda 会计算依赖关系并提示你将要安装的包,输入 y 并回车确认。

- 4. 激活环境: 创建环境后,需要先激活才能使用它。
 - Windows:

conda activate ml_book

o macOS / Linux:

conda activate ml_book

激活成功后,命令行提示符前面会显示环境名称 (ml_book)。

- 5. **在环境中安装库**: 激活环境后,使用 conda install 或 pip install 安装所需的库。我们将在下一节进行安装。
- 6. 查看当前环境已安装的库:

conda list

7. 退出当前环境 (返回 base 环境):

conda deactivate

命令行提示符前面的 (ml_book) 会消失。

8. 删除环境(谨慎操作): 如果不再需要某个环境,可以将其删除。确保你已退出该环境。

 $\verb|conda| env remove --name ml_book|$

强烈建议:为你的每个机器学习项目创建一个独立的 Conda 环境。

2.4 安装机器学习核心库

现在,让我们激活刚刚创建的 ml_book 环境,并安装本教程所需的核心库。

1. 激活环境:

conda activate ml_book

2. **安装 NumPy, Pandas, Matplotlib, Scikit-learn:** 这些库通常可以通过 conda 直接安装, conda 会处理好复杂的依赖关系。

conda install numpy pandas matplotlib scikit-learn joblib

- o numpy: Python 科学计算的基础库,提供强大的 N 维数组对象和相关函数。
- o pandas:基于 NumPy 构建,提供高性能、易用的数据结构(如 DataFrame)和数据分析工具。
- o matplotlib:用于绘制各种静态、动态、交互式图表的基础库。
- o scikit-learn: 核心的通用机器学习库。
- 。 joblib: Scikit-learn 常用的依赖库,尤其用于模型持久化(保存/加载模型),也提供简单的并行计算功能。
- 3. **安装 XGBoost:** XGBoost 也可以通过 conda 安装,通常推荐从 conda-forge 渠道安装,它通常包含最新的稳定版本。

```
conda install -c conda-forge xgboost
```

- -c conda-forge: 指定从 conda-forge 渠道安装。
- 4. 安装 Jupyter Notebook/Lab (推荐): 用于交互式开发和数据探索。

```
conda install jupyterlab notebook
```

- jupyterlab: 下一代 Jupyter 界面,功能更强大。
- o notebook: 经典的 Jupyter Notebook 界面。
- 5. 验证安装: 安装完成后,可以简单验证一下。在激活的 ml_book 环境中,启动 Python 解释器:

```
python
```

然后尝试导入这些库:

```
import numpy
import pandas
import sklearn
import xgboost
import joblib
import matplotlib
import jupyterlab
print("NumPy version:", numpy.__version__)
print("Pandas version:", pandas.__version__)
print("Scikit-learn version:", sklearn.__version__)
print("XGBoost version:", xgboost.__version__)
print("Joblib version:", joblib.__version__)
print("Matplotlib version:", matplotlib.__version__)
print("JupyterLab loaded successfully")
exit() # 退出 Python 解释器
```

如果所有库都能成功导入并打印版本号,说明安装成功。

2.5 Jupyter Notebook/Lab: 交互式编程的首选

Jupyter Notebook 和 JupyterLab 提供了一个基于 Web 的交互式计算环境,允许你创建和共享包含实时代码、公式、可视化和叙述文本的文档。它们在数据探索、模型原型设计和教学演示中非常受欢迎。

启动 JupyterLab:

- 1. 确保你的 ml_book 环境已激活。
- 2. 在 Anaconda Prompt 或终端中输入:

```
jupyter lab
```

这会在你的默认浏览器中打开 JupyterLab 界面。

启动 Jupyter Notebook (经典界面):

- 1. 确保你的 ml_book 环境已激活。
- 2. 在 Anaconda Prompt 或终端中输入:

```
jupyter notebook
```

这会在你的默认浏览器中打开 Jupyter Notebook 界面。

基本使用:

- 界面: 左侧是文件浏览器,右侧是工作区。
- **创建 Notebook:** 点击 "File" -> "New" -> "Notebook", 选择 Python 3 (ipykernel) 内核。
- 单元格 (Cell): Notebook 由单元格组成。主要有两种类型:
 - Code Cell: 用于编写和运行 Python 代码。按 Shift + Enter 运行当前单元格并将焦点移到下一个单元格,按 Ctrl + Enter 运行当前单元格并保持焦点。
 - o Markdown Cell: 用于编写格式化的文本、标题、列表、公式等 (使用 Markdown 语法)。
- 内核 (Kernel): 负责运行单元格中的代码。你可以中断、重启内核。
- 保存: Notebook 会自动保存,也可以手动保存 (Ctrl + S)。文件扩展名为 .ipynb。

我们将在后续章节的代码示例中大量使用 Notebook/Lab。

2.6 集成开发环境 (IDE) : VS Code

虽然 Jupyter 非常适合交互式探索,但对于开发更大型、结构化的 Python 脚本(就像你提供的 .py 文件),使用一个功能更全面的集成开发环境 (IDE) 会更高效。**Visual Studio Code (VS Code)** 是一个免费、开源且功能强大的代码编辑器,对 Python 开发提供了极佳的支持。

安装 VS Code:

- 访问 VS Code 官网: https://code.visualstudio.com/
- 下载并安装适用于你操作系统的版本。

配置 Python 开发环境:

1. 安装 Python 扩展:

- 。 打开 VS Code。
- 。 点击左侧边栏的扩展图标 (Extensions)。
- 搜索 "Python" (由 Microsoft 发布)。
- 。 点击 "Install" 安装。这个扩展提供了代码补全 (IntelliSense)、调试、Linting (代码风格检查)、环境管理等核心功能。

2. 安装 Jupyter 扩展 (推荐):

- 。 在扩展市场搜索 "Jupyter" (由 Microsoft 发布)。
- 点击 "Install" 安装。这个扩展允许你在 VS Code 中直接打开、编辑和运行.ipynb 文件,提供了 类似 Jupyter Notebook/Lab 的体验,并能更好地与.py 脚本开发集成。

3. 选择 Python 解释器 (关联 Conda 环境):

- 打开一个 Python 文件 (.py) 或一个 Notebook 文件 (.ipynb)。
- 。 点击 VS Code 状态栏右下角的 Python 版本号(或者可能显示 "Select Interpreter")。
- 。 VS Code 会自动检测你系统中的 Python 环境,包括 Conda 环境。
- 从列表中选择你创建的 ml_book 环境对应的 Python 解释器。这确保了 VS Code 在运行和调试代码时使用的是正确的环境和库。
- 如果 VS Code 没有自动检测到你的 Conda 环境,可能需要手动配置 settings.json 文件,但这通常在正确安装 Anaconda 并勾选了相关选项后是不需要的。

VS Code 优势:

- 智能代码补全与提示: 提高编码效率。
- 强大的调试功能: 设置断点, 单步执行, 查看变量值。
- 集成终端: 可以直接在 VS Code 中打开终端并运行命令(如 conda activate)。
- Git 集成: 内置了对 Git 版本控制的支持。
- 丰富的扩展: 可以通过安装扩展来增强功能(例如, Docker, Pylint, Black Formatter 等)。

你可以选择主要使用 Jupyter Lab 进行探索和原型设计,然后使用 VS Code 来编写、组织和调试最终的 Python 脚本。

2.7 代码版本控制: Git 与 GitHub

虽然不是严格意义上的机器学习工具,但使用版本控制系统(如 Git)和代码托管平台(如 GitHub)对于管理机器学习项目至关重要。

- Git: 一个分布式版本控制系统,可以跟踪代码的修改历史,方便回溯、协作和分支管理。
- **GitHub (或其他平台如 GitLab, Bitbucket):** 一个基于 Web 的 Git 仓库托管服务,方便备份代码、团队协作和开源分享。

为什么需要版本控制?

- 记录历史: 轻松查看代码的每次修改, 知道谁、何时、为何做了更改。
- 版本回溯: 当代码出错或实验效果不佳时,可以轻松恢复到之前的某个工作版本。
- **分支管理**: 可以创建不同的分支来尝试新的功能或实验,而不影响主代码。例如,可以为不同的模型调优策略创建不同的分支。
- 团队协作: 多人可以同时在同一个项目上工作, Git 负责合并各自的修改。
- 代码备份: 将代码推送到远程仓库 (如 GitHub) 可以有效防止本地数据丢失。

学习资源:

建议花时间学习 Git 的基本命令 (git clone, git add, git commit, git push, git pull, git branch, git checkout, git merge)。 网上有大量优秀的 Git 教程。

实践建议:

- 为你的机器学习项目初始化 Git 仓库。
- 将你的代码 (.py 脚本, .ipynb 文件) 、环境配置文件 (environment.yml 或 requirements.txt) 和 其他重要文件纳入版本控制。
- **注意**: 通常不建议将大型数据文件或模型文件 (.joblib, .npy) 直接提交到 Git 仓库,因为它们体积庞大 且不适合进行差异比较。可以使用 .gitignore 文件来忽略这些文件。对于数据和模型,可以考虑使用 专门的数据版本控制工具(如 DVC)或云存储。
- 经常提交你的修改,并编写清晰的提交信息。
- 将你的本地仓库关联到 GitHub 等远程仓库,并定期推送。

掌握环境搭建、工具使用和版本控制,将为你的机器学习项目打下坚实的基础,让后续的学习和开发更加顺畅 高效。

```markdown

# 机器学习实战: 从 Scikit-learn 到 XGBoost (卷二)

## 第三章: 数据处理基石: NumPy

在深入机器学习算法之前,我们必须掌握处理数据的基本工具。NumPy (Numerical Python) 是 Python 科学计算生态系统的核心库,它为处理大型、多维数组和矩阵提供了高效的数据结构和丰富的函数库。几乎所有 Python 数据科学库,包括 Pandas 和 Scikit-learn,都构建在 NumPy 之上或与之紧密集成。

### 3.1 NumPy 为何如此重要?

- \* \*\*高效的数组运算:\*\* NumPy 的核心是 `ndarray` (N-dimensional array) 对象。这些数组在内存中是连续存储的,并且许多底层运算是用 C 或 Fortran 实现的,使得 NumPy 在处理大型数值数据集时远比 Python 内置的列表 (list) 高效得多。
- \* \*\*向量化操作: \*\* NumPy 支持向量化操作,允许你对整个数组执行操作,而无需编写显式的 Python 循环。这使得代码更简洁、更易读,并且执行速度更快。
- \* \*\*广播功能: \*\* 能够智能地处理不同形状数组之间的运算。
- \* \*\*广泛的数学函数库: \*\* 提供了大量的线性代数、傅里叶变换、随机数生成等数学函数。
- \* \*\*生态系统基础:\*\* 作为 Pandas, SciPy, Matplotlib, Scikit-learn 等库的基础, 理解 NumPy 是理解和高效使用这些库的前提。

在你的天气预测脚本中,特征数据(`X\_flat\_features.npy`)和目标数据(`Y\_flat\_target.npy`)都被存储和加载为 NumPy 数组。模型预测的概率、索引等也都是通过 NumPy 数组来处理和存储的。

### 3.2 核心数据结构: `ndarray`

`ndarray` 是 NumPy 的中心数据结构。它是一个包含相同类型元素的多维网格。

#### 3.2.1 创建数组

有多种方式可以创建 NumPy 数组:

```
从 Python 列表或元组创建:
 ```python
 import numpy as np
 # 一维数组
 arr1d = np.array([1, 2, 3, 4, 5])
 print(arr1d)
 # 输出: [1 2 3 4 5]
 # 二维数组 (矩阵)
 arr2d = np.array([[1, 2, 3], [4, 5, 6]])
 print(arr2d)
 # 输出:
 # [[1 2 3]
 # [4 5 6]]
**使用内置函数创建特定数组:**
 ```python
 # 创建全零数组
 zeros_arr = np.zeros((2, 3)) # 参数是形状元组 (shape)
 print(zeros_arr)
 # 输出:
[[0. 0. 0.]
[0. 0. 0.]]
创建全一数组
 ones_arr = np.ones((3, 2), dtype=int) # 可以指定数据类型 dtype
 print(ones_arr)
 # 输出:
[[1 1]
[1 1]
[1 1]]
 # 创建等差序列数组 (类似 Python range)
 range_arr = np.arange(0, 10, 2) # start, stop (exclusive), step
 print(range arr)
 # 输出: [0 2 4 6 8]
 # 创建等间隔序列数组
 linspace_arr = np.linspace(0, 1, 5) # start, stop (inclusive), num_points
 print(linspace_arr)
 # 输出: [0. 0.25 0.5 0.75 1.]
 # 创建随机数数组 (0到1之间的均匀分布)
 rand_arr = np.random.rand(2, 2)
 print(rand_arr)
 # 输出: (每次运行结果不同)
 # [[0.123 0.456]
[0.789 0.012]]
 # 创建随机数数组 (标准正态分布)
 randn arr = np.random.randn(2, 3)
```

```
print(randn_arr)
 # 输出: (每次运行结果不同)
 # [[-0.5... 1.2... 0.3...]
 # [0.8... -1.0... -0.1...]]
3.2.2 数组属性
`ndarray` 对象包含一些描述其自身的有用属性:
 `shape`:数组的维度大小,返回一个元组。例如 `(n_rows, n_columns)`。
 `dtype`: 数组元素的数据类型 (如 `int64`, `float64`, `bool`) 。
 `ndim`: 数组的维数 (轴的数量)。
 `size`:数组中元素的总数。
```python
arr = np.array([[1., 2., 3.], [4., 5., 6.]])
print("Shape:", arr.shape) # 输出: Shape: (2, 3)
print("Data type:", arr.dtype) # 输出: Data type: float64
print("Dimensions:", arr.ndim) # 输出: Dimensions: 2
print("Size:", arr.size) # 输出: Size: 6
```

3.3 索引与切片: 访问数组元素

NumPy 提供了灵活的方式来访问和修改数组中的元素。

3.3.1 基本索引

与 Python 列表类似,使用方括号 [] 进行索引,索引从 0 开始。

```
arr1d = np.arange(10) # [0 1 2 3 4 5 6 7 8 9]
print(arr1d[0]) # 输出: 0
print(arr1d[5]) # 输出: 5
print(arr1d[-1]) # 输出: 9 (负数索引从末尾开始)

arr2d = np.array([[1, 2, 3], [4, 5, 6]])
print(arr2d[0, 1]) # 输出: 2 (第0行, 第1列)
print(arr2d[1, -1]) # 输出: 6 (第1行, 最后一列)
print(arr2d[0]) # 输出: [1 2 3] (获取第0行, 返回一个一维数组)
```

3.3.2 切片 (Slicing)

切片用于获取数组的子集。语法是 start:stop:step, 其中 stop 索引不包含在内。

```
arr1d = np.arange(10) # [0 1 2 3 4 5 6 7 8 9]
print(arr1d[2:5]) # 输出: [2 3 4] (索引2到4)
print(arr1d[:3]) # 输出: [0 1 2] (从开头到索引2)
print(arr1d[5:]) # 输出: [5 6 7 8 9] (从索引5到末尾)
print(arr1d[::2]) # 输出: [0 2 4 6 8] (步长为2)
```

```
arr2d = np.array([[1, 2, 3], [4, 5, 6], [7, 8, 9]])
print(arr2d[:2, 1:]) # 输出: [[2 3] [5 6]] (前2行, 第1列及之后)
print(arr2d[:, 0]) # 输出: [1 4 7] (所有行的第0列)
```

重要特性:数组切片是原始数组的视图 (View),不是副本 (Copy)。修改切片会影响原始数组。如果需要副本,应使用.copy()方法。

```
arr = np.arange(5)
slice_arr = arr[1:3]
print("Original arr:", arr) # 输出: Original arr: [0 1 2 3 4]
print("Slice:", slice_arr) # 输出: Slice: [1 2]

slice_arr[0] = 99
print("Modified slice:", slice_arr) # 输出: Modified slice: [99 2]
print("Original arr after modification:", arr) # 输出: Original arr after modification: [ 0 99 2 3 4]

# 使用 copy() 创建副本
copy_arr = arr[1:3].copy()
copy_arr[0] = 111
print("Modified copy:", copy_arr) # 输出: Modified copy: [111 2]
print("Original arr after copy modification:", arr) # 输出: Original arr after copy modification: [ 0 99 2 3 4]
```

3.3.3 布尔索引 (Boolean Indexing)

这是 NumPy 中非常强大和常用的功能,允许你使用布尔数组来选择元素。

```
names = np.array(['Bob', 'Sue', 'Bob', 'Will', 'Sue'])
data = np.random.randn(5, 4) # 假设是一些数据

print("Names:", names)
print("Data:\n", data)

# 使用布尔条件选择 names 数组
print(names == 'Bob') # 输出: [ True False True False False] (布尔数组)

# 使用布尔数组索引 data 数组的行
print("Data for Bob:\n", data[names == 'Bob'])
# 输出: (data中第0行和第2行)
# [[ ... ]
# [ ... ]]

# 可以组合多个条件 (&: and, |: or, ~: not)
mask = (names == 'Bob') | (names == 'Will')
print("Mask (Bob or Will):", mask) # 输出: [ True False True True False]
print("Data for Bob or Will:\n", data[mask])
```

```
# 也可以直接修改符合条件的元素
data[data < 0] = 0 # 将 data 中所有负数设为0
print("Data with non-negative values:\n", data)
```

应用关联: 在你的脚本 1_generate_base_predictions.py 中,识别 TP, TN, FP, FN 样本的过程就使用了布尔索引: is_tp = (y_train_true == 1) & (y_train_pred == 1) 创建了一个布尔数组 is_tp, 标记了哪些样本同时满足真实标签为1和预测标签为1。然后 np.where(is_tp)[0] 用来获取这些 True 值对应的索引,这些索引随后被保存 (indices_tp.npy)。

3.3.4 花式索引 (Fancy Indexing)

使用整数数组(或列表)作为索引来选择元素。这允许你选择任意顺序、任意数量的行或列。

```
arr = np.empty((8, 4))
for i in range(8):
   arr[i] = i
print("Array:\n", arr)
# 使用列表选择特定行
print("Rows 4, 3, 0, 6:\n", arr[[4, 3, 0, 6]])
# 使用负数索引
print("Rows -1, -3, -5:\n", arr[[-1, -3, -5]])
# 使用多个索引数组选择元素 (结果是一维数组)
arr2d = np.arange(32).reshape((8, 4))
print("arr2d:\n", arr2d)
print("Selected elements:", arr2d[[1, 5, 7, 2], [0, 3, 1, 2]])
# 输出: [ 4 23 29 10] (元素 (1,0), (5,3), (7,1), (2,2))
# 使用索引数组选择行,并保留所有列
print("Selected rows with all columns:\n", arr2d[[1, 5, 7, 2]])
# 输出:
# [[ 4 5 6 7]
# [20 21 22 23]
# [28 29 30 31]
# [8 9 10 11]]
# 使用索引数组选择列
print("Selected columns:\n", arr2d[:, [3, 0, 1]])
# 输出: (所有行的第3, 0, 1列)
# [[ 3 0 1]
# [7 4 5]
# ...
# [31 28 29]]
```

应用关联: 在脚本 2_train_fn_expert.py, 3_train_fp_expert.py, 6_finetune_base_model.py 和 evaluate_fp_expert_on_tp.py 中,当需要从完整的特征矩阵 X_flat 中加载特定索引(如 indices_fn, indices_tn, combined_indices, finetune_indices, indices_tp) 对应的样本时,就利用了花式索引:

X_expert_data = X_flat[combined_indices] 这行代码直接使用 combined_indices 这个整数数组(包含了 FN 和 TN 样本的索引)从 X flat 中高效地提取出所需的行,构建专家模型的训练数据。

3.4 数组运算

NumPy 的核心优势在于其高效的数组运算。

3.4.1 元素级运算

标准的算术运算符(+,-,*,/,**)可以直接应用于数组,它们会作用于数组的每个元素。

```
arr = np.array([[1., 2.], [3., 4.]])
print("arr * arr:\n", arr * arr)
print("1 / arr:\n", 1 / arr)
print("arr ** 0.5:\n", arr ** 0.5) # 开平方

arr2 = np.array([[0., 5.], [6., 0.]])
print("arr > arr2:\n", arr > arr2) # 元素级比较
# 输出:
# [[ True False]
# [False True]]
```

3.4.2 通用函数 (Universal Functions, ufuncs)

ufuncs 是对 ndarray 中的数据执行元素级操作的函数。它们是 NumPy 运算的核心,速度很快(通常是 C 实现)。

- 一元 ufuncs: 只接受一个数组参数。
 - np.sqrt(arr): 计算平方根
 - np.exp(arr): 计算指数 e^x
 - onp.log(arr), np.log10(arr): 计算自然对数、以10为底的对数
 - np.sin(arr), np.cos(arr), np.tan(arr): 三角函数
 - np.abs(arr): 计算绝对值
 - np.isnan(arr): 判断是否为 NaN (Not a Number)
 - np.isinf(arr): 判断是否为无穷大
 - np.ceil(arr), np.floor(arr), np.rint(arr): 向上取整、向下取整、四舍五入到最近整数
- 二元 ufuncs: 接受两个数组参数。
 - np.add(arr1, arr2): 元素级相加(等同于 arr1 + arr2)
 - o np.subtract(arr1, arr2): 元素级相减
 - np.multiply(arr1, arr2):元素级相乘
 - np.divide(arr1, arr2):元素级相除
 - np.maximum(arr1, arr2), np.minimum(arr1, arr2): 元素级最大/最小值
 - onp.power(arr1, arr2):元素级幂运算(等同于 arr1 ** arr2)
 - o np.greater(arr1, arr2),np.less(arr1, arr2),np.equal(arr1, arr2):元素级比较

```
arr = np.arange(1, 5)
print("sqrt:", np.sqrt(arr))
```

```
print("exp:", np.exp(arr))

x = np.random.randn(4)
y = np.random.randn(4)
print("x:", x)
print("y:", y)
print("y:", y)
print("maximum:", np.maximum(x, y)) # 计算 x 和 y 对应元素的最大值
```

应用关联: 在脚本 4_generate_meta_features.py 中,对加载的 X_test 或 X_chunk 进行 NaN/Inf 检查时,就可能用到 np.isnan() 和 np.isinf()。

3.4.3 聚合函数

这些函数对整个数组或沿某个轴执行计算,返回一个标量值或一个维度降低的数组。

- arr.sum() / np.sum(arr): 计算所有元素的和。
- arr.mean()/np.mean(arr): 计算平均值。
- arr.std()/np.std(arr):计算标准差。
- arr.var() / np.var(arr): 计算方差。
- arr.min()/np.min(arr):找到最小值。
- arr.max()/np.max(arr):找到最大值。
- arr.argmin() / np.argmin(arr): 找到最小值的索引。
- arr.argmax() / np.argmax(arr): 找到最大值的索引。
- arr.cumsum()/np.cumsum(arr):计算累积和。
- arr.cumprod()/np.cumprod(arr):计算累积积。

这些函数通常可以接受一个 axis 参数来指定沿着哪个轴进行计算。对于二维数组:

- axis=0: 沿着行操作(计算每一列的结果)。
- axis=1: 沿着列操作(计算每一行的结果)。

```
arr2d = np.array([[0, 1, 2], [3, 4, 5]])
print("Sum (all):", arr2d.sum()) # 输出: Sum (all): 15
print("Mean (all):", arr2d.mean()) # 输出: Mean (all): 2.5
print("Sum (axis=0):", arr2d.sum(axis=0)) # 输出: Sum (axis=0): [3 5 7] (列和)
print("Mean (axis=1):", arr2d.mean(axis=1))# 输出: Mean (axis=1): [1. 4.] (行均值)
```

应用关联: 在脚本 2_train_fn_expert.py 和 3_train_fp_expert.py 中计算 scale_pos_weight 时,需要统计正负样本的数量,np.sum(y_train == 0) 和 np.sum(y_train == 1) 就利用了聚合函数 sum() (作用在布尔数组上,True 计为 1,False 计为 0) 。

3.5 广播 (Broadcasting)

广播是 NumPy 中一项强大的机制,它描述了 NumPy 在执行算术运算时如何处理不同形状的数组。在满足特定规则的情况下,较小的数组会"广播"到较大数组的形状,从而实现元素级运算,而无需显式创建副本。

广播规则: 在对两个数组进行运算时,NumPy 逐维比较它们的 shape (从末尾开始比较):

- 1. 如果两个维度相等,则继续比较下一个维度。
- 2. 如果其中一个维度是 1,则 NumPy 会在该维度上扩展(复制)这个数组,使其与另一个数组的维度匹配。

3. 如果两个维度不相等, 且没有一个维度是 1, 则会引发错误。

```
arr = np.arange(12).reshape((3, 4))
print("arr:\n", arr)

# 示例1: 数组与标量
print("arr + 5:\n", arr + 5) # 标量 5 被广播到 (3, 4)

# 示例2: 二维数组与一维数组
row_mean = arr.mean(axis=1) # shape (3,)
print("Row mean shape:", row_mean.shape)
# 需要将其 reshape 为 (3, 1) 才能按行广播减去均值
row_mean_reshaped = row_mean.reshape((3, 1))
print("Row mean reshaped:\n", row_mean_reshaped)
print("arr - row_mean_reshaped:\n", arr - row_mean_reshaped) # (3, 4) - (3, 1) ->
(3, 4)

# 示例3: 直接使用 shape (4,) 的一维数组进行列广播
col_vector = np.array([10, 20, 30, 40]) # shape (4,)
print("arr + col_vector:\n", arr + col_vector) # (3, 4) + (4,) -> (3, 4)
```

广播使得向量化代码更加简洁高效。

3.6 线性代数 (np.linalg)

NumPy 包含一个 linalg 模块,提供了许多线性代数运算功能。

- np.dot(a, b) 或 a @ b: 矩阵乘法 (点积)。
- np.linalg.inv(a): 计算矩阵的逆。
- np.linalg.det(a): 计算矩阵行列式。
- np.linalg.eig(a): 计算特征值和特征向量。
- np.linalg.svd(a): 计算奇异值分解 (SVD)。

```
x = np.array([[1., 2.], [3., 4.]])
y = np.array([[5., 6.], [7., 8.]])

print("Dot product (x @ y):\n", x @ y)
# 或者 np.dot(x, y)
```

虽然本教程重点不是线性代数,但了解 NumPy 具备这些能力很有用,因为许多机器学习算法底层依赖于线性代数运算。

3.7 保存与加载数组

将 NumPy 数组持久化到磁盘,以便后续使用或与其他程序共享,是非常常见的需求。

3.7.1 .npy 格式 (单个数组)

这是存储单个 NumPy 数组的标准二进制格式。

- 保存: np.save(filename, arr)
- 加载: arr = np.load(filename)

```
arr = np.arange(10)
np.save('my_array.npy', arr) # 保存到 my_array.npy

loaded_arr = np.load('my_array.npy')
print("Loaded array:", loaded_arr)
```

优点: 高效、保持数据类型和形状信息。 缺点: 只能存储单个数组。

应用关联: 你的脚本中广泛使用 .npy 格式来存储特征矩阵 (X_flat_features.npy)、目标向量 (Y_flat_target.npy)、各种索引 (indices_tp.npy, indices_fn.npy 等) 以及预测概率 (base_probs_train.npy, base_probs_test.npy) 和元特征 (X_meta.npy, y_meta.npy)。

3.7.2 .npz 格式 (多个数组)

可以将多个数组保存到一个未压缩的 .npz 文件中。

- 保存: np.savez(filename, name1=arr1, name2=arr2, ...)
- 加载: 加载 .npz 文件会得到一个类似字典的对象,可以通过保存时使用的名称来访问数组。

```
arr1 = np.arange(5)
arr2 = np.random.randn(2, 3)
np.savez('arrays.npz', a=arr1, b=arr2)

data = np.load('arrays.npz')
print("Array 'a':", data['a'])
print("Array 'b':", data['b'])
data.close() # 建议关闭文件
```

np.savez_compressed()可以将多个数组保存到压缩的.npz 文件,节省磁盘空间。

3.7.3 内存映射 (Memory Mapping) - 处理大型数组

当数组文件非常大,无法完全加载到内存时,内存映射是一个非常有用的技术。它允许你像操作内存中的数组一样操作磁盘上的文件,而无需一次性将整个文件读入 RAM。操作系统会在需要时按需加载文件的部分内容。

- **加载时启用内存映射:** 在 np.load() 中设置 mmap_mode 参数。
 - o mmap_mode='r': 只读模式。允许读取,不允许修改。
 - o mmap_mode='r+': 读写模式。允许读取和修改,修改会直接写入磁盘文件。
 - o mmap mode='w+': 创建或覆盖模式。用于创建新的内存映射文件。
 - o mmap_mode='c': 写时复制 (copy-on-write) 模式。修改不会写入原始文件,而是写入内存副本。

```
# 假设 large_array.npy 是一个非常大的文件
# 创建一个示例大文件(如果不存在)
# if not os.path.exists('large_array.npy'):
   large_data = np.arange(10**8).reshape(10**4, 10**4) # 创建一个大数组
    np.save('large_array.npy', large_data)
     del large_data # 释放内存
print("Loading large array with memory mapping (read-only)...")
try:
   # 使用内存映射加载 (只读)
   mmap_arr = np.load('large_array.npy', mmap_mode='r')
   print("Array loaded (mmap). Shape:", mmap_arr.shape)
   print("Accessing an element:", mmap_arr[100, 100]) # 可以像普通数组一样访问
   # 进行切片操作(仍然是内存映射视图)
   subset = mmap_arr[1000:1010, :5]
   print("Subset shape:", subset.shape)
   # 将子集加载到内存进行计算(如果需要)
   subset_in_memory = np.array(subset)
   print("Subset loaded into memory.")
   #!!!重要:使用完毕后,建议显式关闭内存映射文件句柄!!!
   # 虽然 Python 的垃圾回收通常会处理, 但在处理大文件或长时间运行时显式关闭更安全
   if hasattr(mmap_arr, '_mmap'):
       print("Closing memory map...")
       mmap_arr._mmap.close()
       print("Memory map closed.")
except FileNotFoundError:
   print("Error: large array.npy not found. Please create it first.")
except Exception as e:
   print(f"An error occurred: {e}")
```

应用关联: 你的多个脚本(1_generate_base_predictions.py, 4_generate_meta_features.py, 6_finetune_base_model.py) 在加载 X_flat_features.npy 和 Y_flat_target.npy 时都使用了mmap_mode='r'。这是处理可能非常大的特征和目标数据的关键策略。脚本中也包含了在使用完内存映射对象后关闭它的代码(例如 if hasattr(X_flat, '_mmap'): X_flat._mmap.close()), 这是一个很好的实践。分块处理(Chunking)通常与内存映射结合使用,先通过内存映射访问整个大数组,然后逐块将需要处理的数据实际加载到内存中进行计算(如模型预测)。

掌握 NumPy 是进行后续机器学习任务的基础。你需要熟练地创建、索引、切片和操作 NumPy 数组,并理解其在处理大型数据时的优势和技巧(如内存映射)。

第四章:数据分析利器: Pandas

如果说 NumPy 是处理数值数组的基石,那么 Pandas (Python Data Analysis Library) 就是在 NumPy 基础上构建的、用于数据清洗、处理、分析和可视化的强大武器。Pandas 引入了两种核心数据结构: Series 和 DataFrame,它们使得处理结构化数据(如表格数据、时间序列)变得极其方便。

在机器学习项目中, Pandas 通常用于:

- 加载和探查原始数据 (CSV, Excel, SQL 等)。
- 数据清洗:处理缺失值、异常值、重复数据。
- 数据转换:修改数据类型、创建新特征。
- 数据探索: 计算统计量、分组聚合、可视化。
- 为 Scikit-learn 等库准备输入数据。

虽然你的脚本示例主要使用预处理好的 NumPy 数组 (.npy),但在实际项目中,数据往往是以更原始的格式存在的,这时 Pandas 就派上了大用场。理解 Pandas 的基本操作对于完整的数据科学流程至关重要。

4.1 核心数据结构

4.1.1 Series

Series 是一维带标签的数组,可以存储任何数据类型(整数、浮点数、字符串、Python 对象等)。它有点像一个带有标签(索引)的 NumPy 数组,或者说是一个有序的字典。

```
import pandas as pd
import numpy as np
# 从列表创建 Series (默认整数索引)
s1 = pd.Series([10, 20, 30, 40])
print(s1)
# 输出:
# 0
      10
# 1 20
# 2
     30
# 3 40
# dtype: int64
# 查看值和索引
print("Values:", s1.values) # 输出: Values: [10 20 30 40] (NumPy 数组)
print("Index:", s1.index) # 输出: Index: RangeIndex(start=0, stop=4, step=1)
# 创建带自定义标签索引的 Series
s2 = pd.Series([95.5, 88.0, 73.5], index=['Math', 'English', 'History'])
print(s2)
# 输出:
# Math
          95.5
# English 88.0
# History
          73.5
# dtype: float64
print("Index:", s2.index) # 输出: Index: Index(['Math', 'English', 'History'],
dtype='object')
# 通过标签索引访问
print("Math score:", s2['Math']) # 输出: Math score: 95.5
print("Scores for Math and History:\n", s2[['Math', 'History']])
# 输出:
# Math
           95.5
```

```
# History 73.5
# dtype: float64
# 也可以通过位置索引 (类似 NumPy)
                             # 输出: First element: 95.5
print("First element:", s2[0])
print("Elements 0 and 2:\n", s2[[0, 2]])
# Series 也支持向量化操作和布尔索引
print("Scores > 80:\n", s2[s2 > 80])
print("Square root of scores:\n", np.sqrt(s2))
# 可以像字典一样操作
print("'English' in s2:", 'English' in s2) # 输出: True
print("'Physics' in s2:", 'Physics' in s2) # 输出: False
# 从字典创建 Series
grades = {'Alice': 85, 'Bob': 92, 'Charlie': 78}
s3 = pd.Series(grades)
print(s3)
# 输出:
# Alice
# Bob
# Charlie 78
# dtype: int64
# 可以指定索引顺序, 缺失的值为 NaN
students = ['Alice', 'Bob', 'David']
s4 = pd.Series(grades, index=students)
print(s4)
# 输出:
# Alice 85.0
# Bob
        92.0
# David NaN <-- 注意 NaN (Not a Number)
# dtype: float64
# 检查缺失值
print("Is NaN:\n", pd.isnull(s4)) # 或者 s4.isnull()
# 输出:
# Alice False
# Bob
        False
# David True
# dtype: bool
```

4.1.2 DataFrame

DataFrame 是 Pandas 的核心,它是一个二维的、大小可变的、异构的表格型数据结构,带有标签的轴(行和列)。你可以把它想象成一个 Excel 电子表格、一个 SQL 表,或者一个 Series 对象的字典。

- **行索引** (Index): 每一行有一个唯一的标识符。
- **列索引** (Columns): 每一列有一个名称。
- 数据: 存储在 DataFrame 中的值,每列可以有不同的数据类型。

```
# 从字典创建 DataFrame (字典的 key 是列名, value 是列表或 Series)
data = {'state': ['Ohio', 'Ohio', 'Nevada', 'Nevada'],
       'year': [2000, 2001, 2002, 2001, 2002, 2003],
       'pop': [1.5, 1.7, 3.6, 2.4, 2.9, 3.2]}
df = pd.DataFrame(data)
print(df)
# 输出:
    state year pop
# 0
     Ohio 2000 1.5
# 1 Ohio 2001 1.7
# 2
     Ohio 2002 3.6
# 3 Nevada 2001 2.4
# 4 Nevada 2002 2.9
# 5 Nevada 2003 3.2
# 查看 DataFrame 信息
print("\nHead:\n", df.head()) # 默认显示前 5 行
print("\nTail:\n", df.tail(3)) # 显示后 3 行
print("\nInfo:")
df.info() # 打印 DataFrame 的摘要信息 (索引、列、非空值、类型、内存)
# 输出:
# <class 'pandas.core.frame.DataFrame'>
# RangeIndex: 6 entries, 0 to 5
# Data columns (total 3 columns):
# # Column Non-Null Count Dtype
# --- ----- -----
# 0 state 6 non-null
                          object <-- object 通常指字符串
# 1 year 6 non-null
                          int64
     pop 6 non-null
                          float64
# dtypes: float64(1), int64(1), object(1)
# memory usage: 168.0+ bytes
print("\nDescribe:\n", df.describe()) # 生成数值列的描述性统计信息
# 输出:
#
             year
                       pop
        6.000000 6.000000
# count
# mean 2001.500000 2.550000
        1.048809 0.836062
# std
# min
       2000.000000 1.500000
# 25%
        2001.000000 1.875000
# 50%
       2001.500000 2.650000
# 75%
       2002.000000 3.125000
        2003.000000 3.600000
# max
print("\nShape:", df.shape)
                             # 输出: Shape: (6, 3)
print("\nColumns:", df.columns) # 输出: Columns: Index(['state', 'year', 'pop'],
dtype='object')
print("\nIndex:", df.index)
                             # 输出: Index: RangeIndex(start=0, stop=6,
step=1)
print("\nData types:\n", df.dtypes) # 输出 Series 各列的数据类型
# 指定列的顺序
df2 = pd.DataFrame(data, columns=['year', 'state', 'pop'])
```

```
print("\nDataFrame with specified columns:\n", df2)

# 指定索引

df3 = pd.DataFrame(data, index=['one', 'two', 'three', 'four', 'five', 'six'])

print("\nDataFrame with specified index:\n", df3)

# 从 NumPy 数组创建 DataFrame

arr = np.random.randn(5, 3)

df_from_np = pd.DataFrame(arr, columns=['A', 'B', 'C'],
    index=pd.date_range('20230101', periods=5))

print("\nDataFrame from NumPy array:\n", df_from_np)
```

4.2 数据输入/输出 (I/O)

Pandas 提供了丰富的函数来读取和写入各种格式的数据。

• 读取数据:

- o pd.read_csv(filepath_or_buffer, sep=',', header='infer', index_col=None, ...): 读取 CSV 文件。sep 指定分隔符,header 指定哪一行作为列名,index_col 指定哪一列作为行索引。
- o pd.read_excel(io, sheet_name=0, header=0, ...): 读取 Excel 文件。
- o pd.read_sql(sql, con): 从 SQL 数据库读取数据 (需要数据库连接 con)。
- pd.read_json(path_or_buf, ...): 读取 JSON 文件。
- o pd.read_hdf(path_or_buf, key=None, ...): 读取 HDF5 文件 (常用于存储大型 DataFrame) 。
- o pd.read_parquet(path, ...): 读取 Parquet 文件(列式存储格式, 高效)。

• 写入数据:

- df.to_csv(path_or_buf, sep=',', index=True, header=True, ...): 写入 CSV 文件。index=False 可以禁止写入行索引。
- o df.to_excel(excel_writer, sheet_name='Sheet1', index=True, ...):写入 Excel 文件。
- o df.to_sql(name, con, if_exists='fail', index=True, ...): 写入 SQL 数据库表。 if_exists 参数 ('replace', 'append', 'fail') 很重要。
- df.to_json(path_or_buf, orient='records', ...):写入 JSON 文件。
- df.to_hdf(path_or_buf, key, mode='a', ...):写入 HDF5 文件。
- o df.to_parquet(path, engine='auto', ...): 写入 Parquet 文件。

4.3 数据查看与选择

高效地查看和选择 DataFrame 中的数据是 Pandas 的核心功能。

4.3.1 查看数据

前面已经介绍过 head(), tail(), info(), describe(), shape, columns, index, dtypes 等常用方法。

4.3.2 选择列

```
# 选择单列 (返回 Series)
print("State column:\n", df['state'])
# 等价于 df.state (但不适用于包含空格或与方法名冲突的列名)

# 选择多列 (返回 DataFrame)
print("\nYear and Pop columns:\n", df[['year', 'pop']])
```

4.3.3 选择行

选择行的主要方式有两种:基于标签 (loc) 和基于整数位置 (iloc)。

• loc (Label-based selection): 使用行标签(索引名)和列标签(列名)进行选择。包含结束标签。

```
print("\nUsing loc:")
print("Row 'two':\n", df3.loc['two']) # 选择单行 (返回 Series)
print("\nRows 'two' to 'four':\n", df3.loc['two':'four']) # 选择多行切片 (包含 'four')
print("\nRows 'one', 'three' and columns 'state', 'pop':\n", df3.loc[['one', 'three'], ['state', 'pop']])
print("\nRow 'four', column 'year':", df3.loc['four', 'year'])

# 使用布尔 Series 进行行选择 (非常常用)
print("\nRows where year > 2001:\n", df[df['year'] > 2001]) # 直接用布尔 Series
# 或者使用 loc
print("\nRows where year > 2001 (using loc):\n", df.loc[df['year'] > 2001])
```

• iloc (Integer position-based selection): 使用整数位置进行选择(类似 NumPy 数组索引)。不包含结束位置。

```
print("\nUsing iloc:")
print("Row at index 0:\n", df.iloc[0]) # 选择第0行 (返回 Series)
print("\nRows at index 1 to 3 (exclusive):\n", df.iloc[1:3]) # 选择第1行和第2
行
print("\nRows at index 0, 2, 4 and columns 0, 2:\n", df.iloc[[0, 2, 4], [0, 2]])
print("\nElement at row 1, column 1:", df.iloc[1, 1])
```

选择的最佳实践:

• 优先使用 loc 和 iloc 进行显式索引,避免直接使用 df[...]进行行选择(虽然有时可行,但容易混淆且可能引发 SettingWithCopyWarning)。

• 当需要基于条件选择行时,布尔索引 (df[boolean_condition]) 是最常用的方式。

4.4 数据清洗 (概览)

真实世界的数据往往是"脏"的, Pandas 提供了处理常见问题的工具。

4.4.1 处理缺失值 (NaN)

- isnull() / notnull(): 检测缺失值,返回布尔型的 DataFrame 或 Series。
- dropna(axis=0, how='any', thresh=None, subset=None, ...): 删除包含缺失值的行(axis=0) 或列(axis=1)。
 - how='any': 只要有 NaN 就删除。
 - how='all': 只有所有值都是 NaN 才删除。
 - o thresh: 保留至少有 N 个非 NaN 值的行/列。
 - 。 subset: 指定在哪些列中检查 NaN。
- fillna(value=None, method=None, axis=None, limit=None, ...): 填充缺失值。
 - o value: 用于填充的标量值或字典/Series (指定不同列的填充值)。
 - o method='ffill':使用前一个有效值填充 (forward fill)。
 - o method='bfill':使用后一个有效值填充(backward fill)。

```
s4 = pd.Series({'Alice': 85, 'Bob': 92, 'David': np.nan})
print("Original Series with NaN:\n", s4)
print("Is NaN:\n", s4.isnull())
# 删除 NaN
print("Drop NaN:\n", s4.dropna())
# 填充 NaN
print("Fill NaN with 0:\n", s4.fillna(0))
print("Fill NaN with mean:\n", s4.fillna(s4.mean()))
# DataFrame 示例
df_nan = pd.DataFrame([[1., 6.5, 3.], [1., np.nan, np.nan],
                       [np.nan, np.nan, np.nan], [np.nan, 6.5, 3.]])
print("\nDataFrame with NaN:\n", df_nan)
print("Drop rows with any NaN:\n", df_nan.dropna())
print("Drop rows where all are NaN:\n", df nan.dropna(how='all'))
print("Fill all NaN with -1:\n", df_nan.fillna(-1))
print("Fill NaN in column 1 with mean of column 1:\n", df_nan.fillna({1:
df nan[1].mean()}))
```

4.4.2 处理重复数据

• duplicated(subset=None, keep='first'): 检测重复行,返回布尔型 Series。

- keep='first': 标记除第一次出现之外的重复项为 True。
- o keep='last': 标记除最后一次出现之外的重复项为 True。
- keep=False: 标记所有重复项为 True。
- drop_duplicates(subset=None, keep='first', inplace=False): 删除重复行。参数与duplicated 类似。inplace=True 会直接修改原 DataFrame。

```
df_dup = pd.DataFrame({'col1': ['A', 'B', 'B', 'C', 'A'], 'col2': [1, 2, 2, 3,
1]})
print("\nDataFrame with duplicates:\n", df_dup)
print("Is duplicated:\n", df_dup.duplicated())
print("Drop duplicates:\n", df_dup.drop_duplicates())
print("Drop duplicates keeping last:\n", df_dup.drop_duplicates(keep='last'))
print("Drop duplicates based on col1 only:\n", df_dup.drop_duplicates(subset= ['col1']))
```

4.5 基本数据操作

4.5.1 添加/删除列

```
# 添加列

df['debt'] = np.random.randn(6) * 10 # 添加新列 'debt'

df['eastern'] = (df['state'] == 'Ohio') # 添加基于条件的布尔列

print("\nDataFrame with new columns:\n", df)

# 删除列

del df['eastern'] # 使用 del 关键字

pop_col = df.pop('pop') # 使用 pop 方法 (删除并返回该列)

print("\nDataFrame after deleting columns:\n", df)

print("Popped 'pop' column:\n", pop_col)
```

4.5.2 应用函数

- apply(func, axis=0): 将函数 func 应用于 DataFrame 的行 (axis=1) 或列 (axis=0)。
- map(arg): 用于 Series,将函数或字典/Series 应用于 Series 的每个元素。

```
df = pd.DataFrame(np.random.randn(4, 3) * 5, columns=['A', 'B', 'C'])
print("\nOriginal DataFrame for apply/map:\n", df)

# 定义一个函数, 计算范围 (max - min)
f = lambda x: x.max() - x.min()

# 应用于列 (默认 axis=0)
print("Range of each column:\n", df.apply(f))
# 应用于行
print("Range of each row:\n", df.apply(f, axis=1))
```

```
# 使用 map 对 Series 元素进行格式化
format_func = lambda x: f'{x:.2f}' # 保留两位小数
print("Formatted column A:\n", df['A'].map(format_func))
```

4.6 分组与聚合 (groupby) (简介)

groupby 是 Pandas 中非常强大的功能,用于实现数据的 Split-Apply-Combine 策略:

- 1. Split: 根据某些键 (列名、函数等) 将数据分割成组。
- 2. **Apply:** 对每个组独立地应用一个函数(如聚合函数 sum, mean, count, 或转换函数 transform, 或自定义函数 apply)。
- 3. Combine: 将应用函数后的结果合并成一个新的数据结构。

groupby 对于数据探索和特征工程非常有用。

4.7 合并与连接 (merge, concat) (简介)

Pandas 提供了多种方式来合并来自不同 DataFrame 的数据。

- pd.concat(objs, axis=0, join='outer', ...): 沿着一个轴 (默认 axis=0, 即按行堆叠) 将多个 DataFrame 或 Series 拼接在一起。
 - o join='outer': 保留所有标签 (并集) , 缺失处填充 NaN。
 - o join='inner': 只保留共有的标签(交集)。
- pd.merge(left, right, how='inner', on=None, left_on=None, right_on=None, ...): 类似于 SQL 的 JOIN 操作,根据一个或多个键将两个 DataFrame 的行连接起来。
 - on: 用于连接的列名(左右 DataFrame 中都存在)。
 - left_on, right_on: 分别指定左右 DataFrame 中用于连接的列名。
 - how: 连接方式 ('inner', 'outer', 'left', 'right')。

这些操作在需要整合来自不同数据源的信息时非常关键。

4.8 Pandas 与 Scikit-learn

虽然 Scikit-learn 的函数通常接受 NumPy 数组作为输入,但使用 Pandas DataFrame 作为输入也很常见,尤其是在数据准备阶段。

- 优点: DataFrame 带有列名,更易读;可以方便地进行基于标签的选择和操作。
- 注意事项:
 - 。 **数据类型:** 确保传递给 Scikit-learn 的 DataFrame 只包含数值类型。分类特征需要预先编码(如 One-Hot Encoding)。
 - 特征顺序: Scikit-learn 模型训练后会记住特征的顺序。在预测时,输入的 DataFrame 必须具有与训练时完全相同的列(特征)和顺序。
 - 。 **仅传递特征**: 在调用 fit(X, y) 或 predict(X) 时,确保 X 只包含特征列,不包含目标变量 y。

你的脚本示例中,数据在早期就被转换并保存为 NumPy 数组 (*npy),这可能是为了在后续步骤中简化数据加载和处理,并利用 NumPy 的内存映射功能。但在实际项目中,你会经常看到在脚本的早期阶段使用 Pandas 进行大量的加载、清洗和特征工程工作,最后才将准备好的数据转换为 NumPy 数组输入到 Scikit-learn 或 XGBoost 模型中。

```
# 假设 df_features 是包含特征的 DataFrame, series_target 是包含标签的 Series # df_features 可能包含需要处理的非数值列
# ... (使用 Pandas 进行清洗、编码、特征工程) ...
# 转换为 NumPy 数组以输入 Scikit-learn
# X = df_prepared_features.values
# y = series_target.values
# model.fit(X, y)
```

掌握 Pandas 对于处理和准备机器学习所需的数据至关重要。你需要熟悉其核心数据结构、数据选择方法以及常用的数据清洗和转换操作。

第五章: Scikit-learn 核心 API 与通用功能

Scikit-learn (sklearn) 是 Python 中用于构建机器学习模型的基石库。它提供了一个统一且简洁的接口来访问各种预处理技术、机器学习算法和评估工具。本章将深入探讨 Scikit-learn 的设计哲学、核心 API 以及一些最常用的功能模块。理解这些是高效使用 Scikit-learn 及其生态(包括 XGBoost 的 Scikit-learn 接口)的关键。

5.1 Scikit-learn 生态系统概览

Scikit-learn 的功能覆盖了机器学习工作流的大部分环节:

- **预处理 (sklearn.preprocessing):** 特征缩放、归一化、编码(分类特征、标签)、缺失值填充等。
- 降维 (sklearn.decomposition, sklearn.manifold): PCA, t-SNE 等。
- 模型选择与评估 (sklearn.model_selection, sklearn.metrics): 数据划分 (train_test_split)、交叉验证、超参数调优(网格搜索、随机搜索)、各种评估指标(准确率、召回率、AUC等)。
- 监督学习算法:
 - 分类(sklearn.linear_model, sklearn.svm, sklearn.tree, sklearn.ensemble, sklearn.neighbors, etc.): 逻辑回归, 支持向量机 (SVM), 决策树, 随机森林, 梯度提升树 (GBDT), K

近邻 (KNN) 等。

○ **回归 (**sklearn.linear_model, sklearn.svm, sklearn.tree, sklearn.ensemble, **etc.**): 线性回归, 岭回归 (Ridge), Lasso, SVR, 决策树回归, 随机森林回归等。

- 无监督学习算法 (sklearn.cluster, sklearn.decomposition): K-Means, DBSCAN, PCA 等。
- 模型持久化 (joblib): 保存和加载训练好的模型。

5.2 核心设计原则: 一致性的 Estimator API

Scikit-learn 最强大的地方之一在于其**一致性的 API 设计**。库中的大多数对象都遵循 "Estimator" (估计器) 接口规范。Estimator 是任何可以从数据中学习参数的对象,无论是用于分类、回归、聚类还是数据转换。

这种一致性意味着一旦你理解了如何使用一个 Estimator,你就能很容易地学会使用其他 Estimator,大大降低了学习成本。

关键方法 (Methods):

所有 Estimator 都共享一些核心方法:

- 1. fit(X, y=None): 这是所有 Estimator 的核心学习方法。
 - 作用:根据输入数据 X (特征矩阵)和可选的目标变量 y (标签向量,监督学习需要)来训练模型,调整内部参数。
 - 输入:
 - X: 通常是一个形状为 (n_samples, n_features) 的类数组对象 (NumPy 数组、Pandas DataFrame、稀疏矩阵) 。
 - y: 对于监督学习(分类、回归),是一个形状为 (n_samples,) 的类数组对象 (NumPy 数组、Pandas Series)。对于无监督学习,通常不需要 y。
 - **返回:** self,即 Estimator对象本身。这允许链式调用,例如 model.fit(X_train, y_train).predict(X_test)。
 - o **内部状态:** fit 方法会修改 Estimator 的内部状态,学习到的参数通常存储在以**下划线结尾**的属性中(例如,线性回归的系数存储在 model.coef)。
- 2. predict(X): 用于监督学习模型(分类器、回归器)进行预测。
 - ∘ 作用: 使用 fit 过的模型,对新的输入数据 X 生成预测结果。
 - **输入:** X: 形状为 (n_new_samples, n_features) 的类数组对象,特征必须与训练时的 X 对应。
 - **返回**:一个包含预测结果的 NumPy 数组,形状为 (n_new_samples,)。分类器返回预测的类别标签,回归器返回预测的连续值。
- 3. predict_proba(X):用于分类器,预测每个类别的概率。
 - 作用: 返回每个样本属于每个类别的概率估计。这对于理解模型的不确定性、调整分类阈值或用作 Stacking 的元特征非常有用。
 - **输入:** X: 形状为 (n new samples, n features) 的类数组对象。
 - **返回:** 一个 NumPy 数组,形状为 (n_new_samples, n_classes)。每一行代表一个样本,每一列代表一个类别(按 estimator.classes_ 属性的顺序排列),对应的值是该样本属于该类别的概率。每一行的概率之和应为 1。
 - **应用关联:** 你的脚本中大量使用了 predict proba:
 - 1_generate_base_predictions.py: 计算基模型在训练集和测试集上的概率 (base_probs_train.npy, base_probs_test.npy)。

■ 2_train_fn_expert.py, 3_train_fp_expert.py: 评估专家模型时可能用到(虽然最终评估用了阈值)。

- 4_generate_meta_features.py: FP 专家模型使用 predict_proba 生成元特征的一部分。
- 5_train_evaluate_meta_learner.py: 元学习器使用 predict_proba 进行评估。
- 6_finetune_base_model.py:评估微调模型时使用 predict_proba。
- evaluate_fp_expert_on_tp.py: 使用 FP 专家的 predict_proba 来分析其在 TP 样本上的行为。
- 4. transform(X): 用于数据**预处理或降维**的 Estimator (也称为 Transformer)。
 - o 作用: 使用 fit 过的 Transformer 对新的输入数据 X 进行转换 (例如,缩放、编码、降维)。
 - **输入:** X: 形状为 (n samples, n features) 的类数组对象。
 - 。 **返回:** 转换后的数据,通常是 NumPy 数组,形状可能与输入不同(例如,降维或 One-Hot 编码)。
- 5. fit_transform(X, y=None): 结合了 fit 和 transform。
 - 作用: 在数据 X 上调用 fit, 然后对同一个数据 X 调用 transform。这通常比分开调用 fit 和 transform 更高效。
 - · 常用场景: 在训练集上应用预处理步骤。
 - **注意:** fit_transform 只能应用于训练数据。对于测试数据或新的未知数据,必须使用已经 fit 过的 Transformer 的 transform 方法,以保证使用相同的转换规则(例如,使用训练集的均值和标准差来缩放测试集)。
- 6. score(X, y):用于评估模型在给定数据上的默认性能指标。
 - 。 作用: 对输入数据 X 进行预测,然后将预测结果与真实标签 y 进行比较,计算一个预定义的评分。
 - 。 **返回:** 一个浮点数, 表示得分。
 - 默认指标:
 - 分类器: 通常是平均准确率 (mean accuracy)。
 - 回归器: 通常是决定系数 (Coefficient of Determination, R² score)。
 - 局限性: 默认指标不一定适用于所有场景(例如,不平衡分类问题只看准确率可能具有误导性)。后续章节会详细介绍更多评估指标。

超参数 (Hyperparameters):

Estimator 的行为可以通过在实例化时传递参数来控制,这些参数称为超参数。它们不是从数据中学习的,而是由用户设置的。

- 设置: model = MyEstimator(param1=value1, param2=value2)
- 查看: model.get_params()返回一个包含所有超参数及其值的字典。
- 修改: model.set params(param1=new value)(通常在模型选择或调优中使用)。

例如, XGBoost 的 n estimators, learning rate, max depth 等都是超参数。

5.3 数据划分 (sklearn.model_selection.train_test_split)

将数据集划分为训练集和测试集是评估模型泛化能力的标准做法。

为什么必须划分?

如果使用全部数据来训练模型,然后在**同一份**数据上评估性能,得到的结果会过于乐观,无法反映模型在处理 从未见过的新数据时的真实表现。模型可能会"记住"训练数据(过拟合),但在新数据上效果很差。测试集模 拟了模型在实际部署后遇到的新数据。

train_test_split 函数:

```
from sklearn.model selection import train test split
import numpy as np
X = np.arange(20).reshape((10, 2)) # 示例特征矩阵 (10个样本, 2个特征)
                              # 示例标签向量
y = np.arange(10)
# 基本用法
X_train, X_test, y_train, y_test = train_test_split(
                          # 要划分的数组 (特征和标签)
   random_state=42, # 随机种子, 保证每次划分结果一致 shuffle=True # 是不在因为
                         # 测试集比例(也可以是整数表示绝对数量)
                          # 是否在划分前打乱数据 (默认 True)
)
print("X_train shape:", X_train.shape) # 输出: X_train shape: (7, 2)
print("X_test shape:", X_test.shape) # 输出: X_test shape: (3, 2)
print("y_train shape:", y_train.shape) # 输出: y_train shape: (7,)
print("y_test shape:", y_test.shape) # 输出: y_test shape: (3,)
```

关键参数详解:

- *arrays: 需要划分的一个或多个数组(例如 X 和 y)。它们会沿着第一个轴(行)被同步划分。
- test size: 指定测试集的大小。
 - 如果是浮点数 (0.0 到 1.0 之间), 表示测试集的比例。
 - 。 如果是整数,表示测试集的绝对样本数。
 - 如果未指定 train_size,则 train_size 会自动设为 1.0 test_size。
- train_size: 指定训练集的大小。用法同 test_size。通常只需指定 test_size 或 train_size 中的一个。
- random_state: 控制划分前的随机数生成器。设置一个固定的整数(如 0, 42)可以确保每次运行代码时得到的划分结果完全相同,这对于结果复现非常重要。如果不设置,每次运行结果会不同。
 - **应用关联:** 你的脚本中 RANDOM_STATE = 42 就是用于此目的,确保专家模型训练时的划分 (2 ...py, 3 ...py) 以及后续可能需要划分的地方都使用相同的随机状态。
- shuffle: 是否在划分前打乱样本顺序。默认为 True。对于非时序数据通常需要打乱。
- stratify: 用于分类问题,非常重要!
 - 作用:保证划分后的训练集和测试集中,各个类别标签的比例与原始数据集中大致相同。
 - 输入: 通常传入标签数组 √。
 - **为何重要**: 在类别不平衡(某些类别样本数量远多于其他类别)的情况下,简单的随机划分可能导致测试集(甚至训练集)中某个类别的样本极少或没有,从而无法准确评估模型在这些少数类别上的性能。stratify=y可以避免这种情况。
 - o **应用关联**: 脚本 2_train_fn_expert.py 和 3_train_fp_expert.py 在划分 FN/TN 和 FP/TN 数据时,都正确地使用了 stratify=labels,因为 FN/TN 或 FP/TN 的比例可能是不平衡的,需要确保训练集和验证集都有代表性的样本。

验证集 (Validation Set):

train_test_split 通常用于划分训练集和**测试集**。但在模型选择和超参数调优过程中,我们还需要一个**验证集**。常见的做法有两种:

1. **再划分一次:** 先将数据划分为训练集和测试集。然后,再将训练集进一步划分为新的(较小的)训练集和 验证集。

```
# 划分训练集和测试集 (e.g., 80% train, 20% test)
X_train_full, X_test, y_train_full, y_test = train_test_split(X, y, test_size=0.2, random_state=42, stratify=y)
# 再从训练集中划分出验证集 (e.g., 原始数据的 15% 作为验证, 65% 作为最终训练)
# test_size for the second split = val_size / (1 - test_size_initial) = 0.15
/ (1 - 0.2) = 0.1875
X_train, X_val, y_train, y_val = train_test_split(X_train_full, y_train_full, test_size=0.1875, random_state=42, stratify=y_train_full)
print("Final Train shape:", X_train.shape)
print("Validation shape:", X_val.shape)
print("Test shape:", X_test.shape)
```

- o **应用关联:** 脚本 2_train_fn_expert.py 和 3_train_fp_expert.py 就是采用了这种方法,它们使用 VALIDATION_SIZE = 0.15 从 FN/TN 或 FP/TN 数据中划分出验证集 (X_val, y_val),用于 XGBoost 的 eval_set 和 early_stopping_rounds。
- 2. **交叉验证** (Cross-Validation): 更稳健的方法,我们将在后续章节详细介绍。
- 5.4 数据预处理 (sklearn.preprocessing) (概览)

原始特征数据往往不能直接输入机器学习模型,需要进行预处理。Scikit-learn 的 preprocessing 模块提供了多种工具。

为何需要预处理?

- **算法要求**: 许多算法(如基于距离的 KNN、SVM,以及依赖梯度下降优化的线性模型、神经网络)对特征的尺度非常敏感。如果特征尺度差异很大,尺度较大的特征会主导模型的学习过程。
- 性能提升: 合适的预处理可以提高模型的收敛速度和最终性能。
- 数据兼容: 将非数值特征(如类别)转换为数值表示。

5.4.1 特征缩放 (Feature Scaling)

将特征数值调整到相似的范围。

- StandardScaler (标准化 / Z-score 归一化):
 - 公式: z = (x mean) / std_dev
 - 。 作用: 将特征转换为均值为 0, 标准差为 1 的分布。
 - 适用场景: 适用于数据近似高斯分布的情况,或不确定数据分布时常用的选择。对异常值相对敏感。
- MinMaxScaler (最小-最大规范化):
 - 公式: X_scaled = (X X_min) / (X_max X_min) (默认缩放到 [0, 1])
 - 作用: 将特征缩放到一个给定的范围(通常是[0,1]或[-1,1])。
 - 适用场景: 当需要将特征限定在特定边界内时,或者数据分布不呈高斯分布时。对异常值非常敏感。
- RobustScaler: 使用中位数和四分位数范围 (IQR) 进行缩放,对异常值更鲁棒。

重要原则: Fit on Training, Transform Train & Test

为了防止数据泄露(即防止测试集的信息影响到训练过程),特征缩放器 (Scaler) 应该**只在训练数据上调用** fit **或** fit_transform 来学习缩放参数(如均值、标准差、最大/最小值)。然后,使用**同一个**学习到的缩放器,分别对训练集和测试集(以及任何未来的新数据)调用 transform 方法。

```
from sklearn.preprocessing import StandardScaler, MinMaxScaler

# 假设 X_train, X_test 已经划分好
scaler_std = StandardScaler()
scaler_minmax = MinMaxScaler()

# 在训练集上拟合和转换
X_train_scaled_std = scaler_std.fit_transform(X_train)
X_train_scaled_minmax = scaler_minmax.fit_transform(X_train)

# 使用同一个 scaler 转换测试集
X_test_scaled_std = scaler_std.transform(X_test)
X_test_scaled_minmax = scaler_minmax.transform(X_test)

print("Original X_train sample:\n", X_train[0])
print("StandardScaled X_train sample:\n", X_train_scaled_std[0])
print("MinMaxScaled X_train sample:\n", X_train_scaled_minmax[0])
```

- **应用关联**: 你的脚本示例中似乎没有显式地进行特征缩放。这可能是因为:
 - 1. 数据在生成 .npy 文件前已经预处理过了。
 - 2. 使用的 XGBoost 模型对特征尺度相对不敏感(尤其是基于树的模型),虽然标准化有时也能带来 微小提升。 但在使用其他模型(如逻辑回归、SVM、神经网络)时,特征缩放通常是**必须**的步骤。

5.4.2 编码分类特征 (简介)

机器学习模型通常只能处理数值数据。对于类别特征(如 "性别"={'男', '女'}, "城市"={'北京', '上海', '广州'}), 需要将其转换为数值表示。

- OrdinalEncoder: 将类别特征转换为整数(例如,'小'->0,'中'->1,'大'->2)。适用于类别间存在**有序关 系**的情况。
- OneHotEncoder: 将每个类别转换为一个二进制向量(独热编码)。例如,'颜色'={'红','绿','蓝'},'红'->[1,0,0],'绿'->[0,1,0],'蓝'->[0,0,1]。
 - · **优点**: 避免了引入不存在的顺序关系。
 - · **缺点**: 如果类别数量很多,会导致特征维度急剧增加(维度爆炸)。
 - **常用参数:** handle_unknown='ignore' 可以在转换新数据时忽略训练集中未出现过的类别(否则会报错)。sparse=False 可以返回 NumPy 数组而不是稀疏矩阵。
- LabelEncoder: 通常**只用于编码目标变量 y**,将类别标签转换为 [0, n_classes-1] 范围内的整数。不推荐 直接用于编码输入特征 X(除非是树模型,且只有一个分类特征),因为它引入的序数关系可能误导模型。

```
from sklearn.preprocessing import OneHotEncoder, LabelEncoder
# --- One-Hot Encoding for Features ---
data_cat = pd.DataFrame({'color': ['Red', 'Green', 'Blue', 'Green'],
                         'size': ['S', 'M', 'L', 'S']})
print("\nCategorical DataFrame:\n", data_cat)
encoder_ohe = OneHotEncoder(sparse_output=False, handle_unknown='ignore') #
sparse=False in older versions
# Fit on training data (or all data if appropriate)
encoder ohe.fit(data cat)
# Transform
X encoded = encoder ohe.transform(data cat)
print("One-Hot Encoded Features:\n", X encoded)
print("Feature names:", encoder_ohe.get_feature_names_out())
# Output:
# One-Hot Encoded Features:
# [[1. 0. 0. 0. 1. 0.] <-- Red, S
# [0. 1. 0. 1. 0. 0.] <-- Green, M
# [0. 0. 1. 0. 0. 1.] <-- Blue, L
# [0. 1. 0. 0. 1. 0.]] <-- Green, S
# Feature names: ['color Blue' 'color Green' 'color Red' 'size L' 'size M'
'size S']
# --- Label Encoding for Target Variable ---
y_labels = np.array(['Cat', 'Dog', 'Cat', 'Fish', 'Dog'])
print("\nOriginal labels:", y_labels)
```

```
encoder_label = LabelEncoder()
y_encoded = encoder_label.fit_transform(y_labels)
print("Label Encoded target:", y_encoded) # Output: [0 1 0 2 1]
print("Encoded classes:", encoder_label.classes_) # Output: ['Cat' 'Dog' 'Fish']
# Inverse transform
print("Decoded labels:", encoder_label.inverse_transform(y_encoded))
```

• 应用关联: 你的脚本输入的 X_flat_features.npy 假设已经是纯数值特征,可能在之前的步骤中已经完 成了编码。XGBoost 的 Scikit-learn 接口的 use label encoder=False 参数 (在脚本 2, 3, 5 中使用) 建议关闭内置的标签编码,推荐用户在外部自行处理。

5.5 模型持久化 (joblib)

训练一个好的机器学习模型可能需要花费大量时间和计算资源。为了避免每次使用时都重新训练,我们需要将 训练好的模型保存到磁盘,并在需要时加载回来。joblib 是 Scikit-learn 推荐的用于持久化包含大型 NumPy 数组的对象(如训练好的模型)的库。

- 保存模型: joblib.dump(model, filename)
- 加载模型: loaded_model = joblib.load(filename)

```
from sklearn.linear_model import LogisticRegression
import joblib
import os
# 假设已经训练好了一个模型 model
# model = LogisticRegression().fit(X_train, y_train) # 示例训练
# --- 保存模型 ---
model filename = 'logistic regression model.joblib'
print(f"\nSaving model to {model_filename}...")
try:
   # 确保目录存在(如果需要)
   # model dir = os.path.dirname(model filename)
   # if model dir and not os.path.exists(model dir):
         os.makedirs(model dir)
   # joblib.dump(model, model filename) # 替换为实际 model 变量
   print("Model saved successfully (simulation).")
except Exception as e:
   print(f"Error saving model: {e}")
# --- 加载模型 ---
print(f"\nLoading model from {model_filename}...")
if os.path.exists(model filename): # 检查文件是否存在
       # loaded model = joblib.load(model filename)
       print("Model loaded successfully (simulation).")
       # 现在可以使用 loaded model 进行预测
       # y_pred = loaded_model.predict(X_test)
   except Exception as e:
```

```
print(f"Error loading model: {e}")
else:
   print(f"Error: Model file not found at {model_filename}.")
```

应用关联: 你的所有脚本都广泛应用了 joblib 来保存和加载模型:

- 1_generate_base_predictions.py:加载基模型(xgboost_default_full_model.joblib)。
- 2_train_fn_expert.py: 保存 FN 专家模型 (fn_expert_model.joblib)。
- 3_train_fp_expert.py: 保存 FP 专家模型 (fp_expert_model.joblib)。
- 4_generate_meta_features.py: 加载 FP 专家模型。
- 5_train_evaluate_meta_learner.py: 保存元学习器模型 (meta_learner_model.joblib)。
- debug_predict_test.py: 加载指定模型进行测试。
- 6_finetune_base_model.py: 加载初始基模型,保存微调后的模型 (xgboost_finetuned_model.joblib)。
- evaluate_fp_expert_on_tp.py: 加载 FP 专家模型。

注意事项:

- 版本兼容性: 加载模型时,环境中 Scikit-learn, NumPy, XGBoost 等相关库的版本最好与保存模型时使用的版本一致或兼容,否则可能出现错误或行为不一致。这是使用虚拟环境隔离项目依赖的重要原因之一。
- 安全性: joblib (以及 Python 的 pickle) 加载文件时可能执行任意代码,因此永远不要加载来自不可信来源的模型文件。

掌握 Scikit-learn 的核心 API 和通用功能是进行高效机器学习开发的基石。你需要理解 Estimator 的 fit/predict/transform 模式,知道如何正确划分数据,了解基本的预处理技术,并能够保存和加载你的模型。这将为你学习后续更具体的模型和技术打下坚实的基础。

```
```markdown
机器学习实战: 从 Scikit-learn 到 XGBoost (卷二)
...(接上文)

第六章: 监督学习: 分类模型详解
在监督学习中, 分类任务的目标是预测样本属于哪个预定义的类别。我们在第一章已经介绍了其基本概
```

在监督学习中,分类任务的目标是预测样本属于哪个预定义的类别。我们在第一章已经介绍了其基本概念。本章将深入探讨 Scikit-learn 中几种常用且重要的分类算法。了解这些不同模型的原理、优缺点、关键参数以及适用场景,对于为你的特定问题选择合适的工具至关重要。

#### 我们将重点介绍以下几种模型:

- 1. 逻辑回归 (Logistic Regression)
- 2. K-近邻 (K-Nearest Neighbors, KNN)
- 3. 支持向量机 (Support Vector Machines, SVM / SVC)
- 4. 决策树 (Decision Trees)
- 5. 随机森林 (Random Forests)

### 6.1 逻辑回归 (`sklearn.linear\_model.LogisticRegression`)

尽管名字带有"回归",但逻辑回归是用于解决\*\*分类问题\*\*的一种经典的、广泛使用的线性模型。

### \*\*核心思想:\*\*

- 1. \*\*线性组合:\*\* 首先,逻辑回归计算输入特征的线性组合,与线性回归类似: `z = w1\*x1 + w2\*x2 + ... + wn\*xn + b` (或者用向量表示 `z = w^T \* x + b`)。其中 `w` 是权重 (系数) , `b` 是偏置 (截距) 。
- 2. \*\*Sigmoid 函数: \*\* 然后,将这个线性组合的结果 `z` 输入到一个称为 Sigmoid (或 Logistic)的非线性函数中: `p = 1 / (1 + exp(-z))`。
- 3. \*\*概率输出:\*\* Sigmoid 函数的输出 `p` 在 0 到 1 之间,可以被解释为样本属于\*\*正类 \*\* (通常标记为 1) 的概率。
- 4. \*\*决策边界: \*\* 通过设定一个阈值(通常默认为 0.5),将概率 p 转换为最终的类别预测:如果 p >= 0.5,预测为正类 (1);如果 p < 0.5,预测为负类 (0)。决策边界(概率为 0.5的地方)对应于 p = p 、因此逻辑回归学习的是一个线性的决策边界。

#### \*\*优点:\*\*

- \* \*\*实现简单, 计算速度快: \*\* 特别适合大型数据集和需要快速训练的场景。
- \* \*\*易于理解和解释: \*\* 模型的权重 `coef\_` 可以直观地理解为每个特征对预测概率 (的对数几率) 的影响程度。
- \* \*\*输出概率: \*\* `predict\_proba` 方法可以直接给出样本属于每个类别的概率。
- \* \*\*适合作为基线模型: \*\* 由于其简单性和效率,常被用作更复杂模型比较的基准。
- \* \*\*适合作为 Stacking 的元学习器:\*\* 其简单性使其成为组合基模型预测的良好选择 (如你的脚本 `5\_train\_evaluate\_meta\_learner.py`中可以选择使用)。

### \*\*缺点:\*\*

- \* \*\*线性决策边界:\*\* 只能学习线性可分的数据,对于非线性关系需要手动进行特征工程(如添加多项式特征)。
- \* \*\*对特征尺度敏感: \*\* 与其他基于梯度优化的模型一样,建议在使用前进行特征缩放(如标准化)。
- \* \*\*容易欠拟合: \*\* 对于复杂模式可能无法很好地拟合。
- \*\*关键参数 (`LogisticRegression(...)`):\*\*
- \* `penalty`:指定正则化类型,用于防止过拟合。
  - \* `'12'` (默认): L2 正则化 (岭回归) , 使得权重趋向于较小的值。
- \* `'l1'`: L1 正则化 (Lasso) ,可以产生稀疏权重 (某些权重变为 0) ,有助于特征选择。
  - \* `'elasticnet'`: L1 和 L2 的组合。
  - \* ` 'none'`: 不使用正则化。
- \* `C`: \*\*正则化强度的倒数\*\* (默认 1.0)。\*\*较小\*\*的 `C` 值表示\*\*更强\*\*的正则化。这是最重要的调优参数之一。需要通过交叉验证等方式选择最佳值。
- \* `solver`: 用于优化问题的算法。不同的 `solver` 支持不同的 `penalty`。
  - \* `'liblinear'`(默认):对于小型数据集效果不错,支持 L1 和 L2。
- \* `'lbfgs'`, `'newton-cg'`, `'saga'`: 对于大型数据集通常更快, 支持 L2 或 'none' (`'saga'` 还支持 L1 和 'elasticnet')。
- \* `max\_iter`: 优化算法的最大迭代次数 (默认 100)。如果算法不收敛,可能需要增加此值。
- \* `class\_weight`:用于处理不平衡类别。可以设置为`'balanced'`,算法会自动调整权重,使得样本量少的类别获得更高的权重。也可以手动传入一个字典`{class\_label:weight}`。
- \* `multi\_class`: 处理多分类问题的方式。
  - \* `'auto'` (默认): 通常会选择 `'ovr'`。

```
`'ovr'` (One-vs-Rest): 为每个类别训练一个二元分类器。
 `'multinomial'`: 直接优化多项逻辑损失 (通常与 `'lbfgs'`, `'newton-cg'`,
`'sag'` 求解器配合使用)。
代码示例:
```python
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.datasets import make_classification # 用于生成模拟数据
# 1. 生成模拟分类数据
X, y = make_classification(n_samples=1000, n_features=20, n_informative=15,
                         n_redundant=5, random_state=42)
# 2. 划分数据
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,
random_state=42, stratify=y)
# 3. 特征缩放 (重要!)
scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)
# 4. 实例化并训练模型
# 使用 L2 正则化, 调整 C 值 (可以尝试不同的 C 值)
log_reg = LogisticRegression(C=0.1, solver='liblinear', random_state=42)
log_reg.fit(X_train_scaled, y_train)
# 5. 预测
y_pred = log_reg.predict(X_test_scaled)
y_pred_proba = log_reg.predict_proba(X_test_scaled)[:, 1] # 获取正类的概率
# 6. 评估 (将在下一章详细介绍)
from sklearn.metrics import accuracy_score, roc_auc_score
print("Logistic Regression Accuracy:", accuracy_score(y_test, y_pred))
print("Logistic Regression AUC:", roc_auc_score(y_test, y_pred_proba))
# 查看系数(权重)和截距
# print("Coefficients:", log reg.coef )
# print("Intercept:", log_reg.intercept_)
```

6.2 K-近邻 (sklearn.neighbors.KNeighborsClassifier)

K-近邻 (KNN) 是一种简单直观的**非参数、基于实例**的分类算法。它不做显式的模型训练,而是存储整个训练数据集。

核心思想:

- 1. 存储: 在训练阶段, KNN 只是将所有训练样本 (特征和标签) 存储起来。
- 2. 预测: 当需要预测一个新样本时:

- 计算新样本与训练集中**所有**样本之间的**距离**(常用欧氏距离)。
- 。 找出距离新样本最近的 K 个训练样本 ("K 个近邻")。
- 。 根据这 K 个近邻的**多数类别**来决定新样本的预测类别(投票法)。

优点:

- 算法简单, 易于理解和实现。
- 对数据分布没有假设(非参数)。
- 天然适用于多分类问题。
- 训练阶段非常快 (只是存储数据)。

缺点:

- 预测阶段计算量大: 需要计算新样本与所有训练样本的距离, 对于大型数据集可能非常慢。
- 对特征尺度高度敏感: 必须进行特征缩放, 否则距离计算会被尺度大的特征主导。
- 需要选择合适的 K 值: K 太小容易受噪声影响, K 太大容易忽略局部结构。K 值通常通过交叉验证选择。
- **对高维数据效果可能不佳(维度灾难)**:在高维空间中,所有点之间的距离可能趋于相等,使得"近邻"的概念变得模糊。
- 需要存储整个训练集: 内存开销大。

关键参数 (KNeighborsClassifier(...)):

- n_neighbors: **K 的值** (默认 5)。最重要的参数。
- weights: 投票时邻居的权重。
 - 'uniform' (默认): 所有 K 个邻居权重相同。
 - 。 'distance': 权重与距离成反比 (距离越近,权重越大)。通常效果更好。
- metric: 用于计算距离的度量。
 - 'minkowski' (默认) 配合 p 参数:
 - p=1: 曼哈顿距离 (L1)。
 - p=2: 欧氏距离 (L2)。
 - 'euclidean', 'manhattan', 'chebyshev' 等。
- algorithm: 计算最近邻使用的算法。
 - 'auto'(默认): 自动选择最合适的算法 ('ball tree', 'kd tree', 'brute')。
 - 。 'brute': 暴力搜索, 计算所有距离。
 - o 'kd_tree', 'ball_tree':使用树结构加速近邻搜索,在高维时效率可能下降。
- n jobs: 并行计算使用的 CPU 核心数 (-1 表示使用所有核心)。可以加速距离计算。

代码示例:

```
from sklearn.neighbors import KNeighborsClassifier

# (接上例的 X_train_scaled, X_test_scaled, y_train, y_test)

# 4. 实例化并训练模型 (KNN 的 fit 只是存储数据)
knn = KNeighborsClassifier(n_neighbors=5, weights='distance') # 尝试不同的 K 值
knn.fit(X_train_scaled, y_train)

# 5. 预测
y_pred_knn = knn.predict(X_test_scaled)
y_pred_proba_knn = knn.predict_proba(X_test_scaled)[:, 1] # 获取正类的概率 (基于邻居
```

```
比例)
# 6. 评估
print("\nKNN Accuracy:", accuracy_score(y_test, y_pred_knn))
print("KNN AUC:", roc_auc_score(y_test, y_pred_proba_knn))
```

6.3 支持向量机 (sklearn.svm.SVC)

支持向量机 (Support Vector Machine, SVM) 是一种强大的、灵活的分类算法,尤其在处理高维数据和非线性问题时表现出色。其核心思想是找到一个能够将不同类别样本最大间隔地分开的超平面 (Hyperplane)。

核心思想 (简化版):

- 1. **寻找最大间隔超平面:** 在特征空间中,寻找一个决策边界(在高维空间中称为超平面),使得离这个边界最近的两个不同类别的样本点(称为**支持向量 (Support Vectors)**)到该边界的**间隔 (Margin)** 最大化。
- 2. **核技巧 (Kernel Trick):** 对于线性不可分的数据,SVM 使用核函数(如 RBF 核、多项式核)将数据映射到更高维的空间,使得原本不可分的数据在高维空间中变得线性可分,然后再寻找最大间隔超平面。这使得 SVM 能够学习复杂的非线性决策边界。

优点:

- 在高维空间中表现良好: 特别是当特征数量大于样本数量时。
- 模型鲁棒性: 由于决策边界仅由支持向量决定,模型对远离边界的点不敏感。
- 核技巧: 可以通过选择不同的核函数来学习复杂的非线性决策边界。
- 理论基础扎实: 基于统计学习理论。

缺点:

- **计算复杂度高:** 训练时间随样本数量增加而显著增加(尤其是使用非线性核时)。对于非常大的数据集可能不适用。
- 对特征尺度敏感: 必须进行特征缩放。
- 对参数选择敏感: 需要仔细调整参数 C 和核函数参数 (如 gamma)。
- 解释性较差: 不如逻辑回归或决策树直观。
- **原生不支持概率输出**: predict_proba 是通过内部交叉验证估计的(可能较慢),且可能不如逻辑回归的概率校准得好。

关键参数 (SVC(...)):

- C: **正则化参数** (默认 1.0)。与逻辑回归类似,但含义略有不同。它控制了对误分类样本的惩罚程度。
 - **较大的 C**: 惩罚更大,试图将所有训练样本正确分类,可能导致间隔变小,模型更复杂,容易过拟合。
 - **较小的 C:** 惩罚较小,允许一些样本被误分类或处于间隔内,使得间隔更大,模型更简单,泛化能力可能更好。
- kernel: 指定使用的核函数。
 - 。 'linear': 线性核。适用于线性可分数据。
 - 。 'rbf' (默认): 径向基函数核 (高斯核)。适用于非线性数据,通常是首选尝试的核。
 - 'poly': 多项式核。
 - 'sigmoid': Sigmoid 核。
- gamma: 核系数(仅用于 'rbf', 'poly', 'sigmoid')。

- 定义了单个训练样本的影响范围。
- 。 **较小的** gamma: 影响范围大,决策边界更平滑,模型更简单。
- 。 较大的 gamma: 影响范围小,决策边界更复杂,容易过拟合。
- 'scale'(默认): 使用1 / (n_features * X.var()) 作为 gamma 值。
- 'auto':使用1 / n features。
- degree: 多项式核的次数 (仅用于 'poly')。
- probability: 是否启用概率估计 (默认 False)。设置为 True 才能使用 predict_proba,但这会增加训练时间。
- class_weight: 处理不平衡类别,用法同逻辑回归。

代码示例:

```
from sklearn.svm import SVC

# (接上例的 X_train_scaled, X_test_scaled, y_train, y_test)

# 4. 实例化并训练模型

# 使用 RBF 核、调整 C 和 gamma (需要调优)
svc = SVC(C=1.0, kernel='rbf', gamma='scale', probability=True, random_state=42)
svc.fit(X_train_scaled, y_train)

# 5. 预测
y_pred_svc = svc.predict(X_test_scaled)
y_pred_proba_svc = svc.predict_proba(X_test_scaled)[:, 1] # 获取正类的概率

# 6. 评估
print("\nSVC Accuracy:", accuracy_score(y_test, y_pred_svc))
print("SVC AUC:", roc_auc_score(y_test, y_pred_proba_svc))

# 查看支持向量的数量
# print("Number of support vectors for each class:", svc.n_support_)
```

6.4 决策树 (sklearn.tree.DecisionTreeClassifier)

决策树是一种直观的、基于树状结构进行决策的模型。它通过一系列基于特征的"问题"(节点分裂)来将数据逐步划分到不同的叶子节点,每个叶子节点代表一个类别预测。

核心思想:

- 1. **节点分裂:** 从根节点开始,算法选择一个**最佳特征**和**最佳分裂点**,将当前节点的数据分成两个或多个子集(子节点),目标是使得分裂后的子集尽可能"纯净"(即包含的样本尽量属于同一类别)。
- 2. 纯度度量: 常用的纯度度量(或不纯度度量)包括:
 - · 基尼不纯度 (Gini Impurity): 衡量从数据集中随机选择一个样本,然后随机地将其分类,分错的概率。值越小越纯净。
 - 。 **信息熵 (Entropy):** 衡量数据集混乱程度的指标。值越小越纯净。基于信息熵可以计算**信息增益** (Information Gain),选择信息增益最大的特征进行分裂。
- 3. **递归构建:** 对每个子节点递归地重复分裂过程,直到满足停止条件(如达到最大深度、节点样本数过少、节点已纯净等)。
- 4. 叶子节点: 最终不再分裂的节点称为叶子节点,该节点的预测类别通常是该节点中样本数量最多的类别。

优点:

- 易于理解和解释: 决策过程可以可视化为树状图, 非常直观。
- 能够处理数值型和类别型特征: 不需要像线性模型那样必须编码类别特征(虽然 Scikit-learn 的实现要求输入是数值)。
- 对特征尺度不敏感: 不需要进行特征缩放。
- 能够捕捉特征之间的交互关系。
- 训练速度相对较快。

缺点:

- 容易过拟合: 决策树倾向于生成非常复杂的树来完美拟合训练数据,导致泛化能力差。需要进行剪枝(限制树的生长)来缓解。
- 不稳定性: 数据微小的变动可能导致生成完全不同的树。
- 对最优决策树的搜索是 NP 难问题: 实际算法通常采用贪心策略(每一步选择局部最优分裂),不能保证 找到全局最优树。
- 可能产生偏向: 对于某些特征(如具有较多取值水平的特征)可能存在偏向。

关键参数 (DecisionTreeClassifier(...)):

- criterion: 选择分裂节点的标准。
 - 'gini'(默认): 使用基尼不纯度。
 - 'entropy': 使用信息熵 (计算信息增益)。
- splitter: 选择分裂点策略。
 - 'best' (默认): 选择最佳分裂点。
 - 'random': 选择随机分裂点 (通常用于集成模型中增加随机性)。
- max_depth: 树的最大深度 (默认 None, 不限制)。最重要的控制过拟合的参数之一。 通常需要调优。
- min samples split: 节点继续分裂所需的最少样本数 (默认 2)。增加此值可以防止过拟合。
- min_samples_leaf: **叶子节点所需的最少样本数**(默认 1)。增加此值可以使树更平滑,防止过拟合。
- max_features: 在寻找最佳分裂时考虑的特征数量。可以设置为整数、浮点数(比例)或特定字符串('sqrt','log2')。限制特征数量可以增加树的多样性(用于集成)。
- random_state: 控制随机性 (例如当 splitter='random' 或 max_features < n_features 时) 。
- class_weight: 处理不平衡类别。

代码示例:

```
from sklearn.tree import DecisionTreeClassifier
# from sklearn.tree import plot_tree # 用于可视化
# import matplotlib.pyplot as plt

# (接上例的 X_train, X_test, y_train, y_test - 注意: 决策树不需要缩放特征)

# 4. 实例化并训练模型
# 限制树的深度来防止过拟合
dtree = DecisionTreeClassifier(criterion='gini', max_depth=5, min_samples_leaf=10, random_state=42)
dtree.fit(X_train, y_train) # 可以使用未缩放的数据

# 5. 预测
y_pred_dtree = dtree.predict(X_test)
```

```
y_pred_proba_dtree = dtree.predict_proba(X_test)[:, 1]

# 6. 评估
print("\nDecision Tree Accuracy:", accuracy_score(y_test, y_pred_dtree))
print("Decision Tree AUC:", roc_auc_score(y_test, y_pred_proba_dtree))

# 7. (可选) 可视化决策树 (对于较小的树)
# plt.figure(figsize=(20,10))
# plot_tree(dtree, filled=True, feature_names=[f'feature_{i}' for i in range(X.shape[1])], class_names=['0', '1'], rounded=True, fontsize=10)
# plt.show()
```

6.5 随机森林 (sklearn.ensemble.RandomForestClassifier)

随机森林是一种强大且广泛使用的**集成学习**算法,它通过构建多个决策树并将它们的预测结果进行**集成**(分类问题通常是投票,回归问题通常是平均)来提高模型的性能和鲁棒性。

核心思想 (Bagging + 特征随机):

- 1. **自助采样 (Bootstrap Sampling / Bagging):** 从原始训练集中**有放回地**随机抽取多个子样本集(每个子样本集大小与原始训练集相同)。由于是有放回抽样,每个子样本集中会包含一些重复样本,同时也会遗漏一些原始样本(约 37%)。
- 2. 独立训练决策树: 为每个子样本集独立地训练一棵决策树。
- 3. **特征随机性:** 在每个节点进行分裂时,**不是**考虑所有特征,而是**随机选择一部分特征** (max_features), 然后从这部分特征中选择最佳分裂点。这进一步增加了每棵树之间的差异性。
- 4. **集成预测:** 对于分类任务,将所有决策树的预测结果进行**多数投票**,得到最终的预测类别。对于回归任务,计算所有决策树预测值的**平均值**。

优点:

- 高准确率: 通常比单棵决策树性能更好。
- 强大的抗过拟合能力: 通过 Bagging 和特征随机性,有效地降低了模型的方差,减少了过拟合风险。
- 鲁棒性强: 对噪声和异常值相对不敏感。
- 能够处理高维数据: 特征随机性使其在高维数据上表现良好。
- 能够评估特征重要性: 可以根据特征在所有树中对不纯度减少的贡献来评估其重要性。
- 易于并行化: 每棵树可以独立训练。

缺点:

- 模型复杂度高,解释性差:不如单棵决策树直观,是一个"黑盒"模型。
- 训练时间和内存占用较大: 需要训练多棵树。
- 对于某些非常稀疏或特定结构的数据可能不如其他模型。

关键参数 (RandomForestClassifier(...)):

- n_estimators: 森林中决策树的数量 (默认 100)。通常值越大越好,但会增加计算成本,超过一定数量后性能提升会饱和。
- criterion, max_depth, min_samples_split, min_samples_leaf: 与 DecisionTreeClassifier 中的 参数相同,用于控制**单棵树**的生长和复杂度。是防止过拟合的关键调优参数。

max_features: 每个节点分裂时随机选择的特征数量 (默认 'sqrt', 即 sqrt(n_features))。减小max_features 会增加树之间的随机性,可能降低过拟合,但可能需要更多的树 (n_estimators)。
 常用值: 'sqrt', 'log2',整数,浮点数 (比例)。

- bootstrap: 是否使用自助采样来构建树 (默认 True)。
- oob_score: 是否使用袋外 (Out-of-Bag) 样本来估计泛化准确率 (默认 False)。OOB 样本是指在构建某 棵树时未被该树的自助采样选中的样本。设置为 True 可以提供一种无需额外验证集的模型性能估计。
- random state: 控制自助采样和特征选择的随机性。
- class_weight: 处理不平衡类别。
- n jobs: 并行训练使用的 CPU 核心数 (-1 使用所有核心)。

代码示例:

```
from sklearn.ensemble import RandomForestClassifier
# (接上例的 X_train, X_test, y_train, y_test - 随机森林也不需要缩放特征)
# 4. 实例化并训练模型
# 调整 n estimators 和控制树复杂度的参数
rf = RandomForestClassifier(n_estimators=200, max_depth=10, min_samples_leaf=5,
                           max_features='sqrt', random_state=42, n_jobs=-1,
                           oob score=True, class weight='balanced') # 使用
oob score 和 class weight
rf.fit(X_train, y_train)
# 5. 预测
y_pred_rf = rf.predict(X_test)
y_pred_proba_rf = rf.predict_proba(X_test)[:, 1]
# 6. 评估
print("\nRandom Forest Accuracy:", accuracy_score(y_test, y_pred_rf))
print("Random Forest AUC:", roc auc score(y test, y pred proba rf))
if rf.oob score :
    print("Random Forest OOB Score:", rf.oob score ) # 00B 分数是对训练数据的估计
# 7. 查看特征重要性
# feature_importances = rf.feature_importances_
# importance_df = pd.DataFrame({'Feature': [f'feature_{i}' for i in
range(X.shape[1])],
                               'Importance': feature importances})
# importance_df = importance_df.sort_values(by='Importance', ascending=False)
# print("\nFeature Importances:\n", importance df.head())
```

其他常用分类器 (简介):

- **朴素贝叶斯 (sklearn.naive_bayes)**: 基于贝叶斯定理和特征条件独立性假设的一类简单高效的分类器 (GaussianNB 用于连续特征, MultinomialNB 用于离散计数特征, BernoulliNB 用于二元特征)。速度 快,适合文本分类等。
- **梯度提升决策树 (sklearn.ensemble.GradientBoostingClassifier):** 另一种强大的集成方法 (Boosting),通过迭代地训练新的树来修正前面树的错误。XGBoost 是其更高级、更优化的实现。我 们将在第八章详细讨论 XGBoost。

模型选择小结:

- **基线模型**: 逻辑回归。
- 需要解释性:逻辑回归、决策树。
- **非线性问题:** KNN (需缩放), SVM (需缩放, 调核), 决策树, 随机森林, GBDT/XGBoost。
- 高维数据: 逻辑回归 (配合正则化), SVM (线性核或 RBF 核), 随机森林。
- 大型数据集: 逻辑回归, 朴素贝叶斯, 随机森林 (可并行), GBDT/XGBoost (优化实现)。 SVM (非线性核) 和 KNN 可能较慢。
- 性能优先 (结构化数据): 随机森林, XGBoost (通常是首选)。

选择哪个模型没有绝对的答案,通常需要根据数据特点、问题需求和计算资源进行实验和比较。理解它们的工作原理和关键参数是进行有效选择和调优的基础。

第七章:模型评估:度量你的模型

训练模型只是机器学习流程的一部分,同样重要的是**如何客观、准确地评估模型的性能**。我们需要选择合适的评估指标来衡量模型在解决特定业务问题上的效果,并理解这些指标的含义和局限性。本章将重点介绍用于**分类问题**的常用评估指标。

7.1 为什么需要多种评估指标?

仅仅使用准确率 (Accuracy) 作为唯一的评估指标往往是不够的,甚至可能产生误导,尤其是在以下情况:

- **类别不平衡** (Imbalanced Classes): 当数据集中某个类别的样本数量远超其他类别时。例如,在欺诈检测中,绝大多数交易是正常的,只有极少数是欺诈。一个简单地将所有样本预测为"正常"的模型也能获得极高的准确率(例如 99%),但这显然是一个无用的模型,因为它完全无法识别出欺诈交易。
- 不同错误类型的代价不同: 在某些场景下,不同类型的错误(例如,漏报 vs. 误报)带来的后果严重程度不同。
 - **示例 (医疗诊断):** 将癌症患者误诊为健康(假阴性/FN)的代价通常远高于将健康人误诊为癌症(假阳性/FP)。我们需要更关注召回率。
 - **示例 (垃圾邮件过滤):** 将重要邮件误判为垃圾邮件(假阳性/FP)的代价可能高于漏掉一封垃圾邮件(假阴性/FN)。我们需要更关注精确率。
 - **示例 (天气预报):** 漏报一次大雨 (假阴性/FN) 的代价可能很高,但频繁误报小雨 (假阳性/FP) 也会降低预报的可信度。需要综合考虑。

因此,我们需要根据具体问题选择能够反映业务目标的、更全面的评估指标。

7.2 混淆矩阵 (sklearn.metrics.confusion matrix)

混淆矩阵是理解分类模型性能的基础,它以矩阵形式展示了模型预测结果与真实标签之间的关系。对于二元分类问题,混淆矩阵通常是 2x2 的形式:

	预测为 0 (阴性)	预测为 1 (阳性)
实际为 0 (阴性)	TN (True Negative)	FP (False Positive)
	FN (False Negative)	TP (True Positive)

- TP (True Positive, **真阳性**): 样本**实际为** 1,模型**预测也为** 1。(预测正确)
 - · 天气示例: 实际下雨, 预测也下雨。 (命中)

- TN (True Negative, **真阴性**): 样本**实际为 0**,模型**预测也为 0**。(预测正确)
 - · 天气示例: 实际未下雨, 预测也未下雨。 (正确否定)
- FP (False Positive, 假阳性 / Type I Error): 样本实际为 0, 但模型预测为 1。(预测错误)
 - 。 *天气示例*: 实际未下雨, 但预测下雨。 (误报/虚警)
- FN (False Negative, 假阴性 / Type II Error): 样本实际为 1,但模型预测为 0。(预测错误)
 - 。 *天气示例*: 实际下雨, 但预测未下雨。 (漏报/未命中)

应用关联: 你的脚本 1_generate_base_predictions.py 的核心任务之一就是根据基模型的预测 (y_train_pred) 和真实标签 (y_train_true) 来识别训练集中的 TP, TN, FP, FN 样本,并打印它们的数量,这实际上就是在计算混淆矩阵的四个基本元素。

使用 Scikit-learn 计算混淆矩阵:

```
from sklearn.metrics import confusion matrix
import matplotlib.pyplot as plt
import seaborn as sns # 用于更美观的热力图
# 假设 y_true 是真实标签, y_pred 是模型预测标签 (来自某个模型)
# y_true = y_test
# y_pred = y_pred_log_reg # 或 y_pred_knn, y_pred_svc, ...
# (使用之前的逻辑回归预测结果作为示例)
y_true = y_test
y_pred = log_reg.predict(X_test_scaled) # 来自 6.1 节
cm = confusion_matrix(y_true, y_pred)
print("Confusion Matrix:\n", cm)
# 输出示例:
# Confusion Matrix:
# [[136 14] <-- TN=136, FP=14
# [ 23 127]] <-- FN=23, TP=127
# 可视化混淆矩阵 (使用 seaborn)
plt.figure(figsize=(6, 4))
sns.heatmap(cm, annot=True, fmt='d', cmap='Blues',
           xticklabels=['Predicted 0', 'Predicted 1'],
           yticklabels=['Actual 0', 'Actual 1'])
plt.xlabel('Predicted Label')
plt.ylabel('True Label')
plt.title('Confusion Matrix')
plt.show()
# 提取 TP, TN, FP, FN (如果需要单独使用)
tn, fp, fn, tp = cm.ravel()
print(f"TN: {tn}, FP: {fp}, FN: {fn}, TP: {tp}")
```

7.3 基于混淆矩阵的核心指标

许多重要的分类指标都可以从混淆矩阵的四个基本元素计算得出。

7.3.1 准确率 (Accuracy)

- 公式: Accuracy = (TP + TN) / (TP + TN + FP + FN) (预测正确的样本 / 总样本)
- 含义: 模型预测正确的样本占总样本的比例。
- Scikit-learn: sklearn.metrics.accuracy_score(y_true, y_pred)
- 局限性: 在类别不平衡时具有误导性。

```
from sklearn.metrics import accuracy_score
acc = accuracy_score(y_true, y_pred)
print(f"Accuracy: {acc:.4f}")
# 等价于: (tp + tn) / (tp + tn + fp + fn)
```

7.3.2 精确率 (Precision)

- 公式: Precision = TP / (TP + FP)
- 含义: 在所有被模型预测为正类 (1) 的样本中,实际也为正类的比例。衡量模型预测的"准确性"或"查准率"。
- 关注点: 最小化 FP (假阳性)。当 FP 的代价很高时,精确率是重要指标。
 - · 天气示例: 提高精确率意味着减少误报(预测下雨但实际没下)。
- **Scikit-learn:** sklearn.metrics.precision_score(y_true, y_pred, pos_label=1, average='binary')
 - o pos_label: 指定哪个标签是正类 (默认 1) 。
 - o average: 用于多分类。'binary' 仅用于二分类。

```
from sklearn.metrics import precision_score
precision = precision_score(y_true, y_pred)
print(f"Precision: {precision:.4f}")
# 等价于: tp / (tp + fp) if (tp + fp) > 0 else 0
```

7.3.3 召回率 (Recall) / 敏感度 (Sensitivity) / 命中率 (Hit Rate) / 真阳性率 (TPR) / POD

- 公式: Recall = TP / (TP + FN)
- 含义: 在所有实际为正类 (1) 的样本中,被模型成功预测为正类的比例。衡量模型"找回"正类的能力或"查全率"。
- 关注点: 最小化 FN (假阴性)。当 FN 的代价很高时, 召回率是重要指标。
 - *天气示例*: 提高召回率 (POD) 意味着减少漏报 (实际下雨但预测未下)。
- **Scikit-learn:** sklearn.metrics.recall_score(y_true, y_pred, pos_label=1, average='binary')

```
from sklearn.metrics import recall_score
recall = recall_score(y_true, y_pred)
print(f"Recall (POD): {recall:.4f}")
# 等价于: tp / (tp + fn) if (tp + fn) > 0 else 0
```

精确率 vs. 召回率 (Precision-Recall Trade-off): 通常情况下,精确率和召回率是相互制约的。提高一个指标往往会导致另一个指标下降。例如,如果我们降低分类阈值,使得模型更容易预测为正类,那么召回率会提高(找到更多实际为正的样本),但精确率可能会下降(因为一些实际为负的样本也可能被预测为正,增加了FP)。选择合适的平衡点取决于业务需求。

7.3.4 F1 分数 (F1-Score)

- 公式: F1 = 2 * (Precision * Recall) / (Precision + Recall)
- **含义**: 精确率和召回率的**调和平均数 (harmonic mean)**。它综合考虑了精确率和召回率,当两者都较高时,F1 分数也会较高。
- 适用场景: 当精确率和召回率同等重要时, 或者想找到一个在两者之间取得较好平衡的模型时。
- **Scikit-learn:** sklearn.metrics.f1_score(y_true, y_pred, pos_label=1, average='binary')

```
from sklearn.metrics import f1_score
f1 = f1_score(y_true, y_pred)
print(f"F1 Score: {f1:.4f}")
```

7.3.5 特异度 (Specificity) / 真阴性率 (TNR)

- 公式: Specificity = TN / (TN + FP)
- 含义: 在所有实际为负类 (0) 的样本中,被模型成功预测为负类的比例。衡量模型正确识别负类的能力。
- 关系: Specificity = 1 FPR (FPR 是假阳性率, 见下文 ROC)。
- Scikit-learn: 没有直接的函数,但可以通过混淆矩阵计算 tn / (tn + fp)。

7.3.6 分类报告 (sklearn.metrics.classification report)

提供了一个文本摘要,包含了每个类别的主要分类指标(精确率、召回率、F1 分数)以及总体准确率和宏/加权平均指标。

```
from sklearn.metrics import classification report
report = classification_report(y_true, y_pred, target_names=['Class 0', 'Class
1'])
print("Classification Report:\n", report)
# 输出示例:
# Classification Report:
              precision recall f1-score support
# Class 0 0.86
                          0.91
                                  0.88
                                           150 # 对应实际为 0 的精确率、召
回率、F1
# Class 1
                0.90
                          0.85
                                  0.87
                                           150 # 对应实际为 1 的精确率、召
回率、F1
                                            300 # 总体准确率
                                   0.88
   accuracy
                                            300 # 宏平均 (简单平均各类的指
#
  macro avg
                 0.88
                          0.88
                                  0.88
标)
# weighted avg
                 0.88
                          0.88
                                  0.88
                                            300 # 加权平均(按各类样本数
support 加权)
```

```
#
# support 列显示每个类别的真实样本数量。
```

7.4 基于概率和阈值的指标

之前的指标都是基于模型最终的类别预测(通常使用 0.5 作为概率阈值)。但模型的 predict_proba 方法提供了更丰富的信息,我们可以通过改变阈值来调整模型的行为,并使用 ROC 曲线和 AUC 来评估模型在**所有可能 阈值下**的整体性能。

7.4.1 ROC 曲线 (Receiver Operating Characteristic Curve)

- 绘制: ROC 曲线绘制的是在不同分类阈值下,真阳性率 (TPR / Recall / POD) 相对于 假阳性率 (FPR) 的变化情况。
 - FPR = FP / (FP + TN) (在所有实际为负类的样本中,被错误预测为正类的比例)。
- 解读:
 - 。 曲线越靠近左上角 (TPR 高, FPR 低) 表示模型性能越好。
 - 对角线(从(0,0)到(1,1))代表随机猜测的模型。
 - 曲线完全贴合左边界和上边界代表完美分类器 (TPR=1, FPR=0)。
- Scikit-learn:
 - sklearn.metrics.roc_curve(y_true, y_score): 计算 ROC 曲线的 FPR, TPR 和对应的阈值。
 y_score 通常是模型预测为**正类**的概率 (predict_proba(X)[:, 1])。
 - o sklearn.metrics.RocCurveDisplay.from_estimator(estimator, X, y) 或 from_predictions(y_true, y_pred_proba):直接绘制 ROC 曲线。

7.4.2 AUC (Area Under the ROC Curve)

- **含义**: ROC 曲线下的面积。取值范围在 0 到 1 之间。
- 解读:
 - AUC = 1: 完美分类器。
 - AUC = 0.5: 随机猜测。
 - AUC < 0.5: 模型性能差于随机猜测(可能标签反了)。
 - · AUC 越大,通常表示模型区分正负样本的能力越强。
- 优点: AUC 是一个单一数值,便于比较不同模型的整体性能,且对类别不平衡和阈值选择不敏感。
- **Scikit-learn:** sklearn.metrics.roc_auc_score(y_true, y_score)

代码示例 (ROC & AUC):

```
from sklearn.metrics import roc_curve, auc, roc_auc_score, RocCurveDisplay

# (使用之前的逻辑回归概率预测结果作为示例)
y_score = log_reg.predict_proba(X_test_scaled)[:, 1] # 正类的概率

# 1. 计算 AUC
auc_score = roc_auc_score(y_true, y_score)
print(f"\nAUC Score: {auc_score:.4f}")

# 2. 计算 ROC 曲线的点
fpr, tpr, thresholds = roc_curve(y_true, y_score)
```

```
# 3. 绘制 ROC 曲线 (方法一: 手动绘制)
plt.figure(figsize=(7, 6))
plt.plot(fpr, tpr, color='darkorange', lw=2, label=f'ROC curve (area =
{auc score:.2f})')
plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--') # 对角线
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.05])
plt.xlabel('False Positive Rate (FPR)')
plt.ylabel('True Positive Rate (TPR / Recall / POD)')
plt.title('Receiver Operating Characteristic (ROC) Curve')
plt.legend(loc="lower right")
plt.grid(alpha=0.3)
plt.show()
# 3. 绘制 ROC 曲线 (方法二: 使用 RocCurveDisplay)
# display = RocCurveDisplay.from_estimator(log_reg, X_test_scaled, y_test)
# display = RocCurveDisplay.from predictions(y true, y score, name='Logistic
Regression')
# plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
# plt.title('ROC Curve (using Display)')
# plt.show()
```

7.4.3 精确率-召回率 (PR) 曲线

• **绘制:** PR 曲线绘制的是在不同分类阈值下,**精确率 (Precision)** 相对于 **召回率 (Recall / TPR / POD)** 的变化情况。

• 解读:

- 。 曲线越靠近右上角 (Precision 高, Recall 高) 表示模型性能越好。
- 。 对于类别**高度不平衡**的数据集,PR 曲线比 ROC 曲线更能反映模型在少数正类上的表现。因为 ROC 的 FPR 计算涉及大量的 TN,容易使得 AUC 看起来很高,而 PR 曲线则更关注正类的预测情况 (TP, FP, FN)。

Scikit-learn:

- o sklearn.metrics.precision_recall_curve(y_true, y_score): 计算 PR 曲线的 Precision, Recall 和对应的阈值。
- o sklearn.metrics.PrecisionRecallDisplay.from_estimator(...)或 from_predictions(...):直接绘制 PR 曲线。
- o sklearn.metrics.average_precision_score(y_true, y_score): 计算 PR 曲线下的面积 (Average Precision, AP)。

7.5 领域特定指标(以天气预报为例)

在特定领域,可能会使用一些组合指标或特定名称的指标。你的脚本中多次提到了 POD, FAR, CSI, 它们在气象领域非常常用。

- POD (Probability of Detection): 同召回率 (Recall)。
 - \circ POD = TP / (TP + FN)
 - 。 衡量**实际发生**的事件(如下雨)被**正确预报**的比例。越高越好。
- FAR (False Alarm Ratio):

- \circ FAR = FP / (TP + FP)
- 衡量在所有**预报发生**的事件(如预报下雨)中,**实际未发生**(误报)的比例。越低越好。**注意**: 这与 ROC 中的 FPR (False Positive Rate) 不同,FPR 的分母是所有实际未发生的事件 (FP + TN)。
- CSI (Critical Success Index / Threat Score):
 - \circ CSI = TP / (TP + FN + FP)
 - 综合考虑了命中 (TP)、漏报 (FN) 和误报 (FP)。它衡量的是在所有实际发生或被预报发生的事件中,被正确预报的比例。取值范围 0 到 1, 越高越好。CSI 对 TP, FN, FP 都敏感。

计算示例 (基于混淆矩阵):

```
# (tn, fp, fn, tp 来自之前的混淆矩阵 cm.ravel())

pod = tp / (tp + fn) if (tp + fn) > 0 else 0

far = fp / (tp + fp) if (tp + fp) > 0 else 0

csi = tp / (tp + fn + fp) if (tp + fn + fp) > 0 else 0

print(f"\nPOD (Recall): {pod:.4f}")

print(f"FAR: {far:.4f}")

print(f"CSI: {csi:.4f}")
```

应用关联: 你的脚本 1, 2, 3, 5, 6 中都包含了 calculate_metrics 函数,用于计算并打印包括 Accuracy, POD, FAR, CSI 在内的多个指标。这表明在你的天气预测任务中,综合评估模型的命中能力、误报情况和整体威胁评分是非常重要的。脚本中通过遍历不同的概率阈值来评估这些指标,可以帮助找到一个在 POD, FAR, CSI 之间取得理想平衡的阈值。

7.6 选择合适的评估指标

最终选择哪个或哪些指标来评估和优化模型, 取决于:

- **业务目标**: 最关心什么? 是尽可能找出所有正类(高 Recall/POD),还是确保预测为正类的尽可能准确(高 Precision),或者是减少误报(低 FAR)?
- **数据特征**: 类别是否平衡?
- 模型应用场景: 不同错误的代价是什么?

通常,不会只依赖单一指标,而是会综合考虑多个指标,例如:

- 平衡数据集: Accuracy, AUC 可能是不错的起点,但也应关注 Precision, Recall, F1。
- **不平衡数据集:** 优先考虑 Precision, Recall, F1, AUC, PR-AUC。Accuracy 可能具有误导性。
- **关注减少漏报 (FN):** 优先考虑 Recall (POD)。
- **关注减少误报 (FP):** 优先考虑 Precision 或 FAR。
- 需要综合评估命中、漏报、误报: CSI 是一个常用选择。
- 比较模型整体排序能力: AUC 是常用选择。

理解各种评估指标的计算方式、含义和适用场景,是有效评估、比较和选择机器学习模型的关键一步。

```
```markdown
机器学习实战: 从 Scikit-learn 到 XGBoost (卷三)
...(接上文)
```

---

## 第八章: 梯度提升王者: XGBoost 深入解析

在之前的章节中,我们介绍了随机森林这种基于 Bagging 思想的强大集成模型。本章我们将深入探讨另一种更先进、在许多结构化数据竞赛和实际应用中表现更为出色的集成算法—\*\*梯度提升决策树 (Gradient Boosting Decision Tree, GBDT)\*\*, 特别是其高效、优化的实现: \*\*XGBoost (Extreme Gradient Boosting)\*\*。

XGBoost 是许多数据科学家工具箱中的利器,你的天气预测系列脚本也深度依赖它作为核心模型 (基模型、专家模型、元学习器、微调模型)。理解 XGBoost 的工作原理和关键参数对于有效使用和调优这些脚本至关重要。

### 8.1 从 GBDT 到 XGBoost: 提升思想的演进

#### 8.1.1 提升 (Boosting) 的核心思想

与 Bagging (如随机森林) 并行训练多个独立学习器不同, Boosting 是一种\*\*串行\*\*的集成学习方法。它迭代地训练一系列\*\*弱学习器\*\*(通常是决策树),每一轮训练的新学习器都\*\*专注于修正前面学习器所犯的错误\*\*。最终的模型是所有弱学习器的加权组合。

#### 8.1.2 梯度提升决策树 (GBDT)

GBDT (Gradient Boosting Decision Tree) 是 Boosting 思想的一个具体实现,它巧妙地利用了\*\*梯度下降\*\*的思想来优化损失函数:

- 1. \*\*初始化:\*\* 首先,用一个简单的模型(通常是目标变量的均值)进行初始预测。
- 2. \*\*迭代计算残差:\*\* 在每一轮迭代 `m` 中:
- \* 计算当前集成模型预测结果与真实标签之间的\*\*残差\*\*(或者更准确地说,是损失函数关于当前预测值的\*\*负梯度\*\*)。这些残差代表了模型当前"未能拟合"的部分。
- 3. \*\*训练新树拟合残差:\*\* 训练一棵新的决策树 `h\_m(x)` 来拟合上一轮计算得到的残差 (负梯度)。这棵树的目标是学习如何弥补当前模型的不足。
- 4. \*\*更新模型:\*\* 将新训练的树以一定的\*\*学习率 (learning rate)\*\* 加入到集成模型中,更新预测结果:  $F_m(x) = F_{m-1}(x) + learning_rate * h_m(x)$ 。学习率用于控制每棵树对最终结果的贡献,防止单棵树影响过大。
- 5. \*\*重复:\*\* 重复步骤 2-4 指定的轮数 (树的数量)。

GBDT 通过逐步减少残差(或损失函数的梯度),使得集成模型的预测结果越来越接近真实标签。

#### 8.1.3 XGBoost 的优势: 为何选择它?

XGBoost 在 GBDT 的基础上进行了多项重要改进,使其在性能、速度和灵活性上都更胜一筹:

- 1. \*\*正则化 (Regularization): \*\* GBDT 容易过拟合, XGBoost 在目标函数中加入了 L1 (Lasso) 和 L2 (Ridge) 正则化项(体现在参数 `reg\_alpha` 和 `reg\_lambda`), 惩罚模型的复杂度(树的叶子节点数量和叶子节点权重), 有助于防止过拟合,提高泛化能力。
- 2. \*\*优化目标函数: \*\* XGBoost 对损失函数进行了二阶泰勒展开,引入了二阶导数信息,使得目标函数的优化更精确、收敛更快。
- 3. \*\*内置处理缺失值: \*\* XGBoost 能够自动学习如何处理缺失值,将缺失样本划分到默认方向(在训练过程中学习哪个方向更好)。
- 4. \*\*并行化与缓存优化:\*\*
  - \* \*\*特征粒度并行:\*\* 在节点分裂查找最佳分裂点时,可以并行计算不同特征的增益。
  - \* \*\*块结构设计: \*\* 数据在内存中以特定的块(Block)结构存储,支持列式存储和访问,便

于并行计算,并优化了缓存命中率。

- 5. \*\*高效的树构建算法:\*\*
- \* \*\*预排序 (Pre-sorted) 算法:\*\* GBDT 常用的方法,对特征值预排序,然后遍历查找最佳分裂点。
- \* \*\*直方图算法 (`tree\_method='hist'`):\*\* XGBoost 引入的关键优化。将连续特征值分桶(放入直方图), 然后基于分桶的统计信息查找最佳分裂点。大大减少了候选分裂点的数量,显著提高了训练速度,降低了内存消耗,尤其适用于\*\*大型数据集\*\*。
- \* \*\*应用关联:\*\* 你的脚本 `2`, `3`, `5`, `6` 在定义 XGBoost 模型时都设置了 `tree\_method='hist'`, 这正是利用了该算法处理大数据的优势。
- 6. \*\*内置交叉验证:\*\* 可以在训练过程中执行交叉验证(`xgb.cv`)。
- 7. \*\*支持自定义目标函数和评估指标。\*\*
- 8. \*\*剪枝 (Pruning): \*\* 除了预剪枝 (如 `max\_depth`, `min\_child\_weight`), XGBoost 还使用了后剪枝策略 (基于 `gamma` 参数)。当分裂带来的增益小于 `gamma` 时,不进行分裂。

这些优化使得 XGBoost 成为处理结构化数据 (表格数据) 的强大、快速且可靠的选择。

### 8.2 XGBoost 核心 API (`xgboost.XGBClassifier`)

XGBoost 提供了两种主要的 API: 原生 API 和 Scikit-learn 风格的 API。Scikit-learn API (`XGBClassifier` 用于分类, `XGBRegressor` 用于回归) 与 Scikit-learn 的 Estimator 接口完全兼容,可以无缝集成到 Scikit-learn 的工作流中(如 `Pipeline`, `GridSearchCV`)。你的脚本使用的正是这种 Scikit-learn 风格的 API。

```
实例化:
```python
import xgboost as xgb
# 创建一个 XGBoost 分类器实例
model = xgb.XGBClassifier(
   # --- 核心参数 ---
   n_estimators=100, # 树的数量 ( boosting rounds)
                    # 学习率 (eta)
# 每棵树的最大深度
   learning_rate=0.1,
   max depth=3,
   # --- 控制过拟合 ---
                    # 训练样本采样比例 (行采样)
   subsample=0.8,
   colsample_bytree=0.8, #每棵树的特征采样比例 (列采样)
                    # 分裂最小损失下降 (min split loss)
   gamma=0,
   reg_alpha=0,
                    # L1 正则化系数
                     # L2 正则化系数 (默认)
   reg lambda=1,
   # --- 处理不平衡数据 ---
   scale pos weight=1,
                    # 正样本权重 (默认1,适用于平衡数据)
   # --- 性能与效率 ---
   tree_method='auto', # 建树算法 ('hist' 推荐用于大数据)
                    # 并行核心数 (-1 使用所有)
   n jobs=-1,
   # --- 其他 ---
   objective='binary:logistic', # 目标函数 (二分类对数损失)
   eval_metric='logloss', # 评估指标 (用于早停等)
   random state=42, # 随机种子
   use_label_encoder=False # 重要: 新版本推荐设置为 False
)
```

8.3 关键参数详解

理解并调整 XGBoost 的参数是获得良好性能的关键。

8.3.1 通用参数 (General Parameters)

- booster: 指定使用的提升器类型。通常是 'gbtree' (基于树的模型) 或 'gblinear' (线性模型)。默认 'gbtree'。
- n_jobs: 并行运行的线程数。-1 表示使用所有可用线程。
- random_state / seed: 随机数种子,用于复现结果。

8.3.2 提升器参数 (Booster Parameters - 主要针对 gbtree)

控制模型复杂度与过拟合:

- learning_rate (或 eta, 默认 0.3): 学习率。在每次迭代后缩减新树的权重,以减少每一步的影响,使得模型更加稳健。
 - 典型值: 0.01 到 0.3。
 - **权衡: 较低**的 learning_rate 需要**较多**的 n_estimators 才能达到相同的性能,但通常能获得更好的泛化能力。
 - o **应用关联:** 脚本 6 (微调) 中使用了较低的 **FINETUNE_LR** = 0.02, 这是 Fine-tuning 的常见做法。 脚本 2, 3, 5 使用了 0.1。
- n_estimators: 总共拟合的树的数量 (Boosting 的轮数)。这是最重要的参数之一。值越大模型越复杂,但也越容易过拟合。通常需要与 learning_rate 和 early_stopping_rounds 配合调优。
 - **应用关联:** 脚本 2, 3 设置为 500, 脚本 5 设置为 100, 脚本 6 (微调) 设置为 50。
- max_depth (默认 6): **每棵树的最大深度**。增加深度会使模型能够学习更复杂的特征交互,但也更容易过拟合。
 - 典型值: 3 到 10。
 - o **应用关联**: 脚本 2, 3 设置为 7, 脚本 5 设置为 3。
- min_child_weight (默认 1): 定义了一个子节点所需的**最小样本权重和** (hessian sum) 。如果一个叶子节点的样本权重和小于 min_child_weight,则分裂过程会停止。在线性回归模型中,它大致对应于每个节点所需的最小样本数。值越大,模型越保守(不易过拟合)。
- gamma (或 min_split_loss, 默认 0): 节点分裂所需的**最小损失下降**。只有当分裂带来的损失减少大于 gamma 时,节点才进行分裂。值越大,算法越保守(剪枝效果越强)。
 - 応用关联: 脚本 2, 3, 5 设置为 0.1。
- subsample (默认 1): **训练每棵树时使用的样本(行)比例**。设置为小于 1 的值(例如 0.8)意味着每次 迭代时随机抽取 80% 的训练数据来构建树,这可以防止过拟合,并使训练更快。
 - **应用关联:** 脚本 2, 3, 5 设置为 0.8 或 0.9。
- colsample_bytree (默认 1): **构建每棵树时使用的特征 (列) 比例**。例如, 0.8 表示每棵树随机选择 80% 的特征进行训练。有助于防止过拟合和提高训练速度。
 - **应用关联:** 脚本 2, 3, 5 设置为 0.8 或 0.9。
- colsample_bylevel (默认 1): 树的每一层分裂时使用的特征比例。
- colsample bynode (默认 1): 树的每个节点分裂时使用的特征比例。
- reg_alpha (L1 正则化, 默认 0): L1 正则化项的权重。增加此值会使模型更保守,可能导致叶子权重稀疏。
- reg_lambda (L2 正则化, 默认 1): L2 正则化项的权重。增加此值会使模型更保守,使得叶子权重更平滑。

处理不平衡数据:

- scale_pos_weight (默认 1): 控制正负样本的权重平衡。常用于处理类别不平衡的二分类问题。
 - **推荐设置:** sum(negative instances) / sum(positive instances)。即,负样本总数除以正样本总数。这会给样本量较少的正类赋予更高的权重,使得模型更关注少数类。
 - o **应用关联:** 脚本 2 (FN 专家, FN 是正类)、3 (FP 专家, FP 是正类)、5 (元学习器)、6 (微调) 都计算并设置了 scale_pos_weight,表明这些任务都可能面临类别不平衡的问题,需要特别处理。

建树算法:

- tree_method (默认 'auto'):
 - 'hist': 基于直方图的快速算法,推荐用于大型数据集。
 - 'exact': 精确贪心算法, 计算所有可能的分割点。适用于小型数据集, 结果精确但慢。
 - 。 'approx': 近似贪心算法,使用分位数和梯度直方图。
 - 'gpu_hist': 使用 GPU 加速的直方图算法 (需要支持 GPU 的 XGBoost 版本和硬件)。
 - **应用关联:** 脚本 2, 3, 5, 6 显式使用了 'hist'。

8.3.3 学习任务参数 (Learning Task Parameters)

- objective: 定义学习任务和相应的目标函数。
 - 'reg:squarederror':回归任务,使用平方损失。
 - o 'binary:logistic': 二元分类,输出概率,使用对数损失 (logloss)。这是你的脚本中最常用的。
 - 。 'binary:logitraw': 二元分类,输出逻辑转换前的分数。
 - o 'multi:softmax':多分类,输出类别标签。需要设置 num_class 参数。
 - 'multi:softprob': **多分类**, 输出每个类别的概率。需要设置 num class 参数。
- eval_metric: **在验证集上评估所使用的指标**。可以是一个字符串或一个列表。
 - **回归常用:** 'rmse' (均方根误差), 'mae' (平均绝对误差)。
 - 。 分类常用:
 - 'logloss' (对数损失,常用于概率输出)。
 - 'error'(分类错误率, 1 accuracy)。
 - 'auc' (ROC 曲线下面积)。
 - 'aucpr' (PR 曲线下面积)。
 - 'map' (平均精度均值)。
 - 。 **选择:** 应选择与 objective 相关且能反映最终业务目标的指标。eval_metric 对于**早停 (Early Stopping)** 至关重要。
 - **应用关联:** 脚本 2, 3, 5 都设置了 eval_metric=['logloss', 'auc', 'error'] 或类似组合, 以 便在训练过程中监控多个指标。
- seed / random state: 为了结果可复现。

8.3.4 Scikit-learn API 特定参数

- use_label_encoder (默认 True in older versions, False recommended now): 是否使用 Scikit-learn 的 LabelEncoder 自动编码目标变量 y。 **在新版本的 XGBoost 中,强烈建议设置为 False**,并在调用 fit 之前手动将类别标签(如果是字符串等)编码为 [0,n_classes-1] 范围内的整数。这可以避免一些潜在的问题和警告。
 - o **应用关联:** 脚本 2, 3, 5 都遵循了这个建议,设置了 use label encoder=False。

8.4 模型训练 (fit 方法)

XGBClassifier 的 fit 方法与 Scikit-learn 的标准接口一致,但增加了一些 XGBoost 特有的重要参数,特别是用于监控和控制训练过程的参数。

```
# 假设 X_train, y_train, X_val, y_val 已经准备好
# model 是已实例化的 XGBClassifier 对象
# 1. 定义评估集
eval_set = [(X_train, y_train), (X_val, y_val)] # 包含训练集和验证集
# 2. 调用 fit 方法进行训练
model.fit(
   X_train, y_train,
   eval_set=eval_set,
                            # 提供评估集
                             # 指定在评估集上监控的指标 (也可以在初始化时指定)
   eval_metric='logloss',
   early_stopping_rounds=50,
                            # 早停轮数
   verbose=100
                             # 每隔 100 轮打印一次评估结果
                           # 用于微调, 传入已加载的模型对象
   # xgb_model=None
   # sample_weight=None
                           # 样本权重
)
```

关键参数 (fit 方法):

- X, y: 训练数据和标签。
- eval_set: 一个包含(X, y)元组的列表,用于在训练过程中评估模型性能。通常至少包含一个验证集(X_val, y_val)。如果提供多个评估集,早停会基于**列表中的最后一个**评估集进行。
- eval_metric: 在 eval_set 上使用的评估指标。如果在 XGBClassifier 初始化时已指定,这里可以省略或覆盖。
- early_stopping_rounds: 早停机制。
 - o **作用**: 在训练过程中,持续监控模型在 eval_set(通常是验证集)上的 eval_metric。如果在指定的 early_stopping_rounds 轮内,该指标没有改善(例如,logloss 或 error 没有降低,auc 没有升高),则停止训练。
 - 好处:
 - 防止过拟合: 在模型开始在验证集上性能下降时及时停止。
 - **自动确定最佳** n_estimators: 无需手动猜测树的最佳数量,模型会在验证性能不再提升时停止。实际使用的树的数量存储在 model.best_ntree_limit (如果早停触发)。
 - o 注意: 需要提供 eval set 才能使用早停。
 - o **应用关联:** 脚本 2, 3 中使用了 early_stopping_rounds=30, 脚本 6 (微调) 中为了在特定数据上强制训练指定轮数而移除了早停 (early_stopping_rounds=None)。
- verbose: 控制训练过程中打印信息的频率。True 或一个整数 n 表示每 n 轮打印一次。False 表示不打印。
- xgb_model: 用于**继续训练 (微调)**。传入一个已加载的 XGBoost 模型对象 (.joblib 加载的或 XGBoost 原生 Booster 对象) ,训练将从该模型的状态继续。
 - o **应用关联**: 脚本 6 (微调) 正是利用了这个参数 finetune_model.fit(..., xgb_model=initial_model, ...) 来在初始基模型的基础上,使用 FN/FP 样本进行额外的训练。
- sample_weight: 为训练样本指定不同的权重。

预测方法与 Scikit-learn 标准一致:

- model.predict(X):返回预测的类别标签。对于 binary:logistic, 默认使用 0.5 的概率阈值。
- model.predict_proba(X):返回每个样本属于每个类别的概率。对于 binary:logistic,返回 (n_samples, 2)的数组,通常我们关心正类(第1列)的概率。

8.6 特征重要性 (feature_importances_)

评估哪些特征对模型的预测贡献最大,有助于理解模型、进行特征选择和业务洞察。XGBoost 提供了计算特征 重要性的方法。

- **访问:** model.feature_importances_ 属性返回一个 NumPy 数组,包含了每个特征的重要性得分(顺序与训练时的 X 列一致)。
- 重要性类型 (可以通过 importance_type 参数在 get_score 或 plot_importance 中指定,或者模型可能默认计算某种类型):
 - 'weight' (默认): 特征在所有树中被用作分裂节点的**次数**。
 - 'gain': 特征在所有树中作为分裂节点带来的**平均增益** (对损失函数的平均减少量)。通常比 'weight' 更能反映特征的实际重要性。
 - 'cover': 特征在所有树中作为分裂节点时所覆盖(处理)的平均样本数量。

• 可视化:

- o xgboost.plot_importance(model, max_num_features=...): XGBoost 自带的绘图函数。
- 。 手动绘图: 获取 feature_importances_数组,结合特征名称(如果可用),使用 Matplotlib 或 Seaborn 绘制条形图。
 - **应用关联**: 脚本 5 (元学习器) 中就包含了手动绘制特征重要性条形图的代码,展示了 "Base_Prob" 和 "FP_Expert_Prob" 这两个元特征对最终预测的贡献度。

```
import matplotlib.pyplot as plt
import pandas as pd
# (假设 model 是训练好的 XGBoost 模型, X train 是训练数据)
# feature names = X train.columns if isinstance(X train, pd.DataFrame) else
[f'f{i}' for i in range(X_train.shape[1])]
# # 获取重要性
# importances = model.feature_importances_
# importance df = pd.DataFrame({'Feature': feature names, 'Importance':
importances})
# importance_df = importance_df.sort_values(by='Importance', ascending=False)
# print("\nTop 10 Feature Importances:\n", importance df.head(10))
##绘制条形图
# plt.figure(figsize=(10, 6))
# plt.barh(importance_df['Feature'][:20], importance_df['Importance'][:20]) # 显示
前20个
# plt.xlabel("Importance (Gain/Weight/Cover)")
# plt.ylabel("Feature")
# plt.title("XGBoost Feature Importance")
# plt.gca().invert yaxis()
# plt.tight_layout()
```

```
# plt.show()

# # 使用 XGBoost 内置函数 (需要安装 matplotlib)

# from xgboost import plot_importance

# fig, ax = plt.subplots(figsize=(10, 8))

# plot_importance(model, ax=ax, max_num_features=20, importance_type='gain') # 指
定类型

# plt.show()
```

8.7 保存与加载 XGBoost 模型

有两种主要方式:

1. 使用 joblib:

- 与 Scikit-learn 标准方法一致,可以直接保存和加载 XGBClassifier 或 XGBRegressor 对象。
- 。 优点: 与 Scikit-learn 生态无缝集成。
- 缺点: 可能存在版本兼容性问题(依赖 Python, Scikit-learn, XGBoost 版本)。
- o **应用关联**: 这是你的所有脚本中使用的方法 (joblib.dump, joblib.load)。

2. 使用 XGBoost 原生方法:

- o model.save_model(filename): 将模型保存为 XGBoost 特有的二进制格式。
- loaded_model = xgb.Booster()(创建空 Booster)
- loaded model.load model(filename):加载模型。
- 。 或者直接加载 Scikit-learn wrapper:
 - loaded_xgb_classifier = xgb.XGBClassifier()
 - loaded_xgb_classifier.load_model(filename)
- 优点:通常具有更好的跨版本、跨语言兼容性。文件通常更小。
- 缺点: 需要分別保存和加载 Scikit-learn 包装器和内部的 Booster 模型,或者确保加载时使用正确的包装器。

对于纯 Python 和 Scikit-learn 的工作流,使用 joblib 通常更方便。但如果需要考虑跨语言部署或长期版本兼容性,原生格式可能更好。

XGBoost 是一个功能极其丰富的库,本章介绍了其核心概念、相对于 GBDT 的优势、关键参数 (特别是 Scikit-learn API 中的)、训练与预测方法、特征重要性评估以及模型持久化。掌握这些内容将使你能够更有效地使用 XGBoost 解决实际问题,并理解和调优类似天气预测脚本中的 XGBoost 模型。

```
```markdown
机器学习实战: 从 Scikit-learn 到 XGBoost (卷三)
...(接上文)

第九章: 特征工程: 化数据为神奇
"数据和特征决定了机器学习的上限,而模型和算法只是逼近这个上限而已。" — 这是机器学习领域一句广为流传的话,深刻揭示了特征工程的重要性。
```

特征工程(Feature Engineering)是指利用领域知识和数据分析技术,从原始数据中提取、构造或转换出对模型预测更有用、更具表达力的特征的过程。它往往是机器学习项目中最为耗时、最具创造性,同时也是对最终模型性能影响最大的环节之一。

一个好的特征应该能够清晰地将不同类别或预测目标区分开来。糟糕的特征或未经处理的原始数据可能 会包含噪声、冗余信息或不适合模型学习的格式,从而限制模型的表现。

本章将探讨一些常见的特征工程技术,帮助你理解如何"雕琢"数据,为你的模型提供更好的"养料"。

### ### 9.1 特征工程的目标

- \* \*\*提高模型精度: \*\* 提取更能反映数据内在模式的特征。
- \* \*\*提升模型鲁棒性:\*\*减少噪声和异常值的影响。
- \* \*\*简化模型: \*\* 使用更少的、但表达力更强的特征可能构建出更简单、更易解释的模型。
- \* \*\*降低计算成本: \*\* 合理的特征选择和降维可以减少模型训练和预测的时间。
- \* \*\*满足算法要求: \*\* 将数据转换为特定算法可以接受的格式 (例如,将类别特征转换为数值)。

#### ### 9.2 常见的特征工程技术

特征工程是一个非常广泛的主题,涉及的技术多种多样,很大程度上依赖于具体的数据和问题领域。这里我们介绍一些通用的技术。

### #### 9.2.1 数据清洗 (回顾与深化)

数据清洗是特征工程的基础,我们在第四章(Pandas)已经有所涉及,这里再次强调其重要性并补充一些内容:

- \* \*\*处理缺失值 (Handling Missing Values):\*\*
  - \* \*\*删除: \*\* 删除包含缺失值的样本(行)或特征(列)。简单但可能丢失信息。
  - \* \*\*填充 (Imputation):\*\*
- \* \*\*简单填充: \*\* 使用均值、中位数 (对异常值鲁棒)、众数 (用于类别特征) 进行填充 (`fillna`)。
- \* \*\*模型预测填充:\*\* 使用其他特征训练一个模型 (如 KNN, 回归) 来预测缺失值 (`sklearn.impute.KNNImputer`, `sklearn.impute.IterativeImputer`)。更复杂但可能更准确。
  - \* \*\*常数填充:\*\* 用一个特定的值(如 0, -1, 'Unknown')填充。
- \* \*\*创建指示器特征: \*\* 添加一个新的二元特征,指示原始值是否缺失。这可以保留"缺失"本身可能蕴含的信息。
- \* \*\*处理异常值 (Handling Outliers):\*\*
- \* \*\*识别:\*\* 通过统计方法 (如 Z-score, IQR 四分位距) 或可视化 (箱线图、散点图) 识别异常值。
  - \* \*\*处理:\*\*
    - \* \*\*删除: \*\* 如果确定是错误数据,可以删除。
    - \* \*\*转换:\*\* 使用对数转换、Box-Cox 转换等来减小异常值的影响。
- \* \*\*盖帽 (Capping / Winsorization): \*\* 将超出某个阈值 (如 1% 和 99% 分位数)的值替换为该阈值。
  - \* \*\*视为缺失值: \*\* 将异常值标记为缺失, 然后使用缺失值处理方法。
- \* \*\*处理重复值:\*\* `drop\_duplicates()` (已在 Pandas 章节介绍)。

### #### 9.2.2 特征转换 (Feature Transformation)

改变特征的分布或尺度,使其更适合模型学习。

- \* \*\*特征缩放 (Feature Scaling):\*\* (已在第五章介绍)
  - \* `StandardScaler`: 标准化 (均值0,标准差1)。
  - \* `MinMaxScaler`: 规范化 (缩放到 [0, 1] 或指定范围)。
  - \* `RobustScaler`: 使用中位数和 IQR, 对异常值鲁棒。
- \* \*\*何时需要: \*\* 对基于距离 (KNN, SVM) 或梯度下降 (线性模型,神经网络) 的算法非常 重要。对树模型 (决策树,随机森林,XGBoost) 通常不是必需的,但有时标准化也有助于某些树模型的性能。
- \* \*\*对数转换 (Log Transform):\*\* `np.log(x)` 或 `np.log1p(x)` (处理 x=0 的情况)。
- \* \*\*作用:\*\* 处理\*\*右偏 (positively skewed)\*\* 的数据分布,使其更接近正态分布。压缩较大值的范围,拉伸较小值的范围。
- \* \*\*适用场景: \*\* 当特征值的范围非常广,且大部分值集中在较小区域时(如收入、人口数量、网站点击次数等)。
- \* \*\*Box-Cox 转换 (`scipy.stats.boxcox`):\*\*
- \* 一种更通用的幂转换,可以自动寻找最佳的 lambda 参数,使得转换后的数据最接近正态分布。
  - \* \*\*限制:\*\* 要求数据必须是\*\*正数\*\*。
- \* \*\*分箱/离散化 (Binning / Discretization):\*\* 将连续特征划分为若干个区间 (箱) , 转换为类别特征。
  - \* \*\*方法:\*\*
    - \* \*\*等宽分箱 (`pd.cut`): \*\* 将数据范围分成宽度相等的箱。对异常值敏感。
    - \* \*\*等频分箱 (`pd.qcut`): \*\* 将数据分成包含大致相同数量样本的箱。
    - \* \*\*基于模型的离散化:\*\* 使用决策树等模型找到最佳分裂点进行分箱。
  - \* \*\*优点:\*\* 引入非线性,处理异常值,简化模型。
- \* \*\*缺点: \*\* 损失信息,需要确定合适的分箱数量和方式。分箱后通常需要进行 One-Hot 编码。
  - \* \*\*Scikit-learn:\*\* `sklearn.preprocessing.KBinsDiscretizer`.

#### 9.2.3 特征编码 (Feature Encoding)

将类别特征转换为数值表示。(已在第五章初步介绍)

- \* \*\*`OrdinalEncoder`:\*\* 用于\*\*有序\*\*类别特征。
- \* \*\*`OneHotEncoder`: \*\* 用于\*\*名义\*\*类别特征(无序)。最常用,但注意维度爆炸问题。
- \* \*\*`LabelEncoder`:\*\* \*\*仅用于目标变量 `y`\*\*。
- \* \*\*其他编码方法(适用于高基数类别特征):\*\*
  - \* \*\*频率编码 (Frequency Encoding): \*\* 用类别出现的频率 (或计数) 代替类别标签。
- \* \*\*目标编码 (Target Encoding / Mean Encoding):\*\* 用该类别对应的\*\*目标变量的平均值\*\*(或其他聚合值,如中位数)来代替类别标签。
  - \* \*\*优点: \*\* 可以有效地处理高基数类别, 并直接融入目标信息。
- \* \*\*缺点: \*\* 容易过拟合,需要小心处理(例如,使用交叉验证策略进行编码,加入平滑等)。
  - \* \*\*Hashing Trick: \*\* 使用哈希函数将类别映射到固定数量的特征上。可能存在哈希冲突。
- \*\*应用关联:\*\* 虽然你的脚本直接使用了 `.npy` 数值数据,但在生成这些数据之前的步骤中,很可能进行了特征编码,特别是如果原始数据包含非数值的类别信息。
- #### 9.2.4 特征构造 (Feature Creation / Construction)

这是特征工程中最具创造性的部分,通常需要结合领域知识和数据洞察来手动创建新的特征。

- \* \*\*交互特征 (Interaction Features): \*\* 将两个或多个特征进行组合(相乘、相除、相加、相减等),以捕捉它们之间的交互效应。
  - \* \*\*示例:\*\* `特征A \* 特征B`, `特征A / (特征B + epsilon)`。
  - \* \*\*Scikit-learn:\*\* `sklearn.preprocessing.PolynomialFeatures` 可以自动生成多

项式特征和交互特征。例如,对于特征 `[a, b]`, `PolynomialFeatures(degree=2, interaction\_only=False)` 会生成 `[1, a, b, a^2, a\*b, b^2]`; `interaction\_only=True` 则生成 `[1, a, b, a\*b]`。

- \* \*\*领域特定特征:\*\* 基于对问题领域的理解创建有意义的特征。
  - \* \*\*示例 (天气):\*\*
    - \* 计算温差 (最高温 最低温)。
    - \* 计算湿度变化率。
    - \* 结合风向和风速计算风矢量分量。
    - \* 根据历史数据计算滑动窗口内的平均降雨量、最大降雨量等。
- \* 时间特征:从日期时间戳中提取年、月、日、星期几、小时、是否节假日等。`dt`访问器 (`series.dt.xxx`) 在 Pandas 中非常有用。
- \* \*\*分解特征:\*\* 将复杂特征分解为更简单的部分。
  - \* \*\*示例: \*\* 将地址分解为国家、省份、城市。将 URL 分解为协议、域名、路径。
- \* \*\*聚合特征:\*\* 基于某些分组(例如,用户 ID、商品类别)计算统计特征(均值、标准差、计数、最大/最小值等)。`groupby()`+`agg()`/`transform()`在 Pandas 中是实现聚合特征的关键。

\*\*应用关联:\*\* 你的 `X\_flat\_features.npy` 中的 100 或 101 个特征,很可能就包含了通过上述方法构造出来的特征,而不仅仅是原始的传感器读数。例如,可能包含了不同时间点、不同地点的气象数据,或者经过计算得到的衍生指标。脚本 `3` 和 `4` 中将基模型的概率 (`base\_probs`) 作为额外特征加入到 FP 专家模型的输入中,这也是一种\*\*特征构造\*\*,利用了另一个模型的输出来增强当前模型的输入信息。

#### 9.2.5 特征选择 (Feature Selection)

当特征数量非常多时,选择一个相关的子集可以:

- \* \*\*简化模型,提高可解释性。\*\*
- \* \*\*减少训练时间。\*\*
- \* \*\*降低过拟合风险(减少维度灾难)。\*\*
- \* \*\*提高某些模型的性能(移除无关或冗余特征)。\*\*

### 特征选择方法主要分为三类:

- 1. \*\*过滤法 (Filter Methods):\*\*
- \* 独立于模型,根据特征本身的统计特性(如方差、与目标变量的相关性、互信息等)进行评分和筛选。
  - \* \*\*优点:\*\* 计算速度快,不易过拟合。
- \* \*\*缺点:\*\* 没有考虑特征之间的组合效应,可能移除掉单独看效果一般但组合起来有用的特征。
  - \* \*\*常用方法:\*\*
- \* \*\*移除低方差特征 (`sklearn.feature\_selection.VarianceThreshold`):\*\* 移除方差小于某个阈值的特征(例如,值基本不变的特征)。
  - \* \*\*单变量选择 (`sklearn.feature\_selection.SelectKBest`,
- `SelectPercentile`):\*\* 基于统计检验 (如 F 检验 `f\_classif`/`f\_regression`, 卡方检验 `chi2`, 互信息 `mutual\_info\_classif`/`mutual\_info\_regression`) 选择得分最高的前 K 个或前百分之 P 的特征。
- \* \*\*相关系数: \*\* 计算特征与目标变量的相关系数(如皮尔逊相关系数), 移除相关性低的特征。注意: 线性相关系数无法捕捉非线性关系。
- 2. \*\*包裹法 (Wrapper Methods):\*\*
- \* 将特征选择过程看作是一个搜索问题,使用一个特定的机器学习模型来评估不同特征子集的性能。
  - \* \*\*优点:\*\* 直接优化特定模型的性能,考虑特征组合。

- \* \*\*缺点: \*\* 计算成本非常高,因为需要反复训练模型。容易过拟合特征选择过程本身。
- \* \*\*常用方法:\*\*
  - \* \*\*递归特征消除 (Recursive Feature Elimination, RFE)

(`sklearn.feature\_selection.RFE`, `RFECV`):\*\* 反复训练模型,每次移除最不重要的一个或多个特征(基于模型的 `coef\_` 或 `feature\_importances\_`) ,直到达到所需的特征数量。 `RFECV` 通过交叉验证自动选择最佳特征数量。

- \* \*\*前向选择 (Forward Selection): \*\* 从空集开始,每次加入一个能最大程度提升模型性能的特征。
- \* \*\*后向消除 (Backward Elimination): \*\* 从全集开始,每次移除一个对模型性能影响最小的特征。
- 3. \*\*嵌入法 (Embedded Methods):\*\*
  - \* 将特征选择嵌入到模型训练过程中,模型在训练时自动进行特征选择。
  - \* \*\*优点: \*\* 计算效率通常介于过滤法和包裹法之间, 考虑了特征与模型的交互。
  - \* \*\*缺点: \*\* 特征选择依赖于所使用的模型。
  - \* \*\*常用方法:\*\*
- \* \*\*带 L1 正则化的模型 (Lasso, Logistic Regression with `penalty='l1'`, LinearSVC with `penalty='l1'`):\*\* L1 正则化倾向于产生稀疏权重,可以将不重要特征的系数压缩为 0,从而实现特征选择。可以通过 `SelectFromModel` 配合使用。
- \* \*\*基于树的模型 (Random Forest, Gradient Boosting, XGBoost): \*\* 这些模型在 训练后可以提供特征重要性评分 (`feature\_importances\_`)。可以根据重要性阈值选择特征,或者 直接使用模型本身进行预测(模型内部已经考虑了特征的重要性)。 `SelectFromModel` 也可以配合 树模型的重要性评分使用。
- \* \*\*应用关联: \*\* XGBoost 模型本身就具有一定的嵌入式特征选择能力,因为它在构建树的过程中会优先选择带来更大增益的特征。通过查看 `feature\_importances\_`,可以了解哪些特征对模型贡献更大。

```
`sklearn.feature_selection.SelectFromModel`:
这是一个元转换器 (meta-transformer),可以配合任何提供了`coef_`或
`feature_importances_`属性的评估器使用。它会根据指定的阈值选择特征。
```

#### python

from sklearn.feature\_selection import SelectFromModel
from sklearn.ensemble import RandomForestClassifier

- # (假设 model 是训练好的带 feature\_importances\_ 的模型, 如 RandomForest 或 XGBoost)
- # model = RandomForestClassifier(n\_estimators=100, random\_state=42).fit(X\_train,
  y train)
- # 使用 SelectFromModel 进行特征选择
- # threshold='median' 表示选择重要性高于中位数的特征
- # 也可以设置为浮点数阈值,或 '1.25\*mean'等

selector = SelectFromModel(model, threshold='median', prefit=True) # prefit=True 表示 model 已训练好

```
X_train_selected = selector.transform(X_train)
X_test_selected = selector.transform(X_test)

print("Original feature shape:", X_train.shape)
print("Selected feature shape:", X_train_selected.shape)
print("Selected features indices:", selector.get support(indices=True))
```

# 9.3 特征工程流程与实践

特征工程不是一个线性的过程,而是一个需要不断迭代和实验的循环:

- 1. 理解数据和业务: 这是最重要的第一步。深入了解每个特征的含义、数据的来源、业务目标是什么。
- 2. **初步探索 (EDA):** 使用统计和可视化方法检查数据分布、缺失情况、异常值、特征间关系、特征与目标关系。
- 3. **头脑风暴与假设:** 基于领域知识和 EDA 结果,思考哪些原始特征可能有用,哪些可以组合,哪些需要转换。提出关于特征与目标关系的假设。
- 4. **实施特征工程**: 应用数据清洗、转换、编码、构造等技术。
- 5. 特征选择 (可选但推荐): 如果特征维度很高,使用过滤法、包裹法或嵌入法选择一个子集。
- 6. 模型训练与评估: 使用处理后的特征训练模型,并使用合适的指标进行评估。
- 7. **分析结果与迭代:** 分析模型的错误,查看特征重要性,思考哪些特征工程步骤有效,哪些无效,是否有新的特征可以尝试。回到步骤 3 或 4 进行改进。

#### 工具:

- Pandas: 数据操作、清洗、转换、构造的基础。
- NumPy: 底层数值计算。
- Scikit-learn (preprocessing, impute, feature\_selection): 提供标准化的特征工程工具。
- Matplotlib / Seaborn: 数据可视化,用于 EDA 和结果分析。
- 领域知识: 无可替代!

特征工程的好坏直接决定了机器学习项目的成败。它需要耐心、细致、创造力以及对数据和业务的深刻理解。 不要害怕尝试不同的方法,并始终以提升模型在目标任务上的最终性能为导向。

```
```markdown
```

机器学习实战: 从 Scikit-learn 到 XGBoost (卷四)

...(接上文)

- - -

第十章: 集成学习进阶: Stacking 与 Blending

在前面的章节中,我们已经接触并深入探讨了两种强大的集成学习方法: Bagging (以随机森林为代表) 和 Boosting (以 GBDT 和 XGBoost 为代表)。这些方法通过结合多个学习器的预测来获得比单一学习器更好的性能。本章我们将探讨一种更高级、也可能更强大的集成技术—**Stacking (堆叠泛化)**,以及其简化版 **Blending**。

Stacking 的思想是训练多个不同的**基学习器 (Base Learners)**, 然后将它们的预测结果作为**新的特征 (Meta-Features)**, 再用这些新特征去训练一个**元学习器 (Meta-Learner)** 来做出最终的预测。你的天气预测脚本(特别是 1-5)实际上实现了一种 Stacking 的变体,理解 Stacking 的原理对于弄清这些脚本的工作流程至关重要。

10.1 集成学习回顾

在深入 Stacking 之前,简单回顾一下其他集成方法的核心思想:

- * **Bagging (Bootstrap Aggregating):**
 - * **核心思想:** 通过自助采样 (Bootstrap Sampling) 创建多个训练数据的子集, 在每个

子集上独立训练同一种类型的基学习器 (通常是决策树) , 最后通过投票或平均来合并预测结果。

- * **目的:** 主要用于**降低模型的方差 (Variance)**, 提高模型的稳定性和泛化能力。
- * **代表:** 随机森林 (Random Forest)。
- * **Boosting:**
- * **核心思想: ** 串行地训练一系列弱学习器,每个新的学习器都重点关注前一轮学习器预测错误的样本,最终将所有学习器加权组合。
 - * **目的:** 主要用于**降低模型的偏差 (Bias)**, 提升模型的准确率。
 - * **代表:** AdaBoost, GBDT, XGBoost, LightGBM。

10.2 Stacking (堆叠泛化)

Stacking 的目标是学习如何**智能地组合**多个不同基学习器的预测,而不是简单地投票或平均。它试图利用不同模型可能从数据中学到的不同模式或犯的不同错误。

核心架构 (多层结构):

- 1. **Level 0 (基学习器层):**
- * 包含多个**不同的**基学习器。这些学习器可以是不同类型的算法(如逻辑回归、SVM、KNN、随机森林、XGBoost),也可以是相同算法但使用不同超参数的多个实例。
- * **关键:** 基学习器之间最好具有**多样性 (Diversity)**。如果所有基学习器都非常相似,或者犯的错误都一样,那么 Stacking 的效果可能不会比最好的单个基学习器好多少。多样性使得元学习器有机会学习如何取长补短。
- 2. **Level 1 (元特征层):**
- * 这一层的输入是 Level 0 基学习器对**原始数据**的预测结果。这些预测结果被称为**元特征 (Meta-Features)**。
- * 对于分类问题,元特征通常是基学习器预测的**类别概率**(`predict_proba`的输出)。 对于回归问题,则是预测的数值。
- 3. **Level 2 (元学习器层):**
- * 使用 Level 1 生成的元特征作为**输入**,原始数据的**真实标签**作为**输出**,训练一个新的学习器,称为**元学习器 (Meta-Learner)**。
- * 元学习器的任务是学习如何根据不同基学习器的预测(元特征)来做出最终的、更准确的预测。
- * **选择:** 元学习器通常选择**相对简单**的模型,如**逻辑回归、岭回归、线性 SVM 或层数较浅的 XGBoost/LightGBM**,以防止在元特征上过拟合。

Stacking 训练流程 (防止数据泄露的关键):

为了让元学习器能够有效地学习如何组合基模型的预测,生成 Level 1 的元特征 (特别是用于**训练元学习器**的元特征) 时必须非常小心,以**防止数据泄露**。数据泄露是指在训练过程中,模型(在这里是元学习器)不应接触到它在预测时无法获得的信息(例如,基模型在预测某个样本时使用了该样本的真实标签)。

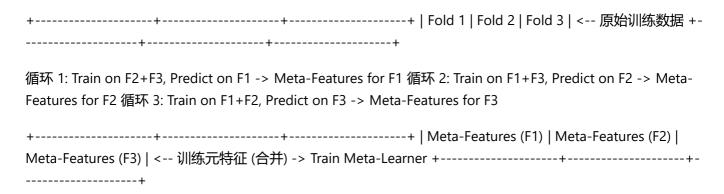
标准的 Stacking 流程通常使用 **K-Fold 交叉验证** 来生成用于训练元学习器的元特征:

- 1. **划分数据: ** 将原始训练数据划分为 K 个互不重叠的折 (Folds)。
- 2. **循环 K 次:**
- * 在第 `k` 次循环中,将第 `k` 折作为**验证集 (Hold-out Fold)**,其余 K-1 折作为 **训练集**。
 - * **训练 Level 0 模型:** 使用这 K-1 折数据训练所有的 Level 0 基学习器。
- * **生成验证集的元特征:** 使用刚刚训练好的 Level 0 基学习器对**第 `k` 折验证集** 进行预测,得到这一折样本的元特征。
- 3. **合并元特征:** 将 K 次循环中生成的每一折的元特征按原始顺序合并起来,形成完整的**训练元特征矩阵 (Meta-Features for Meta-Learner Training)**。这个矩阵的行数与原始训练集相同,列数等于所有基学习器预测输出的总维度(例如,如果是 N 个二分类基学习器都输出概率,则

有 N*2 或 N 列元特征)。

4. **训练 Level 0 模型 (全量):** 使用**全部**原始训练数据重新训练所有的 Level 0 基学习器。这一步是为了得到最终用于预测测试集的基模型。

- 5. **生成测试集的元特征:** 使用在第 4 步中训练好的 Level 0 基学习器对**原始测试集**进行预测,得到**测试元特征矩阵 (Meta-Features for Final Prediction)**。
- 6. **训练 Level 2 元学习器:** 使用第 3 步生成的**训练元特征矩阵**和原始训练数据的**真实标签**来训练元学习器。
- 7. **最终预测:** 将第 5 步生成的**测试元特征矩阵**输入到训练好的元学习器中,得到对测试集的最终预测结果。
- **图示 (简化版 K=3):**



Train Base Models on F1+F2+F3 (Full Train Data) | V Predict on Test Data -> Test Meta-Features -> Predict using Trained Meta-Learner -> Final Predictions

Stacking 的优点:

- * **潜力巨大:** 通过学习组合不同模型的优势,有可能达到比任何单一模型或简单集成(如投票、平均)更高的性能。
- * **灵活性: ** 可以组合任意类型的基学习器。
- **Stacking 的缺点:**
- * **实现复杂: ** K-Fold 流程相对复杂, 容易出错。
- * **训练时间长: ** 需要多次训练基学习器和一次元学习器训练。
- * **调优困难: ** 不仅需要调整每个基学习器的超参数,还需要选择合适的元学习器及其超参数。
- * **过拟合风险: ** 如果元特征生成不当(数据泄露)或元学习器过于复杂,容易过拟合。
- * **信息损失:** 从基学习器的复杂决策过程压缩到几个预测值(元特征),可能损失一些信息。
- **连接用户脚本 (1-5) 一种 Stacking 变体:**

你的脚本实现了一种 Stacking 的变体,其特点是引入了**针对特定错误类型的专家模型**,并且在生成元特征的方式上有所简化(可能带来风险)。

- 1. **Level 0 (基模型):**
 - * `xgboost_default_full_model.joblib` (由 `xgboost1.py` 假定生成)。这是基础。
- * `1_generate_base_predictions.py`使用它在**训练集**上预测,识别出TP/TN/FP/FN,并保存训练集和测试集的**预测概率**(`base probs train.npy`,

`base probs test.npy`)。这些概率是后续步骤的关键输入。

- 2. **Level 0.5 (专家模型): ** 这可以看作 Stacking 结构的一个扩展或修改。
- * `2_train_fn_expert.py`: 使用训练集中的 **FN** 和 **TN** 样本及其**原始特征** 训练一个 FN 专家模型。目标是识别那些被 Level 0 漏掉的正类样本。
- * `3_train_fp_expert.py`: 使用训练集中的 **FP** 和 **TN** 样本及其**原始特征 + Level 0 的训练集预测概率 (`base_probs_train`)** 训练一个 FP 专家模型。目标是识别那些被Level 0 误报的正类样本。加入基模型概率作为特征,让 FP 专家可以学习"哪些看起来概率高的预测实际上是错的"。
 - * **关键点:**
 - * 专家模型是针对基模型的**特定弱点**训练的。
 - * 它们使用的训练数据是基模型在**训练集**上犯错的子集。
 - * FP 专家模型利用了基模型的输出作为**额外特征**。
- 3. **Level 1 (元特征生成 脚本 `4_generate_meta_features.py`):**
 - * 这一步生成**测试集**的元特征 `X meta`。
 - * **输入:**
 - * 测试集的原始特征 (`X_test`)。
 - * 基模型在测试集上的预测概率 (`base_probs_test`)。
 - * 训练好的 **FP 专家模型**。
 - * **过程:** 对测试集中的每个样本:
 - * 获取基模型的预测概率 (来自 `base_probs_test`)。
 - * 准备 FP 专家的输入: `[测试样本原始特征,基模型预测概率]`。
 - * 使用 FP 专家模型预测该样本"是 FP 的概率"

(`fp_expert_model.predict_proba(...)[;, 1]`).

- * **輸出 `X_meta`:** 形状为 `(n_test_samples, 2)`, 包含两列:
 - * `X_meta[:, 0]`: 基模型在测试集上的预测概率 (`Base_Prob`)。
 - * `X_meta[:, 1]`: FP 专家模型在测试集上预测"是 FP"的概率
- (`FP_Expert_Prob`).
 - * **注意:**
 - * 这里**没有使用** FN 专家模型的预测结果作为元特征,这是一种设计选择。
- * 这里生成的元特征是直接用于**训练和评估**元学习器的(因为 `y_meta` 是测试集的真实标签)。这与标准的 K-Fold Stacking 不同,标准的做法是用 K-Fold 在训练集上生成元特征来训练元学习器。这种简化做法(直接在测试集元特征上训练/评估元学习器)可能会导致对模型最终性能的**评估略微乐观**,因为它没有模拟严格的"先训练后预测"流程。
- 4. **Level 2 (元学习器训练与评估 脚本 `5_train_evaluate_meta_learner.py`):**
 - * **输入:** `X_meta` (测试集的元特征), `y_meta` (测试集的真实标签)。
- * **过程:** 选择一个元学习器 (脚本中是 Logistic Regression 或 XGBoost) , 使用 `X_meta` 和 `y_meta` 进行训练。
 - * **评估:** 在**相同**的 `X meta` 和 `y meta` 上进行评估。
- * **输出: ** 训练好的元学习器模型 (`meta_learner_model.joblib`) 和最终的性能评估结果。
- **总结这种变体:**
- * 它不是使用多个独立的基学习器,而是围绕一个基学习器,并训练专家模型来"修正"它的特定错误。
- * FP 专家模型利用了基模型的概率输出,这是一种特征工程。
- * 最终的元特征融合了基模型的直接预测和 FP 专家对"误报可能性"的判断。
- * 在元特征生成和元学习器训练/评估的流程上,相比标准 K-Fold Stacking 有所简化,需要注意 其对最终性能评估可能带来的乐观偏差。

10.3 Blending

Blending 是 Stacking 的一种简化方法,它避免了 K-Fold 交叉验证的复杂性。

Blending 训练流程:

- 1. **划分数据: ** 将原始训练数据划分为两部分: 一个**训练子集 (Training Subset) ** (例如 70-80%) 和一个**验证/留出集 (Validation/Holdout Set) ** (例如 20-30%) 。测试集保持不 变。
- 2. **训练 Level 0 模型:** 使用**训练子集**训练所有的 Level 0 基学习器。
- 3. **生成验证集的元特征:** 使用在第 2 步中训练好的 Level 0 基学习器对**验证/留出集**进行预测,得到**训练元特征矩阵 (Meta-Features for Meta-Learner Training)**。
- 4. **生成测试集的元特征:** 使用**相同**的 Level 0 基学习器对**原始测试集**进行预测, 得到**测试元特征矩阵 (Meta-Features for Final Prediction)**。
- 5. **训练 Level 2 元学习器:** 使用第 3 步生成的**训练元特征矩阵**和**验证/留出集**的** 真实标签**来训练元学习器。
- 6. **最终预测: ** 将第 4 步生成的**测试元特征矩阵**输入到训练好的元学习器中,得到对测试集的最终预测结果。
- **Blending vs. K-Fold Stacking:**
- * **优点:**
 - * **实现简单:** 比 K-Fold 流程简单得多。
- * **数据泄露风险低:** 元学习器的训练数据(来自验证集)与基学习器的训练数据(来自训练子集)是完全分开的,有效避免了数据泄露。
- * **缺点:**
 - * **数据使用效率较低:**
 - * 基学习器只使用了部分训练数据(训练子集)。
 - * 元学习器也只使用了部分数据(验证/留出集)。
- * **性能可能稍差:** 由于使用的数据量较少,最终模型的性能可能略低于精心实现的 K-Fold Stacking。
 - * **对划分敏感: ** 最终性能可能受最初训练集/验证集划分比例和随机状态的影响。

如果计算资源有限或追求快速实现,Blending 是一个不错的替代方案。如果追求极致性能且计算资源允许,K-Fold Stacking 通常是更好的选择。

10.4 Stacking/Blending 实践考量

- * **基学习器选择: ** 多样性是关键。尝试不同类型的算法(线性模型、树模型、核方法、近邻方法等)。即使是同一种算法,使用不同的超参数或不同的特征子集也能增加多样性。
- * **元特征选择:**
- * 通常使用 `predict_proba` 的输出作为元特征。对于二分类,可以选择只使用正类的概率,或者使用所有类的概率。
 - * 可以考虑加入原始特征到元学习器的输入中(虽然会增加维度和复杂性)。
 - * 可以像你的脚本一样,加入"专家模型"的预测。
- * **元学习器选择: ** 简单模型 (如 Logistic Regression, Ridge) 通常足够, 且不易过拟合。如果元特征数量较多或关系复杂, 也可以尝试 LightGBM 或浅层 XGBoost。
- * **特征工程: ** 特征工程对于基学习器和元学习器都可能很重要。
- * **调优:** Stacking 的调优很复杂。可以分步进行:先独立调优每个基学习器,然后固定基学习器,再调优元学习器。
- * **计算成本:** Stacking(尤其是 K-Fold)的计算成本很高,需要有耐心和足够的计算资源。

Stacking 和 Blending 是强大的集成技术,能够通过学习组合多个模型的预测来突破单一模型的性能瓶颈。理解其原理、流程(特别是防止数据泄露的方法)以及与你的脚本实现的关联,将有助于你更好地应用和改进这些高级技术。

_ _ _

第十一章: 模型调优与选择

我们已经学习了多种机器学习模型,并了解了如何评估它们的性能。然而,几乎所有模型都包含一些**超参数(Hyperparameters)**,这些参数不是模型从数据中学习得到的,而是需要我们在训练前设置的(例如,逻辑回归的 `C`,SVM 的 `C` 和 `gamma`,决策树的 `max_depth`,随机森林和 XGBoost 的 `n_estimators`,`learning_rate`,`max_depth`等)。

超参数的选择对模型的性能有着至关重要的影响。选择不当可能导致模型欠拟合或过拟合。**模型调优 (Model Tuning)** 或 **超参数优化 (Hyperparameter Optimization)** 的目标就是找到一组能 够使模型在**未见过的数据上**表现最佳的超参数组合。**模型选择 (Model Selection)** 则是在 多个不同的模型算法 (或同一算法的不同超参数组合) 之间进行比较,选出最终最优的模型。

本章将介绍 Scikit-learn 中进行模型调优和选择的常用技术,主要是**交叉验证 (Cross-Validation)** 和基于交叉验证的**网格搜索 (Grid Search)** 与**随机搜索 (Random Search)**。

11.1 交叉验证 (Cross-Validation)

我们在第五章提到了使用**验证集(Validation Set)** 来进行模型评估和超参数调优,以避免在测试集上过拟合。然而,将数据固定地划分为训练集和验证集存在一个问题:模型的性能评估结果可能**高度依赖于这次具体的划分**。如果某次划分恰好使得验证集比较"容易"或比较"困难",那么我们得到的评估结果就可能无法代表模型在其他未见过数据上的真实泛化能力。

交叉验证 (Cross-Validation, CV) 提供了一种更稳健、更可靠的模型评估和选择方法。它通过多次不同的数据划分和模型训练/评估,来获得对模型性能更鲁棒的估计。

K- 折交叉验证 (K-Fold Cross-Validation):

这是最常用的交叉验证技术。

- 1. **划分:** 将原始**训练数据集**随机划分为 K 个大小相似、互不重叠的子集(称为"折", Fold)。
- 2. **循环 K 次:**
 - * 在第 `k` 次迭代中:
 - * 将第 `k` 折作为**验证集 (Validation Fold)**。
 - * 将其余的 K-1 折合并作为**训练集 (Training Fold)**。
 - * 使用这个训练集训练模型。
 - * 在第 `k` 折验证集上评估模型性能,记录评估指标(例如,准确率、AUC 等)。
- 3. **计算平均性能:** 将 K 次迭代得到的评估指标进行平均 (有时也会看标准差) , 得到最终的 交叉验证性能估计。

图示 (K=5):

Data: | Fold 1 | Fold 2 | Fold 3 | Fold 4 | Fold 5 |

Iter 1: Train on [F2,F3,F4,F5], Validate on F1 -> Score 1 Iter 2: Train on [F1,F3,F4,F5], Validate on F2 -> Score 2 Iter 3: Train on [F1,F2,F4,F5], Validate on F3 -> Score 3 Iter 4: Train on [F1,F2,F3,F5], Validate on F4 -> Score 4 Iter 5: Train on [F1,F2,F3,F4], Validate on F5 -> Score 5

Final CV Score = Average(Score 1, Score 2, Score 3, Score 4, Score 5)

K 的选择:

- * 常用的 K 值为 5 或 10。
- * **较大的 K:**
 - * 优点:训练集更大(接近原始训练集),验证集更小,性能估计的偏差较低。
- * 缺点:训练次数更多,计算成本更高。不同迭代的训练集重叠度高,可能导致性能估计的方差较大。
- * **较小的 K:**
 - * 优点:计算成本较低。
 - * 缺点:训练集相对较小,验证集较大,性能估计的偏差可能较高,方差较低。
- * **留一交叉验证 (Leave-One-Out Cross-Validation, LOOCV): ** K 等于样本数 N。每次只留一个样本作验证,其余 N-1 个作训练。计算成本极高,通常只用于小数据集。
- **分层 K 折交叉验证 (`StratifiedKFold`):**

对于**分类问题**,特别是类别不平衡时,应该使用**分层 K 折交叉验证**。它确保在每次划分时,训练集和验证集中的类别比例都与原始数据集大致相同。这可以防止某个折中缺少少数类样本的情况。

Scikit-learn 中的交叉验证工具:

- * **交叉验证分割器 (Cross-validation iterators) (`sklearn.model_selection`):**
- * `KFold(n_splits=5, shuffle=False, random_state=None)`: 标准 K 折。 `shuffle=True` 可以在划分前打乱数据。
- * `StratifiedKFold(n_splits=5, shuffle=False, random_state=None)`: **分层 K 折,推荐用于分类问题**。
 - * `LeaveOneOut()`: 留一法。
- * `ShuffleSplit(n_splits=10, test_size=0.2, random_state=None)`: 每次随机划分指定比例作为测试集,可以有重叠。
 - * `StratifiedShuffleSplit`:分层的随机划分。
- * **便捷函数 `cross_val_score`:**
 - * 用法: `scores = cross_val_score(estimator, X, y, cv=5, scoring=None)`
 - * 作用:快速计算模型在 K 折交叉验证下的得分。
 - * `estimator`: Scikit-learn 模型对象。
 - * `X`, `y`: 完整的训练数据和标签。
- * `cv`: 指定交叉验证策略。可以是整数 K (默认使用 `StratifiedKFold` for 分类,
- `KFold` for 回归) , 也可以是交叉验证分割器对象 (如 `StratifiedKFold(n_splits=5)`) 。
- * `scoring`: 指定评估指标。可以是字符串(如 `'accuracy'`, `'roc_auc'`, `'f1'`, `'neg_mean_squared_error'`等,注意回归损失函数通常用负值),也可以是自定义评分函数。如果不指定,使用 `estimator.score()`的默认指标。
 - * 返回: 一个包含 K 次评估得分的 NumPy 数组。

```python

from sklearn.model selection import cross val score, StratifiedKFold

from sklearn.linear model import LogisticRegression

from sklearn.preprocessing import StandardScaler

from sklearn.pipeline import Pipeline # 用于串联步骤

- # (假设 X, y 是完整的训练数据)
- # X, y = make_classification(n_samples=1000, n_features=20, random_state=42)

```
# 创建包含缩放和模型的管道 (Pipeline)
# Pipeline 可以将多个处理步骤链接起来,像一个单一的 Estimator
# 这在交叉验证中尤其重要,确保每次都在当前折的训练数据上 fit scaler
pipe = Pipeline([
    ('scaler', StandardScaler()),
    ('logistic', LogisticRegression(solver='liblinear', random_state=42))
1)
# 定义分层 K 折交叉验证策略
cv_strategy = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
# 使用 cross_val_score 计算 AUC 分数
# scoring='roc_auc' 会自动调用 predict_proba
auc_scores = cross_val_score(pipe, X, y, cv=cv_strategy, scoring='roc_auc',
n jobs=-1
print(f"\nCross-Validation AUC Scores: {auc_scores}")
print(f"Mean AUC: {auc scores.mean():.4f}")
print(f"Std Dev AUC: {auc_scores.std():.4f}")
# 使用 cross_val_score 计算 Accuracy 分数
acc_scores = cross_val_score(pipe, X, y, cv=cv_strategy, scoring='accuracy',
n_{jobs}=-1)
print(f"\nCross-Validation Accuracy Scores: {acc_scores}")
print(f"Mean Accuracy: {acc_scores.mean():.4f}")
```

交叉验证的重要性: 它是进行可靠模型评估和后续超参数优化的基础。相比单一的训练/验证集划分,交叉验证的结果更能代表模型的平均泛化能力。

11.2 网格搜索 (GridSearchCV)

知道了如何通过交叉验证来评估一个**给定超参数组合**的模型性能后,下一步就是系统性地寻找**最佳的超参数组合**。网格搜索是最常用、最直接的方法之一。

核心思想:

- 1. **定义参数网格:** 为你想要调优的每个超参数,定义一个候选值的列表或范围。这些候选值构成了一个多维的"网格"。
- 2. 遍历所有组合: 网格搜索会穷举参数网格中所有可能的超参数组合。
- 3. **交叉验证评估:** 对于**每一个**超参数组合,使用 K 折交叉验证来评估模型的性能(使用指定的 scoring 指标)。
- 4. 选择最佳组合: 选择在交叉验证中平均性能表现最佳的那一组超参数组合。
- 5. 重新训练: (可选但推荐) 使用找到的最佳超参数组合, 在完整的训练数据集上重新训练最终的模型。

sklearn.model_selection.GridSearchCV:

```
from sklearn.model_selection import GridSearchCV
from sklearn.svm import SVC
from sklearn.pipeline import Pipeline
from sklearn.preprocessing import StandardScaler
```

```
# (假设 X, y 是完整的训练数据)
# X, y = make_classification(n_samples=500, n_features=10, random_state=42) # 使用
稍小的数据集演示
# 1. 创建包含预处理和模型的管道
pipe_svc = Pipeline([
   ('scaler', StandardScaler()),
   ('svc', SVC(probability=True, random_state=42)) # SVC 对参数敏感, 适合演示
])
# 2. 定义参数网格 (字典)
# 键是参数名 (对于 Pipeline, 格式是 '步骤名_参数名')
# 值是候选值列表
param_grid = {
   'svc__C': [0.1, 1, 10, 100],
                                          # 正则化参数 C
   'svc__gamma': [0.001, 0.01, 0.1, 1, 'scale', 'auto'], # RBF 核系数 gamma
   'svc__kernel': ['rbf', 'linear']
                                   # 尝试 RBF 和线性核
}
# 3. 实例化 GridSearchCV
# refit=True (默认) 表示在找到最佳参数后, 在整个 X, y 上重新训练最佳模型
# verbose 控制输出信息的详细程度
# n_jobs=-1 使用所有 CPU 核心并行计算
grid_search = GridSearchCV(
   estimator=pipe_svc,
   param_grid=param_grid,
                      # 使用 AUC 作为评估指标
   scoring='roc_auc',
   cv=StratifiedKFold(n_splits=3, shuffle=True, random_state=42), # 3 折 CV 演示
   refit=True,
   verbose=2,
   n jobs=-1
)
# 4. 执行搜索 (这步会比较耗时)
print("Starting Grid Search...")
grid_search.fit(X, y)
print("Grid Search finished.")
# 5. 查看最佳结果
print("\nBest parameters found:", grid_search.best_params_)
print(f"Best cross-validation AUC score: {grid search.best score :.4f}")
# 6. 获取最佳模型
# grid search.best estimator 是在整个 X, y 上用最佳参数训练好的模型
best_model = grid_search.best_estimator_
# 7. (可选) 在测试集上评估最佳模型 (如果之前划分了测试集)
# X_test_scaled = best_model.named_steps['scaler'].transform(X_test) # 需要先缩放
# y_pred_best = best_model.predict(X_test)
# y_pred_proba_best = best_model.predict_proba(X_test)[:, 1]
# print("\nPerformance of best model on test set:")
# print("Test Accuracy:", accuracy_score(y_test, y_pred_best))
# print("Test AUC:", roc_auc_score(y_test, y_pred_proba_best))
```

```
# 8. (可选) 查看所有结果
# cv_results_df = pd.DataFrame(grid_search.cv_results_)
# print("\nCV Results Summary (Top 5):\n",
cv_results_df.sort_values(by='rank_test_score').head())
```

网格搜索的优点:

- 系统性: 能够保证找到指定网格内的最佳参数组合。
- 易于理解和实现。

网格搜索的缺点:

• **计算成本高**: 随着超参数数量和候选值数量的增加,需要评估的组合数量会呈指数级增长("维度灾难")。对于有很多超参数的模型(如 XGBoost),穷举搜索可能不可行。

11.3 随机搜索 (RandomizedSearchCV)

当超参数空间很大时,随机搜索通常是比网格搜索更高效的替代方案。

核心思想:

- 1. **定义参数分布**: 对于每个想要调优的超参数,**不是**定义一个固定的候选值列表,而是定义一个**概率分布** (例如,均匀分布、对数均匀分布、离散值列表)。
- 2. **随机采样: 不是**遍历所有组合,而是从指定的分布中**随机采样**指定数量 (n_iter) 的超参数组合。
- 3. **交叉验证评估:** 对每个随机采样的组合,使用 K 折交叉验证评估模型性能。
- 4. 选择最佳组合: 选择在交叉验证中表现最佳的那个随机采样的组合。
- 5. **重新训练:** (可选) 使用最佳组合在完整训练集上重新训练模型。

直观理解: 在高维超参数空间中,可能只有少数几个超参数对模型性能影响最大。随机搜索通过在整个空间中广泛撒点,更有可能在有限的尝试次数内找到那些关键参数的较好值,而网格搜索可能会花费大量时间在非关键参数的细微调整上。

sklearn.model_selection.RandomizedSearchCV:

用法与 GridSearchCV 非常相似,主要区别在于 param_grid 变为 param_distributions,并且需要指定 n iter。

```
# 3. 实例化 RandomizedSearchCV
# n iter 控制采样的组合数量
random_search = RandomizedSearchCV(
   estimator=pipe_svc,
   param distributions=param distributions,
                              # 尝试 50 种随机组合 (远少于网格搜索)
   n iter=50,
   scoring='roc_auc',
   cv=StratifiedKFold(n splits=3, shuffle=True, random state=42),
   refit=True,
   verbose=1,
   n_{jobs=-1},
                     # 控制采样过程的随机性
   random_state=42
)
# 4. 执行搜索
print("\nStarting Randomized Search...")
random_search.fit(X, y)
print("Randomized Search finished.")
# 5. 查看最佳结果
print("\nBest parameters found (Randomized):", random_search.best_params_)
print(f"Best cross-validation AUC score (Randomized):
{random_search.best_score_:.4f}")
# 6. 获取最佳模型
best_model_random = random_search.best_estimator_
```

随机搜索的优点:

- **计算效率高**: 在给定计算预算(尝试次数 n_iter)下,通常能比网格搜索更快地找到一个相当好(甚至可能更好)的超参数组合,尤其是在高维空间中。
- 灵活性: 可以为连续参数指定更灵活的分布。

随机搜索的缺点:

• 不保证找到全局最优: 由于是随机采样,可能错过网格中的最佳点。但实践中通常能找到足够好的解。

选择网格搜索还是随机搜索?

- 参数空间较小: 网格搜索可以系统性地探索所有组合。
- 参数空间很大,计算资源有限: 随机搜索通常是更明智的选择。可以先用随机搜索缩小范围,再对有希望的区域进行更精细的网格搜索。

更高级的超参数优化库:

除了 Scikit-learn 自带的工具,还有一些更高级的库专门用于超参数优化,它们使用更智能的搜索策略(如贝叶斯优化、遗传算法等),可能比随机搜索更高效:

- Hyperopt
- Optuna
- Scikit-optimize (skopt)

11.4 模型选择流程总结

- 一个完整的模型选择和调优流程通常如下:
 - 1. 准备数据: 清洗、预处理、特征工程。
 - 2. **划分数据:** 将数据划分为**训练集和测试集** (Test Set)。测试集放在一边,**绝不**用于训练或调优。
 - 3. **选择候选模型和参数范围:** 确定要尝试的模型算法 (例如 Logistic Regression, Random Forest, XGBoost) 以及每个模型要调优的超参数及其搜索范围 (列表或分布)。
 - 4. 选择交叉验证策略: 通常是 StratifiedKFold (用于分类)。
 - 5. 选择评估指标: 例如 'roc_auc', 'f1', 'neg_mean_squared_error'。
 - 6. **执行超参数搜索:** 对每个候选模型,使用 GridSearchCV 或 RandomizedSearchCV 在**训练集**上进行超参数搜索,使用交叉验证进行评估。
 - 7. **比较最佳模型:** 比较不同算法在交叉验证中得到的**最佳性能** (best score)。
 - 8. **选择最终模型:** 选择交叉验证性能最好的那个模型 (即 grid_search.best_estimator_ 或 random_search.best_estimator_) 。这个模型通常是已经在完整训练集上用最佳参数重新训练过的 (refit=True)。
 - 9. 最终评估: 使用测试集对最终选定的模型进行一次性评估, 得到其泛化性能的最终估计。

这个流程确保了超参数调优和模型选择过程是在模拟未见数据的场景下进行的(通过交叉验证),并且最终的性能评估是在完全独立的数据上完成的。

掌握交叉验证、网格搜索和随机搜索是进行有效模型调优和选择的关键技能,它们能帮助你从众多可能的模型和参数组合中找到最优解,从而显著提升机器学习应用的实际效果。

- ```markdown
- # 机器学习实战: 从 Scikit-learn 到 XGBoost (卷五)
- ...(接上文)

- - -

第十二章: 处理不平衡数据

在许多现实世界的分类问题中,不同类别之间的样本数量可能存在巨大差异。这种情况被称为**类别不平衡 (Class Imbalance)**。例如:

- * **欺诈检测: ** 欺诈交易远远少于正常交易。
- * **罕见病诊断: ** 患有某种罕见病的病人数量远少于健康人。
- * **设备故障预测: ** 设备发生故障的情况远少于正常运行的情况。
- * **广告点击率预测 (CTR):** 用户点击广告的次数远少于看到广告但未点击的次数。
- * **天气预报(极端事件): ** 发生极端天气事件(如强降水、台风)的次数远少于正常天气。

类别不平衡会对机器学习模型的训练和评估带来严峻挑战。

挑战:

- 1. **模型偏向多数类: ** 大多数标准分类算法(如逻辑回归、SVM、决策树等)的目标是最大化总体准确率。在类别不平衡的情况下,模型很容易通过简单地将所有样本预测为多数类来获得很高的准确率,而完全忽略了少数类。这样的模型在实际应用中是无用的。
- 2. **评估指标失效: ** 准确率 (Accuracy) 不再是衡量模型性能的可靠指标。我们需要依赖其他指标,如精确率 (Precision)、召回率 (Recall/POD)、F1 分数、AUC、PR-AUC、CSI 等,这些指标

更能反映模型在少数类上的表现。

3. **难以学习少数类模式:** 由于少数类样本数量稀少,模型可能没有足够的信息来学习其内在模式,容易将少数类样本误判为多数类。

本章将介绍处理类别不平衡问题的常用策略。

12.1 认识不平衡: 数据探索

在开始处理不平衡问题之前,首先要确认数据确实存在不平衡,并了解其严重程度。

* **计算类别频率: ** 使用 Pandas 的 `value_counts()` 或 NumPy 的 `bincount()` 查看每个类别的样本数量和比例。

```
```python
import pandas as pd
import numpy as np
假设 y 是包含标签的 Series 或 NumPy 数组
\# y = pd.Series(...) 或 y = np.array(...)
使用 Pandas
class_counts = y.value_counts()
class_percentages = y.value_counts(normalize=True) * 100
print("Class Counts:\n", class_counts)
print("\nClass Percentages:\n", class_percentages)
使用 NumPy (适用于整数标签 0, 1, ...)
class_counts_np = np.bincount(y)
total samples = len(y)
class_percentages_np = (class_counts_np / total_samples) * 100
print("Class Counts (NumPy):\n", class_counts_np)
print("\nClass Percentages (NumPy):\n", class_percentages_np)
```

\* \*\*可视化: \*\* 绘制条形图或饼图来直观展示类别分布。

通常,如果少数类占比低于 10%-20%,就可以认为是显著的不平衡,需要特别处理。极端情况下,少数类占比可能低于 1%。

### 12.2 处理不平衡数据的策略

处理类别不平衡没有万能的方法,通常需要根据具体问题和数据尝试多种策略。主要可以分为三大类: 数据层面方法、算法层面方法和评估层面方法。

#### 12.2.1 数据层面方法: 重采样 (Resampling)

这类方法通过直接修改训练数据集来平衡类别分布。

- 1. \*\*欠采样 (Undersampling):\*\*
  - \* \*\*核心思想: \*\* \*\*减少多数类\*\*的样本数量,使其与少数类的数量相当或接近某个比例。
  - \* \*\*常用方法:\*\*
- \* \*\*随机欠采样 (Random Undersampling):\*\* 随机地从多数类中删除样本。最简单, 但可能丢失多数类中的重要信息。
- \* \*\*Tomek Links:\*\* 查找并移除属于 Tomek Links 的多数类样本。Tomek Link 指的是一对不同类别但互为最近邻的样本点,移除多数类样本有助于清晰化类别边界。
  - \* \*\*NearMiss: \*\* 基于距离选择要保留的多数类样本。例如,选择与 K 个最近的少数类

样本平均距离最小的多数类样本。

\* \*\*Edited Nearest Neighbours (ENN):\*\* 移除那些其类别与 K 个近邻中大多数样本类别不同的多数类样本。

- \* \*\*聚类欠采样:\*\* 对多数类样本进行聚类,然后使用聚类中心或按比例从每个簇中采样。
  - \* \*\*优点:\*\* 可以显著减少训练数据量,加快模型训练速度。
- \* \*\*缺点:\*\* 可能丢失多数类的重要信息,导致模型欠拟合或泛化能力下降。\*\*不推荐用于数据量本身就不大的情况。\*\*
  - \* \*\*实现库: \*\* `imbalanced-learn` (imblearn) 库提供了上述多种欠采样方法的实现。
- 2. \*\*过采样 (Oversampling):\*\*
  - \* \*\*核心思想: \*\* \*\*增加少数类\*\*的样本数量,使其与多数类的数量相当或接近某个比例。
  - \* \*\*常用方法:\*\*
- \* \*\*随机过采样 (Random Oversampling):\*\* 随机地\*\*复制\*\*少数类的样本。最简单,但容易导致模型在少数类上过拟合,因为它只是重复了现有信息。
- \* \*\*SMOTE (Synthetic Minority Over-sampling Technique):\*\* \*\*最常用和推荐的过采样方法之一。\*\* 不是简单复制,而是\*\*合成新的少数类样本\*\*。
- \* 工作原理:对于每个少数类样本,找到其 K 个最近的少数类邻居。然后,在该样本与其随机选择的一个邻居之间的连线上,随机生成一个新的合成样本。
  - ' 优点:可以生成更多样化的少数类数据,缓解随机过采样带来的过拟合问题。
- \* 变种: Borderline-SMOTE (更关注边界附近的少数类样本), ADASYN (Adaptive Synthetic Sampling, 根据样本学习难度自适应地生成更多样本)。
  - \* \*\*优点: \*\* 不会丢失多数类信息。通常比欠采样效果更好。
- \* \*\*缺点:\*\* 可能增加训练时间和计算成本。如果合成样本的方式不当,可能引入噪声或模糊 类别边界。SMOTE 对高维数据效果可能下降。
  - \* \*\*实现库: \*\* `imbalanced-learn` (imblearn) 库提供了 SMOTE 及其变种的实现。
- \*\*`imbalanced-learn` (imblearn) 库:\*\*

这是一个与 Scikit-learn 兼容的 Python 库,专门用于处理不平衡数据集。它提供了各种重采样技术和一些对类别不平衡友好的算法。

```
```python
```

- #安装 imblearn (如果未安装)
- # conda install -c conda-forge imbalanced-learn
- # 或者 pip install -U imbalanced-learn

from imblearn.under sampling import RandomUnderSampler, TomekLinks

from imblearn.over_sampling import RandomOverSampler, SMOTE

from imblearn.pipeline import Pipeline as ImbPipeline # 使用 imblearn 的 Pipeline

from sklearn.linear model import LogisticRegression

from sklearn.model_selection import train_test_split, StratifiedKFold

from sklearn.preprocessing import StandardScaler

from collections import Counter

- # (假设 X, y 是不平衡的原始数据)
- # X, y = make_classification(n_samples=1000, n_features=20, weights=[0.95, 0.05],
 random_state=42)
- # print("Original dataset shape %s" % Counter(y))
- # 划分训练集和测试集 (!!! 重采样只应在训练集上进行 !!!)

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,

random_state=42, stratify=y)

print("Train dataset shape %s" % Counter(y_train))

```
# print("Test dataset shape %s" % Counter(y_test))
# --- 示例: 使用 SMOTE 过采样 ---
print("\nApplying SMOTE...")
smote = SMOTE(random state=42, k neighbors=5) # k neighbors 是 SMOTE 的参数
X_train_resampled, y_train_resampled = smote.fit_resample(X_train, y_train)
print("Resampled dataset shape (SMOTE) %s" % Counter(y_train_resampled))
# --- 示例: 使用随机欠采样 ---
print("\nApplying Random Undersampling...")
rus = RandomUnderSampler(random_state=42)
X_train_undersampled, y_train_undersampled = rus.fit_resample(X_train, y_train)
print("Resampled dataset shape (Random Undersampling) %s" %
Counter(y_train_undersampled))
# --- 在交叉验证中使用重采样 (推荐方式) ---
# 使用 imblearn 的 Pipeline 可以确保重采样只在当前折的训练数据上进行
print("\nUsing SMOTE within Cross-Validation Pipeline...")
# 创建模型和 SMOTE 实例
model = LogisticRegression(solver='liblinear', random_state=42)
smote_cv = SMOTE(random_state=42)
scaler_cv = StandardScaler()
# 创建 imblearn Pipeline
# 注意: 重采样步骤应该在模型之前
pipeline_smote = ImbPipeline([
   ('scaler', scaler_cv),
    ('sampling', smote_cv), # 重采样步骤
    ('classification', model)
1)
# 使用 cross_val_score 进行评估 (使用原始未重采样的 X_train, y_train)
cv_strategy_imb = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
#选择合适的评估指标,如 roc_auc 或 f1
scores_smote_cv = cross_val_score(pipeline_smote, X_train, y_train,
                                scoring='roc_auc', cv=cv_strategy_imb,
n jobs=-1
print(f"Cross-Validation AUC Scores with SMOTE: {scores_smote_cv}")
print(f"Mean AUC with SMOTE: {scores smote cv.mean():.4f}")
# 对比不使用重采样的结果
pipeline no sampling = Pipeline([
    ('scaler', StandardScaler()),
    ('classification', LogisticRegression(solver='liblinear', random_state=42))
])
scores_no_sampling_cv = cross_val_score(pipeline_no_sampling, X_train, y_train,
                                      scoring='roc_auc', cv=cv_strategy_imb,
n_jobs=-1)
print(f"\nCross-Validation AUC Scores without Sampling: {scores no sampling cv}")
print(f"Mean AUC without Sampling: {scores_no_sampling_cv.mean():.4f}")
```

重采样注意事项:

• **只对训练集重采样**: 绝对不要对测试集进行重采样。测试集必须保持原始的数据分布,以真实反映模型在实际未见数据上的表现。

- **在交叉验证中正确应用**: 如果使用交叉验证(如 GridSearchCV, cross_val_score),必须确保重采样 步骤发生在**每个折内部的训练数据上**。使用 imblearn.pipeline.Pipeline 是实现这一点的最佳方 式。
- **过采样 vs. 欠采样**: 通常优先尝试过采样(特别是 SMOTE),因为它不丢失信息。只有在数据集非常大,训练时间是主要瓶颈时才考虑欠采样。也可以尝试结合使用(例如,先 SMOTE 再 Tomek Links)。

12.2.2 算法层面方法: 调整算法和权重

这类方法通过修改学习算法本身或调整其参数来应对不平衡。

1. 调整类别权重 (Class Weighting):

- 核心思想:在模型的损失函数中为不同类别的样本赋予不同的权重。通常给少数类样本赋予更高的权重,使得模型在优化时更关注对少数类的错误分类。
- 。 **实现:** 许多 Scikit-learn 模型 (如 LogisticRegression, SVC, DecisionTreeClassifier, RandomForestClassifier) 和 XGBoost (XGBClassifier) 都提供了 class_weight 参数。
 - 可以设置为 'balanced', 算法会自动根据类别频率计算权重 (weight = n_samples / (n_classes * n_samples_with_class))。
 - 也可以手动传入一个字典, 例如 {0:1,1:10} 表示类别 1 的权重是类别 0 的 10 倍。
- XGBoost 特定参数: scale_pos_weight (用于二分类), 我们在第八章已经详细介绍过,推荐设置为 sum(negative instances) / sum(positive instances)。
- · **优点:** 实现简单,不改变数据本身,计算效率较高。
- 缺点: 需要选择合适的权重 ('balanced' 不一定最优) , 效果可能不如重采样显著。
- **应用关联:** 你的脚本 2, 3, 5, 6 中都计算并使用了 scale_pos_weight 或 class_weight, 这是处理不平衡数据的关键步骤。

2. 代价敏感学习 (Cost-Sensitive Learning):

- 。 **核心思想**: 直接将不同误分类的**代价**引入到模型的学习过程中。例如,将 FN 的代价设置得远高于 FP 的代价。
- **实现:** 某些算法(如一些 SVM 实现)允许直接指定误分类代价矩阵。更通用的方法是通过调整类别权重或修改损失函数来间接实现。

3. **集成方法:**

- 集成不平衡学习器:可以专门设计用于不平衡数据的集成方法、例如:
 - EasyEnsemble: 多次从多数类中进行欠采样,训练多个基于欠采样数据的模型,然后集成它们的预测。
 - BalanceCascade: 迭代地进行欠采样和模型训练,每次迭代中被正确分类的多数类样本会被 移除,使得后续模型更关注难分的多数类样本。
 - RUSBoost / SMOTEBoost: 将 Boosting 算法与随机欠采样 (RUS) 或 SMOTE 结合。
- **实现库:** imbalanced-learn 提供了这些集成方法的实现 (EasyEnsembleClassifier, RUSBoostClassifier 等)。

12.2.3 评估层面方法:选择合适的指标和阈值调整

即使使用了数据或算法层面的方法,选择正确的评估指标仍然至关重要。

- 1. 使用合适的评估指标: (已在第七章详细介绍)
 - 避免 Accuracy。
 - 关注 Recall (POD), Precision, F1-Score, AUC, PR-AUC, CSI 等对少数类敏感的指标。
 - 。 根据业务需求确定优先关注哪个指标(例如,更怕漏报还是误报?)。

2. **阈值调整 (Threshold Moving):**

- **背景:** 大多数分类器默认使用 0.5 的概率阈值来区分正负类。但在不平衡数据或不同错误代价场景下,0.5 不一定是最佳阈值。
- **核心思想:** predict_proba 提供了每个样本属于正类的概率。我们可以不使用默认的 0.5,而是根据验证集上的性能(例如,最大化 F1 分数,或者在满足某个 Recall/POD 要求下最大化 Precision)来选择一个**不同的最优阈值**。
- 方法:
 - 在验证集上获取模型预测的概率 v val proba。
 - 尝试一系列不同的阈值 (例如从 0.1 到 0.9) 。
 - 对于每个阈值,根据 y_val_proba >= threshold 得到预测类别 y_val_pred_threshold。
 - 计算每个阈值下的目标评估指标(如 F1, Precision, Recall)。
 - 选择使得目标指标最优的那个阈值。
 - **最终预测:** 在测试集上获取概率 y_test_proba, 然后使用找到的最优阈值进行最终类别判断 y_test_pred = (y_test_proba >= best_threshold).astype(int)。
- o **可视化:** 可以绘制 Precision-Recall 曲线,并在曲线上寻找合适的平衡点。
- · **优点:** 简单有效,不改变模型训练过程,直接调整决策边界。
- 缺点: 找到的阈值是针对特定数据集和指标的,可能在新数据上不是最优。
- o **应用关联**: 你的脚本 2, 3, 5, 6 中都包含了遍历多个阈值 (thresholds_to_evaluate) 并计算相应 指标的代码。这实际上就是在探索不同阈值对模型性能的影响,帮助用户根据打印出的表格选择 一个在特定指标(如 POD, FAR, CSI)上表现符合要求的阈值,用于最终的决策或模型比较。

阈值调整代码示例:

```
from sklearn.metrics import precision_recall_curve, f1_score, precision_score, recall_score import numpy as np

# (假设 y_val_true 是验证集真实标签, y_val_proba 是模型在验证集上预测的正类概率)
# y_val_proba = model.predict_proba(X_val_scaled)[:, 1]

# 1. 计算 Precision, Recall 和 Thresholds
# precision, recall, thresholds_pr = precision_recall_curve(y_val_true, y_val_proba)

# 2. 寻找最大化 F1 分数的阈值
# f1_scores = []
# thresholds_to_check = thresholds_pr # PR 曲线的阈值, 或者自己定义一个范围 np.linspace(0, 1, 101)
# for thresh in thresholds_to_check:
# y_pred_thresh = (y_val_proba >= thresh).astype(int)
```

```
f1 = f1_score(y_val_true, y_pred_thresh)
     f1 scores.append(f1)
# if len(f1_scores) > 0:
     best f1 idx = np.argmax(f1 scores)
      best threshold f1 = thresholds to check[best f1 idx]
      print(f"\nBest Threshold for F1 Score: {best_threshold_f1:.4f} (F1=
{f1 scores[best f1 idx]:.4f})")
# else:
     print("\nCould not determine best threshold.")
      best_threshold_f1 = 0.5 # Fallback to default
# 3. (可选) 绘制 PR 曲线并标记最佳 F1 点
# plt.figure()
# plt.plot(recall[:-1], precision[:-1], label='Precision-Recall curve') #
thresholds 比 p,r 少一个
# plt.scatter(recall[best_f1_idx], precision[best_f1_idx], marker='o',
color='red', label=f'Best F1 Threshold ({best threshold f1:.2f})')
# plt.xlabel('Recall (POD)')
# plt.ylabel('Precision')
# plt.title('Precision-Recall Curve')
# plt.legend()
# plt.grid(True)
# plt.show()
# 4. 应用最佳阈值到测试集
# y_test_proba = model.predict_proba(X_test_scaled)[:, 1]
# y_test_pred_best = (y_test_proba >= best_threshold_f1).astype(int)
# print("\nFinal metrics on test set using best F1 threshold:")
# print(classification_report(y_test, y_test_pred_best))
```

12.3 策略选择

面对类别不平衡问题,通常建议:

1. 从简单的开始:

- o 首先,使用合适的**评估指标** (AUC, PR-AUC, F1, Recall, Precision, CSI 等) 来评估你的基线模型性能,了解问题的严重性。
- 。 尝试使用算法内置的**类别权重** (class_weight='balanced' 或 scale_pos_weight),这通常是最简单有效的第一步。

2. 尝试数据层面的方法:

- 如果调整权重效果不佳,可以尝试**过采样**,特别是 **SMOTE** 及其变种。使用 imblearn 库,并确保在交叉验证流程中正确应用。
- 如果数据集非常大,可以谨慎地尝试欠采样或结合方法。
- 3. **考虑阈值调整:** 在使用概率输出的模型(如 Logistic Regression, XGBoost)上,阈值调整是一种非常灵活 且有效的策略,可以根据业务需求在 Precision 和 Recall (或 POD 和 FAR) 之间进行权衡。
- 4. 尝试集成方法: 如果上述方法效果仍不理想,可以考虑 imblearn 中专门为不平衡数据设计的集成方法。
- 5. **收集更多数据 (如果可能):** 特别是收集更多**少数类**的数据,往往是解决不平衡问题的最根本、最有效的方法,尽管通常也是最困难或成本最高的。

处理类别不平衡是一个需要反复实验和细致评估的过程。理解各种方法的优缺点,并结合具体问题选择合适的策略组合,是提升模型在不平衡数据上表现的关键。

第十三章: 处理大型数据集策略

随着数据量的爆炸式增长,机器学习项目经常需要处理远超单机内存容量的大型数据集。当数据无法一次性加载到内存中时,标准的内存计算方法就会失效,我们需要采用特殊的策略来应对这一挑战。

你的脚本示例中已经使用了其中两种关键技术: **内存映射** (Memory Mapping) 和 **分块处理** (Chunking/Batching)。本章将更详细地探讨这些策略以及其他相关方法。

13.1 挑战: 内存限制

传统的机器学习工作流(包括许多 Scikit-learn 和 Pandas 的操作)通常假设整个数据集可以被加载到 RAM中。当数据集大小达到几十 GB 甚至 TB 级别时,这显然是不可行的。强行加载会导致 Memory Error,程序崩溃。

13.2 核心策略

13.2.1 内存映射 (Memory Mapping)

我们在第三章介绍 NumPy 时已经详细讲解了内存映射。这里回顾并强调其在处理大型数据集中的作用。

• **原理**: 通过 np.load(filename, mmap_mode='r') 加载 .npy 文件时,并不会将文件的全部内容读入内存。它创建了一个指向磁盘文件数据的"指针"或"视图"。当你通过索引或切片访问这个内存映射数组时,操作系统负责按需从磁盘加载相应的数据块到内存页面。

• 优点:

- 。 允许访问和操作远超内存大小的数据文件。
- 对于只需要读取部分数据(例如,按行索引、切片)的操作非常高效,避免了加载整个文件的开销。
- 多个进程可以映射同一个文件进行只读访问,节省内存。

• 缺点:

- 。 磁盘 I/O 速度远慢于内存访问速度。如果需要频繁地、随机地访问数组的大部分区域,性能可能会受到磁盘速度的限制。
- 某些需要全量数据的 NumPy 或 Scikit-learn 操作(例如,计算整个数组的标准差、某些复杂的矩阵运算)仍然可能需要将数据块加载到内存中,如果操作本身内存需求过大,仍可能失败。
- 。 需要手动管理文件句柄的关闭(如 mmap_arr._mmap.close()),尤其是在长时间运行的程序或循环中。
- **应用关联:** 你的脚本 1, 4, 6 在加载 X_flat_features.npy 和 Y_flat_target.npy 时广泛使用了mmap_mode='r'。这是能够处理可能非常大的特征和目标数组的基础。

13.2.2 分块处理 (Chunking / Batching)

当内存映射本身不足以解决问题 (例如,需要对数据进行复杂计算或模型预测,而单次操作的内存需求仍然过大) 时,分块处理是关键策略。

• **原理**: 将大型数据集(无论是来自内存映射数组还是直接从磁盘文件读取)分成若干个较小的、可以方便地加载到内存中的数据块(Chunk 或 Batch)。然后,对每个数据块依次进行处理(如预处理、模型预

测、指标计算等)。最后,根据需要合并各个块的处理结果。

- 实现:通常使用循环和切片来实现。
 - 确定合适的块大小 (CHUNK_SIZE)。这个大小取决于可用内存、单次处理的计算复杂度等因素。需要进行实验和调整。
 - 使用 range(0, n total samples, CHUNK SIZE) 循环遍历数据。
 - o 在循环内部,加载或切片出当前的数据块(X_chunk = data[start:end])。
 - 对 X chunk 执行所需操作 (例如 model.predict proba(X chunk)) 。
 - 。 存储或累积当前块的结果。
 - 。 (重要) 在循环结束时释放当前块占用的内存 (del X_chunk, chunk_result),以便为下一个块腾出空间。

• 优点:

- 。 内存占用可控,可以处理任意大小的数据集,只要单个块能放入内存。
- 。 可以与内存映射结合使用,先通过内存映射访问整个数据集,再分块加载进行计算。

缺点:

- 代码逻辑更复杂,需要管理循环、索引和结果合并。
- 某些需要全局信息的操作(例如,计算整个数据集的精确均值或标准差)可能需要多次遍历数据或使用近似算法。
- 总处理时间可能会因为频繁的数据加载和重复计算(如果块之间有重叠或需要共享状态)而增加。

• 应用关联:

- 。 脚本 4_generate_meta_features.py: 在使用 FP 专家模型对整个测试集进行预测以生成元特征时,采用了分块处理。它循环读取 X_test (来自内存映射或已加载的切片) 和 base_probs_test的块,对每个块调用 fp_expert_model.predict_proba,然后将结果存入预分配的 X_meta 数组中。
- 。 脚本 6_finetune_base_model.py: 在评估微调后的模型时,对测试集的预测也使用了分块。它循环从内存映射的 X_flat_test_mmap 中加载特征块 X_chunk,调用finetune_model.predict_proba,并将结果追加到列表 y_test_pred_proba_list 中,最后再将列表中的所有块结果合并 (np.concatenate)。

13.2.3 选择性加载数据

- **只加载需要的列**: 如果原始数据是 CSV、Parquet 或数据库表等格式,并且你只需要其中的一部分特征 (列),那么在读取数据时就只指定加载这些列。这可以大大减少内存占用。
 - o pd.read_csv(..., usecols=['col1', 'col5', ...])
 - o pd.read_parquet(..., columns=['col1', 'col5', ...])
- **指定更节省内存的数据类型 (dtype):** 在加载数据时(例如使用 Pandas),如果知道某列的数值范围,可以指定更精确、占用内存更小的数据类型。
 - 例如,如果某整数列的值都在 0-255 之间,可以使用 np.uint8 (1字节) 而不是默认的 np.int64 (8字节)。
 - 对于浮点数,可以使用 np.float32 (4字节) 而不是 np.float64 (8字节)。
 - o 对于类别特征,如果类别数量不多,Pandas 的 Categorical 类型比 object (字符串) 更节省内存。
 - o pd.read csv(..., dtype={'col int': np.int32, 'col float': np.float32})

13.2.4 使用更高效的数据格式

• NumPy .npy: 对于纯数值数组,非常高效。支持内存映射。

- Parquet (.parquet): 一种列式存储格式,非常适合存储表格数据(如 Pandas DataFrame)。
 - **优点**: 读取速度快(特别是只读取部分列时),压缩效率高,支持复杂数据类型,跨平台/语言兼容性好。pandas.read_parquet, pandas.to_parquet。
- **HDF5** (.h5, .hdf5): 一种层次化数据格式,可以存储多个数据集(类似文件系统),支持压缩、分块和部分读写。适合存储大型、复杂的科学数据集。pandas.read_hdf, pandas.to_hdf, h5py 库。
- **Feather (.feather):** 一种快速、轻量级的二进制列式存储格式,专为 Pandas DataFrame 设计,读写速度极快。但在压缩和跨平台兼容性方面可能不如 Parquet。

选择哪种格式取决于具体需求(读写速度、压缩率、兼容性、数据复杂度等)。对于大型表格数据,Parquet 是一个非常好的通用选择。

13.2.5 增量学习 / 在线学习 (Incremental / Online Learning)

对于某些模型算法,存在支持**增量学习**或**在线学习**的实现。这意味着模型可以**在不重新加载整个数据集的情况下,根据新到达的数据块进行更新**。

- Scikit-learn: 提供了一些支持 partial fit 方法的模型,允许增量学习。
 - linear model.SGDClassifier, linear model.SGDRegressor (随机梯度下降)
 - naive_bayes.MultinomialNB, naive_bayes.BernoulliNB
 - cluster.MiniBatchKMeans
 - decomposition.IncrementalPCA
- 其他库: 有些库专门设计用于处理流式数据和在线学习(如 Vowpal Wabbit, River)。
- 优点: 不需要存储整个数据集,可以处理无限的数据流,模型可以适应数据分布的变化。
- 缺点: 不是所有算法都支持增量学习。实现和管理可能更复杂。模型的最终性能可能不如在整个数据集上进行批量训练。

应用关联: 你的脚本没有使用增量学习,而是采用了批量训练(可能结合内存映射和分块)的方式。

13.2.6 使用基于直方图的算法

正如在第八章提到的,XGBoost 和 LightGBM 等现代梯度提升库提供了基于直方图的树构建算法(tree method='hist')。这种算法通过将连续特征值分箱来减少需要考虑的分裂点数量,从而:

- 显著降低内存消耗: 不再需要存储排序后的特征值。
- 大幅提高训练速度: 查找最佳分裂点的计算量减少。

这是处理大型数据集时,这些库相比传统 GBDT 或 Scikit-learn 的 GradientBoostingClassifier 的一个巨大优势。

13.2.7 分布式计算 (Distributed Computing)

当单台机器的内存和计算能力都无法满足需求时(例如,需要处理 TB 级别的数据或训练极其复杂的模型), 就需要借助分布式计算框架。

- 核心思想: 将数据和计算任务分布到多台机器组成的集群上并行处理。
- 常用框架:
 - o **Dask:** 一个灵活的 Python 并行计算库,可以很好地与 NumPy, Pandas, Scikit-learn 集成,将它们的计算扩展到多核或多机。提供了 Dask Array, Dask DataFrame 等类似 NumPy/Pandas 但能处理大型数据的接口,并支持分布式任务调度。

o Apache Spark (with PySpark & MLlib): 一个强大的、通用的分布式计算引擎。Spark MLlib 提供了用于大规模机器学习的算法库,包括数据处理、特征工程和模型训练(支持多种常用算法的分布式实现)。需要搭建 Spark 集群环境。

- 。 **Ray:** 一个用于构建分布式应用程序的框架,包含用于大规模机器学习的库 (Ray Tune for HPO, Ray Train for distributed training)。
- 。 **XGBoost / LightGBM on Dask/Spark:** XGBoost 和 LightGBM 都提供了与 Dask 和 Spark 集成的版本,可以在分布式环境上训练模型。
- 优点: 可以处理单机无法处理的超大规模数据和计算任务。
- 缺点: 系统搭建和维护更复杂,需要分布式编程的知识,网络通信可能成为瓶颈。

13.3 策略组合

在实践中,处理大型数据集通常需要组合使用上述策略:

- 1. 选择高效的数据格式存储数据 (如 Parquet)。
- 2. 只加载需要的特征列,并指定节省内存的数据类型。
- 3. 如果数据仍然过大,使用内存映射 (如果格式支持,如 npy) 或分块读取 (如 Pandas 的 chunksize 参数)。
- 4. 对于需要全量数据的计算或模型预测,使用分块处理,确保每个块能放入内存。
- 5. 优先选择内存效率高的算法实现 (如 XGBoost/LightGBM 的 hist 模式)。
- 6. 如果单机资源仍然不足,考虑使用 Dask 或 Spark 进行分布式计算。

你的脚本主要依赖于 .npy 格式、内存映射和分块预测,这对于处理单个大型特征/目标数组是一个有效的策略组合。如果在数据预处理阶段(生成 .npy 文件之前)数据量就非常大,那么选择性加载列、使用 Parquet 格式、利用 Pandas 的分块读取等策略会更加重要。

理解这些处理大型数据集的策略,可以帮助你在面对内存限制时选择合适的工具和方法,确保机器学习项目能够顺利进行。

- ```markdown
- # 机器学习实战: 从 Scikit-learn 到 XGBoost (卷五)
- ...(接上文)

- - -

第十四章: 模型微调 (Fine-tuning) 实战

在机器学习项目中,我们往往不会从零开始训练每一个模型。有时,我们已经有了一个在大型通用数据集上预训练好的模型,或者一个在先前任务中表现不错的模型。**模型微调 (Fine-tuning)** 就是利用这些已有的模型,在新的、特定的任务或数据上**继续训练**,以适应新需求的过程。

这种技术在深度学习领域(特别是自然语言处理和计算机视觉)极为常见,例如使用在 ImageNet 上预训练的 ResNet 模型来分类特定种类的花朵。但微调的概念同样适用于传统的机器学习模型,尤其是像 XGBoost 这样支持**增量训练**的模型。

你的脚本 `6_finetune_base_model.py` 就展示了对一个已有的 XGBoost 基模型进行微调的过程。本章将深入探讨模型微调的概念、动机、实现方式以及注意事项,并结合你的脚本进行分析。

14.1 微调的动机与优势

为什么不直接在新数据上重新训练一个模型, 而要进行微调呢?

1. **利用预训练知识: ** 预训练模型可能已经从大量数据中学到了通用的模式或特征表示。微调可以利用这些知识,使其在新任务上更快地收敛或达到更好的性能,尤其是在新任务的数据量有限时。

- 2. **节省训练时间和计算资源:** 从一个已经训练过的模型开始,通常比从随机初始化开始需要更少的训练迭代次数和计算量。
- 3. **适应特定数据分布: ** 原始模型可能是在与当前任务略有不同的数据分布上训练的。微调可以帮助模型适应新数据的特性。
- 4. **针对性改进:** 可以针对模型在某些特定子集(例如,之前分类错误的样本)上的表现进行强化训练,试图修正模型的弱点。
- * **应用关联:** 脚本 `6` 的核心思想就是这个。它加载了初始的基模型 (`xgboost_default_full_model.joblib`),然后专门在使用该基模型预测时出错的**训练集样本 ** (即 False Positives 和 False Negatives) 上进行额外的训练,目标是让模型在这些"困难" 样本上表现更好。

14.2 微调的实现方式 (以 XGBoost 为例)

XGBoost 的 Scikit-learn API (`XGBClassifier`, `XGBRegressor`) 通过 `fit` 方法的 `xgb_model` 参数来支持微调。

核心流程:

1. **加载预训练/初始模型:** 使用 `joblib.load()` (或者 XGBoost 的原生加载方法) 加载你想要微调的模型对象。

```
```python
 import joblib
 import xgboost as xgb
 import os
 # 加载之前训练好的模型
 initial model path = 'path/to/your/initial model.joblib'
 if os.path.exists(initial model path):
 try:
 initial_model = joblib.load(initial_model_path)
 print("Initial model loaded successfully.")
 except Exception as e:
 print(f"Error loading initial model: {e}")
 # exit() # Handle error appropriately
 else:
 print(f"Error: Initial model not found at {initial model path}.")
 # exit()
 应用关联: 脚本 `6` 中 `initial_model =
joblib.load(INITIAL MODEL PATH)`.
```

- 2. \*\*准备微调数据:\*\* 确定用于微调的数据集 (`X\_finetune`, `y\_finetune`)。这可以是:
  - \* 一个新的、相关的任务数据集。
  - \* 原始训练数据的一个子集(例如,模型之前预测错误的样本)。
  - \* 原始训练数据加上一些新的数据。
- \* \*\*应用关联:\*\* 脚本 `6`中通过加载之前保存的 FP 和 FN 样本的索引 (`indices\_fp`, `indices\_fn`), 然后使用这些索引从完整的 `X\_flat` 和 `Y\_flat` (内存映射) 中提取出对应的特征和标签,构成了微调数据集 (`X\_finetune`, `y\_finetune\_true`)。

```
3. **定义微调参数:**
```

\* 通常,微调时会使用\*\*更小的学习率 (`learning\_rate`)\*\*。这是因为我们不希望剧烈地改变已经学到的知识,而是希望在原有基础上进行微小的调整。

- \* 通常,微调时的\*\*迭代次数(`n\_estimators`)\*\* 会比从头训练时少。我们只是在添加一些额外的树来适应新数据或修正错误。
- \* \*\*其他参数: \*\* 可以沿用初始模型的参数,也可以根据微调数据的特点进行调整(例如,`scale\_pos\_weight`可能需要根据微调数据集的类别比例重新计算)。
- \* \*\*早停:\*\* 在微调时是否使用早停取决于目标。如果希望强制模型在特定数据上学习指定轮数(如脚本 6),则可以禁用早停(`early\_stopping\_rounds=None`)。如果希望避免在微调数据上过拟合,则仍应使用早停(需要提供 `eval\_set`)。
- \* \*\*应用关联:\*\* 脚本 `6`中定义了 `FINETUNE\_LR = 0.02` (较小学习率) 和 `FINETUNE\_N\_ESTIMATORS = 50` (较少迭代次数)。它还根据微调数据集重新计算了 `scale\_pos\_weight\_finetune`,并显式地将 `early\_stopping\_rounds` 设置为 `None`。
- 4. \*\*执行微调:\*\* 创建一个新的 XGBoost Estimator 实例 (可以先获取 `initial\_model.get\_params()` 并修改学习率、估计器数量等), 然后在调用 `fit` 方法时,将 \*\*加载的初始模型\*\*传递给 `xgb\_model` 参数。
  - ```python

# 获取初始模型参数

```
params = initial_model.get_params()
```

# 更新微调参数

```
params['learning_rate'] = 0.02 # 使用较小的学习率 params['n_estimators'] = 50 # 添加 50 棵新树 params['early_stopping_rounds'] = None # 示例: 禁用早停 # 根据需要重新计算 scale_pos_weight (假设已计算好 finetune_scale_pos) # params['scale_pos_weight'] = finetune_scale_pos # 创建新的模型实例用于微调(或者直接修改 initial_model 的参数) finetune_model = xgb.XGBClassifier(**params)
```

```
使用微调数据和初始模型进行 fit
finetune_model.fit(
 X_finetune,
 y_finetune,
 xgb_model=initial_model, # !!! 传入初始模型 !!!
 verbose=False # 或根据需要设置
 # 如果需要早停, 还需传入 eval_set 和设置 early_stopping_rounds
)
print("Fine-tuning complete.")
```

print(f"Starting fine-tuning for {params['n\_estimators']} rounds...")

- \* \*\*应用关联:\*\* 脚本 `6` 中的 `finetune\_model.fit(X\_finetune, y\_finetune\_true, xgb\_model=initial\_model, verbose=False)` 正是执行微调的核心步骤。
- 5. \*\*保存微调后的模型:\*\* 将 `finetune\_model` 保存到新的文件。

```python

finetuned_model_path = 'path/to/your/finetuned_model.joblib'
trv:

joblib.dump(finetune_model, finetuned_model_path)
print(f"Fine-tuned model saved to {finetuned_model_path}")
except Exception as e:

print(f"Error saving fine-tuned model: {e}")

. . .

- * **应用关联:** 脚本 `6` 中 `joblib.dump(finetune_model, FINETUNED_MODEL_PATH)`。
- 6. **评估微调后的模型:** 使用独立的**测试集**来评估微调后的模型性能,并与初始模型的性能进行比较,判断微调是否带来了预期的改进。
- * **应用关联: ** 脚本 `6` 的后半部分就是加载原始测试集(分块处理),使用 `finetune_model` 进行预测,并计算各种指标,最后打印性能表格。

14.3 微调的注意事项

- * **学习率是关键:** 微调时通常需要显著降低学习率,以避免破坏预训练模型学到的知识。典型值可能是原始学习率的 1/10 或更小。
- * **微调轮数:** 通常不需要很多轮迭代。具体轮数需要实验确定,或者通过在合适的验证集上使用早停来自动决定。
- * **数据相关性: ** 微调的效果很大程度上取决于初始模型训练数据与微调数据之间的相关性。如果两者差异巨大,微调的效果可能不佳,甚至不如从头训练。
- * **参数调整:** 除了学习率和轮数,其他超参数(如树的深度、正则化参数)是否需要调整取决于具体情况。有时保持初始模型的结构参数不变效果更好,有时则需要根据微调数据的特点进行调整。
- * **冻结层(深度学习概念):** 在深度学习中,微调有时会"冻结"模型的部分底层(这些层被认为学到了通用特征),只训练顶部的几层来适应新任务。这个概念在 XGBoost 微调中不直接适用,但降低学习率和限制轮数起到了类似的作用—主要调整后期添加的树,对早期树的影响较小。
- * **评估: ** 必须在独立的测试集上评估微调效果,并与原始模型进行公平比较。如果在微调中使用的数据(或其一部分)也用于评估,结果会过于乐观。脚本 `6` 在原始测试集上评估是正确的做法。
- * **过拟合风险: ** 即使进行微调,仍然存在在微调数据上过拟合的风险,特别是当微调数据量很少或者微调轮数过多/学习率过高时。使用验证集进行监控(如果适用)或谨慎选择微调参数很重要。

14.4 微调 vs. Stacking/专家模型

脚本 $^{\circ}$ (微调) 和脚本 $^{\circ}$ (专家模型 + Stacking 变体) 代表了两种不同的提升基模型性能的思路:

- * **微调 (脚本 6):**
 - * **目标:** 直接改进**原始模型本身**,让它在"困难"样本上表现更好。
 - * **方法:** 在困难样本上**继续训练**原始模型。
 - * ** 最终产物: ** 一个 ** 单一的、增强的 ** XGBoost 模型

(`xgboost_finetuned_model.joblib`).

- * **预测: ** 直接使用这个微调后的模型进行预测。
- * **专家模型 + Stacking (脚本 2-5):**
- * **目标:** 保持原始模型不变,额外训练**辅助模型(专家)**来识别和纠正原始模型的特定错误,然后训练一个**元学习器**来智能地**组合**原始模型和专家模型的(或其衍生的)信息。
- * **方法:** 分别训练 FN 专家、FP 专家 (利用原始特征和基模型概率), 然后生成包含基模型概率和 FP 专家概率的元特征,最后训练元学习器。
 - * **最终产物: ** **多个模型** (基模型、FN 专家、FP 专家、元学习器)。
- * **预测:** 需要一个**多阶段**的过程: 获取原始特征 -> 基模型预测概率 -> FP 专家输入准备 -> FP 专家预测概率 -> 构建元特征 -> 元学习器最终预测。

**选择哪种方法? **

- * **微调**通常更简单直接,如果目标是略微提升现有模型在特定数据子集上的表现,并且不希望增加预测时的复杂度,这是一个好选择。但它直接修改了原始模型。
- * **Stacking/专家模型**方法更复杂,但可能更强大,因为它试图从不同角度(原始预测、错误模式分析)来组合信息。它保留了原始模型,并通过元学习器进行"决策融合"。预测流程更复杂。

选择哪种方法取决于具体需求、计算资源、对预测流程复杂度的接受程度以及实验效果。有时甚至可以结合使用,例如,先微调基模型,再将微调后的模型作为 Stacking 的一部分。

理解模型微调的概念和在 XGBoost 中的实现方式,为你提供了一种重要的模型优化策略,特别是当你希望利用现有模型或针对特定数据子集进行改进时。

第十五章: 案例分析: 天气预测模型代码解读

现在,我们已经学习了 NumPy, Pandas, Scikit-learn, XGBoost, Stacking, Fine-tuning, 不平衡数据处理,大数据策略等核心概念和技术。是时候回过头来,以更深入的理解,系统性地分析你提供的 8 个 Python 脚本所构成的天气预测模型工作流了。

本章将逐一解读每个脚本的核心目标、关键步骤、使用的技术以及它们之间的联系,帮助你将前面学到的知识融会贯通,并理解这个特定案例是如何应用这些技术的。

(注意: 由于`xgboost1.py`未提供,我们将假定它存在并完成了基础 XGBoost 模型的训练,并保存为`xgboost_default_full_model.joblib`。同时,假定数据预处理(特征工程、编码等)已在生成`X_flat_features.npy`和`Y_flat_target.npy`之前完成。)

15.1 脚本 `1_generate_base_predictions.py`

- * **目标:**
- 1. 加载一个预训练好的**基线 XGBoost 模型** (`xgboost_default_full_model.joblib`)。
 - 2. 使用该模型对**完整的训练集**进行预测。
- 3. 基于预测结果和真实标签,识别训练集中的 **TP, TN, FP, FN** 样本,并保存它们的**索引**。这些索引是后续训练专家模型的基础。
- 4. 保存基线模型在**训练集**和**测试集**上的**预测概率**。这些概率将作为后续 FP 专家模型和元学习器的重要输入(元特征)。
- 5. (可选地)保存基线模型在测试集上的**二元预测结果**(基于 0.5 阈值),可用于直接评估基线性能。
- * **关键步骤与技术:**
- * **配置路径: ** 定义所有输入文件、输出中间数据和最终预测文件的路径。保持路径结构清晰对复杂工作流很重要。
- * **加载数据 (内存映射):** 使用 `np.load(..., mmap_mode='r')` 加载大型特征 (`X_flat`) 和目标 (`Y_flat`) 数组,避免内存溢出。**(第 3,13 章)**
- * **数据划分 (隐式): ** 通过索引 `split_idx` 区分训练集和测试集范围,确保与模型训练时使用的划分一致。 **(第 5 章) **
 - * **加载模型 (`joblib`):** 加载预训练的基线 XGBoost 模型。**(第 5 章)**
 - * **训练集处理:**
 - * 切片获取训练数据 `X_train` 和 `Y_train_raw`。 **(第 3 章)**
- * 将原始目标 `Y_train_raw` 根据 `RAIN_THRESHOLD` 转换为二元标签 `y train true`。
- * 使用 `base_model.predict_proba(X_train)` 获取训练集概率。**(第 5,8 章)**
 - * 保存训练集概率 (`base_probs_train.npy`)。**(第 3 章)**
 - * 根据 `PREDICTION_THRESHOLD` (0.5) 计算二元预测 `y_train_pred`。
- * **布尔索引识别 TP/TN/FP/FN:** 使用 `(y_train_true == 1) & (y_train_pred == 1) `等逻辑运算创建布尔掩码。**(第 3 章)**
- * **`np.where()` 获取索引:** 找到布尔掩码中 `True` 值对应的索引。**(第 3 章)**
 - * 保存 TP/TN/FP/FN 索引 (`indices *.npy`)。**(第 3 章)**

- * **测试集处理:**
 - * 切片获取测试数据 `X test`。**(第 3 章)**
 - * 使用 `base_model.predict_proba(X_test)` 获取测试集概率。**(第 5, 8 章)**
 - * 保存测试集概率 (`base_probs_test.npy`)。**(第 3 章)**
 - * 根据 `PREDICTION_THRESHOLD` 计算二元预测 `y_test_pred`。
 - * 保存测试集二元预测 (`base_preds.npy`)。**(第 3 章)**
- * **内存管理:** 使用 `del` 清理不再需要的变量,并注意在脚本结束或不再需要时关闭内存映射 (`_mmap.close()`)。**(第 3,13 章)**
- * **脚本间联系:**
 - * 输出的 **TP/TN/FP/FN 索引** 是脚本 `2`, `3`, `6`,
- `evaluate_fp_expert_on_tp.py` 的输入,用于选择特定的训练/分析样本。
- * 输出的 **训练集概率 (`base_probs_train.npy`)** 是脚本 `3` (FP 专家) 和 `evaluate_fp_expert_on_tp.py` 的输入,作为额外特征或用于分析。
- * 输出的 **测试集概率 (`base_probs_test.npy`)** 是脚本 `4` (元特征生成) 的关键输入。

15.2 脚本 `2_train_fn_expert.py`

- * **目标:** 训练一个 XGBoost 模型,专门用于**识别基线模型漏报 (FN) 的样本**。目标是提高整体的召回率 (POD)。
- * **核心思想: ** 模型试图在"实际下雨但基模型预测未下雨 (FN)"和"实际未下雨基模型也预测未下雨 (TN)"这两类样本中,学习区分前者的模式。
- * **关键步骤与技术:**
 - * **加载索引:** 加载 `indices_fn.npy` 和 `indices_tn.npy`。
 - * **准备训练数据:**
 - * 合并 FN 和 TN 索引 `combined_indices`。
 - * 创建目标标签 `labels`: FN 样本标记为 1 (目标是识别它们), TN 样本标记为 0。
- * **花式索引加载特征:** 使用 `combined_indices` 从 `X_flat` (内存映射) 中加载 FN 和 TN 样本的**原始特征** `X_expert_data`。**(第 3 章)**
- * **划分训练/验证集:** 使用 `train_test_split` 将 FN/TN 数据划分为训练集和验证集,并使用 `stratify=labels` 确保类别比例。**(第 5, 11 章)**
 - * **处理类别不平衡(算法层面):**
- * 计算 `scale_pos_weight`: `num_tn / num_fn`, 给少数类 (FN 样本, 标签为 1) 更高的权重。**(第 8, 12 章)**
 - * 在 `XGBClassifier` 参数中设置 `scale_pos_weight`。
 - * **定义和训练 XGBoost 模型:**
 - * 设置参数,如 `objective='binary:logistic'`, `eval_metric`,

`n_estimators`, `learning_rate`, `max_depth`, `subsample`, `colsample_bytree`, `gamma` 等。**(第 8 章)**

- * 使用 `tree_method='hist'` 以提高效率。**(第 8, 13 章)**
- * 使用 `eval_set=[(X_val, y_val)]` 和 `early_stopping_rounds` 防止过拟 合。**(第 8 章)**
 - * **评估模型:**
 - * 在验证集 `X_val` 上使用 `predict_proba`。
 - * 遍历不同的概率阈值 `thresholds_to_evaluate`。
 - * 对每个阈值,计算二元预测,并调用 `calculate_metrics` 函数计算混淆矩阵、

Accuracy, POD (召回率), FAR, CSI 等指标。**(第 7 章)**

- * 使用 Pandas DataFrame 展示不同阈值下的性能对比。**(第 4 章)**
- * **保存模型:** 使用 `joblib.dump` 保存训练好的 FN 专家模型
- (`fn_expert_model.joblib`)。**(第 5 章)**
- * **脚本间联系:**
 - * 输入来自脚本 `1` (FN/TN 索引)。
- * 输出的 `fn_expert_model.joblib` **在后续提供的脚本中似乎没有被直接使用** (例如,脚本 `4` 没有加载它来生成元特征)。这可能是该工作流的一个待完善之处,或者 FN 专家模型

的效果不佳,或者设计者决定只关注 FP 的修正。

15.3 脚本 `3_train_fp_expert.py`

- * **目标:** 训练一个 XGBoost 模型,专门用于**识别基线模型误报 (FP) 的样本**。目标是提高整体的精确率或降低 FAR。
- * **核心思想: ** 模型试图在"实际未下雨但基模型预测下雨 (FP)"和"实际未下雨基模型也预测未下雨 (TN)"这两类样本中,学习区分前者的模式。**关键区别在于: ** 它利用了基线模型预测的概率作为额外信息。
- * **关键步骤与技术:**
 - * **加载索引:** 加载 `indices_fp.npy` 和 `indices_tn.npy`。
 - * **准备训练数据:**
 - * 合并 FP 和 TN 索引 `combined_indices`。
 - * 创建目标标签 `labels`: FP 样本标记为 1, TN 样本标记为 0。
 - * **加载基模型训练集概率: ** 加载 `base_probs_train.npy`。 **(第 3 章) **
- * **花式索引选择概率:** 使用 `combined_indices` 从 `base_probs_train_full` 中选取 FP 和 TN 样本对应的基模型概率 `base_probs_expert`。**(第 3 章)**
- * **花式索引加载原始特征:** 使用 `combined_indices` 从 `X_flat` 加载 FP 和 TN 样本的原始特征 `X_expert_data_orig`。**(第 3 章)**
- * **特征构造:** 使用 `np.concatenate` 将原始特征和基模型概率 (reshape 后) 合并,构成 FP 专家模型的最终输入特征 `X_expert_data` (维度变为 n_features + 1)。**(第3,9章)**
- * **划分训练/验证集:** 同脚本 `2`, 使用 `train_test_split` 和 `stratify=labels`。**(第 5, 11 章)**
- * **处理类别不平衡 (算法层面):** 计算并设置 `scale_pos_weight = num_tn / num_fp`。**(第 8, 12 章)**
- * **定义和训练 XGBoost 模型:** 同脚本 `2`, 设置参数、使用 `eval_set` 和早停。** (第 8 章)**
 - * **评估模型:** 同脚本 `2`, 在验证集上评估不同阈值下的性能指标。**(第 7 章)**
- * **保存模型:** 保存训练好的 FP 专家模型 (`fp_expert_model.joblib`)。**(第 5 章)**
- * **脚本间联系:**
 - * 输入来自脚本 `1` (FP/TN 索引, 训练集概率)。
- * 输出的 `fp_expert_model.joblib` 是脚本 `4` (元特征生成) 和 `evaluate_fp_expert_on_tp.py` 的关键输入。

15.4 脚本 `4_generate_meta_features.py`

- * **目标: ** 为 ** 测试集 ** 生成 ** 元特征 (Meta-Features) ** ,用于后续训练元学习器。
- * **核心思想:** 结合**基线模型**的预测和 **FP 专家模型**对误报可能性的判断,构建新的特征表示。
- * **关键步骤与技术:**
- * **加载数据标识符/形状:** 获取 `X_flat`, `Y_flat`, `base_probs_test` 的形状信息, 计算测试集大小 `n_test_samples`。
- * **加载测试集数据到内存(或映射):** 加载 `X_test`, `Y_test_raw`(并转换为 `y_test_true`), `base_probs_test`。脚本尝试直接加载,如果内存不足会失败。**(第 3, 13 章)**
- * 注意: 这里 `base_probs_test` 是直接从文件加载,没有再进行切片,假设它只包含测试集的概率。
 - * **加载 FP 专家模型:** `joblib.load(FP_EXPERT_MODEL_PATH)`。**(第 5 章)**
- * **分块处理 (Chunking):** 为了处理可能很大的测试集,使用循环按 `CHUNK_SIZE` 遍 历测试集样本。**(第 13 章)**
 - * **元特征生成(块内):**
 - * 获取当前块的原始特征 `X chunk` 和基模型概率 `base prob chunk`。

* **准备 FP 专家输入:** `np.concatenate((X_chunk, base_prob_chunk.reshape(-1, 1)), axis=1)`。**(第 3, 9 章)**

* **FP 专家预测:**

`fp_expert_model.predict_proba(X_fp_expert_input_chunk)[:, 1]`, 得到当前块"是 FP"的概率 `pred_fp_expert_chunk`。**(第 5, 8 章)**

* **组合元特征: ** 将 `base_prob_chunk` (基模型概率) 和

`pred_fp_expert_chunk` (FP 专家概率) 按列堆叠 `np.stack(..., axis=1)`。

- * 将当前块的元特征存入预分配的 `X meta` 数组中。
- * **保存元特征和标签:**
- * 保存测试集的元特征矩阵 `X_meta` (形状 `(n_test_samples, 2)` 到 `X_meta.npy`)。**(第 3 章)**
 - * 保存测试集的真实标签 `y_test_true` 到 `y_meta.npy`。**(第 3 章)**
- * **脚本间联系:**
 - * 输入来自脚本 `1` (测试集概率) 和脚本 `3` (FP 专家模型)。
 - * 输出的 `X_meta.npy` 和 `y_meta.npy` 是脚本 `5` (元学习器训练) 的输入。

15.5 脚本 `5_train_evaluate_meta_learner.py`

- * **目标:** 训练并评估**元学习器 (Meta-Learner)**, 得到 Stacking 集成模型的最终性能。
- * **核心思想: ** 利用前面生成的元特征 (基模型概率, FP 专家概率) 来学习如何做出最终的分类决策。
- * **关键步骤与技术:**
 - * **加载元数据:** 加载 `X_meta.npy` 和 `y_meta.npy`。**(第 3 章)**
- * **数据划分(简化):** 脚本中直接将加载的`X_meta`,`y_meta` 同时用作训练和评估数据(`X_meta_train = X_meta`,`y_meta_train = y_meta_true`, etc.)。这是一种简化,如前所述,可能导致性能评估略微乐观。更严格的做法是将`X_meta`,`y_meta` 再划分为元训练集和元测试集。**(第 5 章)**
 - * **选择并定义元学习器:**
 - * 提供两种选择: `LogisticRegression` 或 `XGBClassifier`

(`META_LEARNER_TYPE` 参数控制)。**(第 6, 8 章)**

- * **逻辑回归:** 设置 `solver`, `max iter`。**(第 6 章)**
- * **XGBoost:**
- * 计算并设置 `scale_pos_weight` (基于元标签 `y_meta_train` 的不平衡性)。**(第 8, 12 章)**
- * 使用较浅的树 (`max_depth=3`) 和较少的估计器 (`n_estimators=100`), 这是元学习器的常见做法, 防止过拟合元特征。**(第 8 章)**
- * 定义特征名称 `feature_names = ["Base_Prob", "FP_Expert_Prob"]` 用于 后续分析。
 - * **训练元学习器:** `meta_learner.fit(X_meta_train, y_meta_train)`。
- * **特征重要性分析 (XGBoost):** 如果元学习器是 XGBoost, 计算并打印/绘制 `feature_importances_`, 了解基模型概率和 FP 专家概率对最终预测的相对贡献。**(第 8 章)**
 - * **评估元学习器(即最终模型):**
 - * 使用 `meta_learner.predict_proba(X_meta_eval)` 获取最终预测概率。
 - * 遍历不同的概率阈值 `thresholds to evaluate`。
- * 计算每个阈值下的性能指标 (Accuracy, POD, FAR, CSI 等) 并用 Pandas DataFrame 展示。**(第 7 章)**
 - * **保存元学习器模型:** `joblib.dump(meta_learner,

META_LEARNER_MODEL_PATH)`。**(第 5 章)**

- * **脚本间联系:**
 - * 输入来自脚本 `4` (元特征 `X_meta`, 元标签 `y_meta`)。
 - * 输出训练好的元学习器模型和最终的集成模型性能评估。

15.6 脚本 `debug_predict_test.py`

- * **目标:** 提供一个简单的冒烟测试 (Smoke Test) 工具,用于验证保存的某个模型文件是否能够成功加载,并能够对符合预期形状的虚拟数据执行 `predict_proba` 方法,而不会引发错误。
- * **作用:** 在复杂的工作流中,模型保存和加载环节容易出错(例如,环境不兼容、文件损坏、模型对象结构问题)。这个脚本可以快速检查某个关键模型文件是否基本可用。
- * **关键步骤与技术:**
 - * **配置模型路径: ** `MODEL PATH` 变量允许用户指定要测试的模型文件。
 - * **打印环境信息: ** 输出 Python 和关键库的版本, 有助于调试兼容性问题。
 - * **加载模型 (`joblib.load`): ** 尝试加载指定路径的模型。 **(第 5 章) **
- * **创建虚拟数据:** 使用 `np.random.rand` 创建一个符合模型预期输入特征数量 (`N_FEATURES`) 的小型随机数据集。**(第 3 章)**
- * **调用 `predict_proba`:** 尝试使用加载的模型对虚拟数据进行预测。**(第 5,8 章)**
- * **异常捕获:** 使用 `try...except BaseException` 捕获加载或预测过程中可能出现的任何错误,并打印详细的错误信息和堆栈跟踪(`traceback`),方便定位问题。
- * **脚本间联系:** 这是一个独立的调试工具,可以用于测试脚本 1, 2, 3, 5, 6 输出的任何 `.joblib` 模型文件。
- ### 15.7 脚本 `6 finetune base model.py`
- * **目标: ** 探索**另一种**提升基模型性能的策略—**微调 (Fine-tuning) **。它直接在基模型出错的训练样本上继续训练基模型。
- * **核心思想: ** 强化基模型在"困难"样本 (原始的 FP 和 FN) 上的学习,以期提高整体性能,特别是改善错误分类的情况。
- * **关键步骤与技术:**
 - * **加载初始模型:** `joblib.load(INITIAL_MODEL_PATH)`。**(第 5 章)**
- * **加载 FP/FN 索引:** 加载脚本 `1` 输出的 `indices_fp.npy`, `indices_fn.npy`。
 - * **准备微调数据:**
 - * 合并 FP 和 FN 索引 `finetune_indices`。
- * **花式索引加载特征和标签:** 使用 `finetune_indices` 从 `X_flat`, `V_flat` (内存映射) 中提取特征 `X_finetune` 和原始目标 `v_finetune_raw` **(\$
- `Y_flat` (内存映射) 中提取特征 `X_finetune` 和原始目标 `y_finetune_raw`。**(第 3 章)**
 - * 将原始目标转换为二元标签 `y_finetune_true`。
 - * **定义微调参数:**
- * 设置较小的学习率 `FINETUNE_LR` 和较少的迭代次数 `FINETUNE_N_ESTIMATORS`。
 (第 14 章)
- * 根据微调数据集 `y_finetune_true` 重新计算 `scale_pos_weight`。**(第 8, 12, 14 章)**
- * **禁用早停** (`early_stopping_rounds=None`), 强制模型在微调数据上训练指定的轮数。**(第 8, 14 章)**
 - * **执行微调:**
 - * 创建新的 `XGBClassifier` 实例 `finetune_model`, 使用更新后的参数。
- * 调用 `finetune_model.fit()`, 关键在于传入 `xgb_model=initial_model`。** (第 8, 14 章)**
- * **保存微调模型:** `joblib.dump(finetune_model, FINETUNED_MODEL_PATH)`。** (第 5 章)**
 - * **评估微调模型(在原始测试集上):**
 - * 加载原始测试集的真实标签 `y_test_true`。
- * **分块预测:** 由于测试集可能很大,使用循环和内存映射分块加载测试特征
- `X_chunk`, 调用 `finetune_model.predict_proba` 进行预测, 并将结果累积到列表 `y_test_pred_proba_list`。**(第 13 章)**
 - * **合并结果: ** 使用 `np.concatenate` 合并所有块的预测概率。**(第 3 章) **

* **多阈值评估:** 遍历不同阈值, 计算性能指标 (Accuracy, POD, FAR, CSI 等) 并展示。**(第 7 章)**

- * **脚本间联系:**
 - * 输入来自脚本 `1` (FP/FN 索引, 初始模型路径)。
- * 输出是微调后的模型 `xgboost_finetuned_model.joblib` 和其在测试集上的性能评估。它提供了一条与 Stacking 不同的优化路径。

15.8 脚本 `evaluate_fp_expert_on_tp.py`

- * **目标:** 这是一个**分析性**脚本,旨在评估 **FP 专家模型** (脚本 3 训练的模型) 在** 真实的正类样本 (TP)** 上的行为。
- * **核心思想:** FP 专家模型的训练目标是识别"误报"(实际为 0 但预测为 1)。理想情况下,当输入一个"真实命中"(实际为 1 预测也为 1,即 TP)的样本时,FP 专家模型应该**不认为**它是 FP,即输出较低的"是 FP 的概率"。通过检查 FP 专家在 TP 样本上的概率分布,可以了解它是否会错误地将一些真实的命中识别为潜在的误报,这可能影响 Stacking 流程中元学习器的决策。
- * **关键步骤与技术:**
 - * **加载 TP 索引:** 加载脚本 `1` 输出的 `indices_tp.npy`。
 - * **加载 FP 专家模型:** `joblib.load(FP_EXPERT_MODEL_PATH)`。**(第 5 章)**
 - * **准备 TP 样本数据 (含基模型概率):**
 - * 加载 TP 样本的原始特征 `X_tp_data_orig`。**(第 3 章)**
 - * 加载训练集的基模型概率 `base_probs_train_full`。**(第 3 章)**
 - * 选取 TP 样本对应的基模型概率 `base_probs_tp`。 **(第 3 章)**
 - * **特征构造:** 合并原始特征和基模型概率,得到 FP 专家模型的输入

`X_tp_data`。**(第 3, 9 章)**

- * **FP 专家预测:** 使用加载的 `fp_expert_model` 对 `X_tp_data` 调用 `predict_proba`, 获取"是 FP 的概率" `tp_pred_proba_is_fp`。**(第 5, 8 章)**
 - * **分析概率分布:**
 - * 计算概率的基本统计量 (min, max, mean, median, std)。**(第 3 章)**
 - * 计算超过不同阈值的 TP 样本比例。**(第 3 章)**
- * **绘制直方图:** 使用 Matplotlib 可视化 `tp_pred_proba_is_fp` 的分布。** (第 2 章)**
- * **结果解读:** 如果直方图显示大量 TP 样本的"是 FP 概率"很高(例如,集中在 0.5 以上),说明 FP 专家模型可能区分度不够好,会将一些真实的命中也判断为潜在的误报。如果概率大多集中在较低的值,说明 FP 专家模型在这方面表现较好。
- * **脚本间联系:**
 - * 输入来自脚本 `1` (TP 索引,训练集概率) 和脚本 `3` (FP 专家模型)。
 - * 输出是分析结果 (统计数据、直方图) , 用于理解 FP 专家模型的行为和局限性。

15.9 整体工作流总结与思考

这 8 个脚本(加上假定的 `xgboost1.py`) 共同构成了一个相对复杂的机器学习工作流,其核心目标是通过一种定制化的 Stacking 变体(结合基模型和针对错误的专家模型)来提升天气预测(二分类)的性能。

关键流程:

基础模型训练 -> 基础模型错误分析 & 概率输出 -> 专家模型训练 (FN, FP - FP利用基础概率) -> 测试集元特征生成 (基础概率 + FP专家概率) -> 元学习器训练 & 评估。同时,提供了单独调试模型加载和微调基模型的备选路径,以及分析专家模型行为的工具。

体现的技术点:

- * **核心模型: ** XGBoost 被深度应用。
- * **集成学习: ** 实现了 Stacking 的变体。

* **数据处理:** NumPy (数组操作,npy文件,内存映射,索引),Joblib (模型持久化),Pandas (结果展示)。

- * **模型评估:** 多种指标 (Accuracy, POD, FAR, CSI), 混淆矩阵, 多阈值评估。
- * **不平衡处理: ** `scale_pos_weight` 的应用。
- * **大数据策略: ** 内存映射, 分块处理。
- * **模型优化: ** 微调 (Fine-tuning) 作为备选方案。
- * **工作流管理:** 清晰的文件路径配置和模块化脚本。

可思考和改进之处:

- * **FN 专家的利用: ** 当前流程似乎没有将 FN 专家的预测纳入元特征,可以考虑加入 (例如, `X_meta` 变为 3 列: Base_Prob, FN_Expert_Prob, FP_Expert_Prob) 并重新训练元学习器,看是否能进一步提升性能(特别是召回率/POD)。
- * **元特征生成方式: ** 当前直接在测试集元特征上训练/评估元学习器。可以考虑实现标准的 K-Fold Stacking 或 Blending 来生成训练元特征,以获得更可靠的性能评估。
- * **特征工程: ** 脚本假设特征工程已完成。实际项目中,对 `X_flat_features.npy` 中的特征进行更深入的探索、构造和选择可能会带来显著提升。
- * **超参数调优: ** 脚本中的 XGBoost 参数是固定的。对基模型、专家模型和元学习器进行系统的超参数调优(使用 `GridSearchCV`或 `RandomizedSearchCV`) 可能会找到更好的配置。
- * **基学习器多样性:** 当前 Stacking 主要围绕一个基模型及其变种。可以尝试引入完全不同类型的基学习器 (如 SVM,随机森林等) 到 Level Ø,看是否能通过更多样化的信息源提升元学习器的效果。

通过对这个案例的深入分析,你应该能够更清晰地看到前面章节学习到的各种技术是如何在实际项目中组合应用,以解决具体问题的。这也为你将来设计和实现自己的机器学习工作流提供了宝贵的参考。

机器学习实战: 从 Scikit-learn 到 XGBoost (卷六)

...(接上文)

- - -

第十六章: 模型部署与上线 (简介)

我们已经花费了大量篇幅讨论如何定义问题、准备数据、训练模型、评估性能以及优化超参数。然而, 机器学习项目的最终目标通常是将训练好的模型****部署** (Deploy)** 到实际生产环境中, 使其能够接收新的输入数据并提供预测服务, 从而真正创造价值。

模型部署是将机器学习模型集成到现有业务流程或应用程序中的过程。虽然本教程的重点在于模型的构建和评估,但了解部署的基本概念和常见方式对于形成完整的机器学习项目认知至关重要。

本章将简要介绍模型部署的一些关键考虑因素和常见模式。

16.1 为何需要模型部署?

训练好的模型本身(例如,一个`.joblib`或 XGBoost 模型文件)只是一个静态的文件。为了让它能够发挥作用,我们需要:

- * **提供接口: ** 创建一个可以让外部应用程序或用户发送数据并接收预测结果的方式。
- * **运行环境: ** 提供一个稳定、可靠的环境来加载模型、处理输入数据并执行预测计算。
- * **可扩展性:** 能够根据请求量的大小调整计算资源。

* ****监控与维护:**** 持续监控模型的性能、稳定性和资源使用情况。

16.2 部署的关键考虑因素

在将模型部署到生产环境之前,需要考虑以下因素:

- 1. **预测延迟要求 (Latency): ** 应用场景对预测响应时间的要求有多高?
- * ****实时预测** (Real-time / Online Prediction):** 需要在毫秒或秒级内返回结果(例如,在线欺诈检测、实时推荐)。
 - * ****近实时预测 (Near Real-time):** ** 响应时间要求稍宽松(例如,几秒到几分钟)。
- * **批量预测 (Batch Prediction):** 对一批数据进行一次性预测,通常在后台定期运行,对响应时间要求不高(例如,每日销售预测、每周客户流失预测)。
- 2. **预测吞吐量要求 (Throughput): ** 系统需要处理多少预测请求 (例如,每秒请求数 QPS)?
- 3. ****可用性与可靠性** (Availability & Reliability):** 预测服务需要多高的可用性? 能否容忍短暂中断?
- 4. **资源成本 (Cost): ** 部署和维护预测服务的硬件、软件和人力成本。
- 5. ****可维护性与更新** (Maintainability & Updates): ****** 更新模型(重新训练后部署新版本)的 流程是否方便?如何进行版本控制和回滚?
- 6. **安全性 (Security): ** 如何保护模型和预测接口不被滥用或攻击? 如何处理敏感数据?
- 7. ****监控与日志** (Monitoring & Logging): ** 如何监控模型的预测性能(例如,与真实结果对比)、资源使用、错误率等?如何记录请求和预测日志?

16.3 常见的模型部署模式

根据不同的需求和场景,有多种部署模式可供选择:

16.3.1 嵌入式部署 (Embedded Deployment)

- * ****方式:**** 将模型文件直接打包到应用程序 (例如,移动 App、桌面软件、甚至硬件设备)中。 预测在本地设备上执行。
- * **优点:**
 - * 低延迟(无需网络诵信)。
 - * 离线可用。
 - * 数据隐私性好(数据通常不离开本地设备)。
- * **缺点:**
 - * 模型更新困难 (需要更新整个应用程序)。
 - * 受限于本地设备的计算能力和存储空间。
 - * 难以进行集中的监控和管理。
- * **适用场景: ** 对延迟要求极高、需要离线运行、模型相对较小的场景(例如,移动端的图像识别、智能设备的语音命令识别)。

16.3.2 API 服务部署 (API Service Deployment) - 最常用

- * **方式:** 将模型封装在一个独立的 Web 服务中,通过 API (通常是 RESTful API) 对外提供预测接口。应用程序通过发送 HTTP 请求 (包含输入数据) 到 API 端点,并接收包含预测结果的响应。
- * **优点:**
 - * **解耦: ** 模型服务与应用程序逻辑分离,易于独立开发、更新和扩展。
 - * **集中管理:** 模型更新、监控和维护都在服务端进行。
 - * **语言无关: ** 任何能够发送 HTTP 请求的应用程序都可以调用该服务。
 - * **可扩展性好:** 可以通过增加服务器实例或使用负载均衡来处理高并发请求。
- * **缺点:**
 - * 引入网络延迟。
 - * 需要维护额外的服务基础设施。

需要考虑 API 的安全性、认证和授权。

放

- **Web 框架:** 使用 Python Web 框架 (如 Flask, FastAPI, Django) 来构建 API 服 务。
 - **Flask:** 轻量级框架,易于上手,适合快速构建简单的 API。
- **FastAPI:** 现代、高性能框架,基于 ASGI,支持异步操作,自动生成 API 文 档,性能优异,推荐用于构建生产级 API。
 - **Django:** 功能全面的框架,更重,适合构建大型 Web 应用,但也常用于 API。
 - **模型加载:** 在服务启动时加载模型文件 (`joblib.load` 或原生加载)。
- ****请求处理:**** API 端点接收请求数据(通常是 JSON 格式),进行必要的预处理(可能 需要与训练时的预处理步骤一致),调用加载的模型进行 `predict` 或 `predict_proba`,将结果 格式化 (例如 JSON) 并返回。
- ****部署:**** 可以将服务部署到虚拟机、容器 (Docker)、Serverless 平台或专门的机器学 习平台 (如 AWS SageMaker, Google AI Platform, Azure Machine Learning) 。
- ****适用场景:**** 大多数需要将模型集成到 Web 应用、后端服务或其他系统的场景。

简单的 Flask API 示例 (示意):

```
```python
from flask import Flask, request, jsonify
import joblib
import numpy as np
import pandas as pd # 如果需要用 Pandas 处理输入
--- 全局加载模型 (服务启动时加载一次) ---
try:
 model = joblib.load('path/to/your/final_model.joblib') # 替换为最终模型路径(可
能是元学习器)
 # scaler = joblib.load('path/to/your/scaler.joblib') # 如果需要加载预处理器
 print("Model loaded successfully.")
except Exception as e:
 print(f"Error loading model: {e}")
 model = None # 或者退出程序
app = Flask(name)
@app.route('/predict', methods=['POST'])
def predict():
 if model is None:
 return jsonify({'error': 'Model not loaded'}), 500
 try:
 # 1. 获取输入数据 (假设为 JSON 格式)
 data = request.get json(force=True)
 # 示例: {'features': [f1, f2, f3, ...]}
 features = np.array(data['features']).reshape(1, -1) # 转换为 NumPy 数组
 # 2. (可选)数据预处理(必须与训练时一致!)
 # features_scaled = scaler.transform(features) # 如果使用了缩放器
 # 3. 模型预测 (假设使用缩放后的数据)
 # prediction_proba = model.predict_proba(features_scaled)[0] # 获取概率
 prediction_proba = model.predict_proba(features)[0] # 如果模型训练时未使用缩
```

```
positive_class_proba = float(prediction_proba[1]) # 假设是二分类,取正类概率

也可以直接获取类别预测
prediction_class = int(model.predict(features_scaled)[0])

4. 格式化并返回结果
result = {
 'predicted_probability_class_1': positive_class_proba
 # 'predicted_class': prediction_class
 }
 return jsonify(result)

except Exception as e:
 # 记录错误日志
 app.logger.error(f"Prediction error: {e}")
 return jsonify({'error': str(e)}), 400 # 返回错误信息

if __name__ == '__main__':
 # 启动 Flask 开发服务器 (生产环境需要使用 Gunicorn/uWSGI 等)
 app.run(host='0.0.0.0', port=5000, debug=False) # debug=False 用于生产
```

# (注意: 这是一个非常基础的示例, 实际生产部署需要考虑并发、错误处理、日志、安全、性能优化等更多因素。)

#### 16.3.3 批量预测部署 (Batch Scoring)

- **方式**: 创建一个独立的脚本或作业,该作业定期(例如,每天晚上)运行,从数据存储(如数据库、数据湖、文件系统)中读取一批新的输入数据,使用模型进行预测,并将预测结果写回到另一个数据存储或生成报告。
- 优点:
  - 。 实现相对简单。
  - 。 可以处理大量数据。
  - 。 对计算资源的要求通常不如实时服务高 (可以在非高峰时段运行)。
- 缺点:
  - 。 无法提供实时预测。
  - 。 结果具有滞后性。
- **实现**: 通常是一个 Python 脚本,使用 Pandas 或 Dask/Spark 读取数据,加载模型,进行预测(可能需要分块处理),然后将结果写回。可以使用 Airflow, Prefect, Luigi 等工作流管理工具来调度和管理这些批量作业。
- 适用场景: 不需要实时结果的场景, 如客户流失预测、销售预测、离线推荐生成等。

#### 16.3.4 流式处理部署 (Streaming Deployment)

- **方式:** 将模型集成到流处理系统(如 Apache Kafka + Spark Streaming, Flink, Kafka Streams)中。模型对实时到达的数据流进行预测。
- 优点:
  - 。 能够处理持续不断的数据流,实现近乎实时的预测。
- 缺点:
  - 。 系统搭建和维护复杂。

- 。 需要处理状态管理、容错等问题。
- 适用场景: 需要对高速数据流进行实时或近实时分析的场景,如实时异常检测、物联网设备数据分析。

#### 16.4 MLOps: 机器学习运维

将模型成功部署只是开始。为了确保模型在生产环境中持续、稳定、有效地运行,需要引入 **MLOps (Machine Learning Operations)** 的理念和实践。MLOps 结合了机器学习 (ML)、数据工程 (Data Engineering) 和 DevOps (开发运维一体化) ,旨在:

- 自动化: 实现模型训练、测试、部署、监控流程的自动化。
- 版本控制: 对数据、代码和模型进行版本控制。
- 持续集成/持续部署 (CI/CD): 为机器学习模型建立 CI/CD 流水线,实现快速、可靠的模型更新。
- 监控: 持续监控模型的预测性能(准确率、漂移等)、运行状况(延迟、错误率)和资源消耗。
- 可复现性: 确保模型训练和预测过程的可复现性。
- 协作: 促进数据科学家、机器学习工程师、运维工程师之间的协作。

MLOps 是一个庞大且仍在快速发展的领域,涉及许多工具和平台(如 MLflow, Kubeflow, Seldon Core, 以及各大云厂商提供的 MLOps 服务)。对于复杂的、需要频繁更新的生产级机器学习系统,采用 MLOps 实践至关重要。

#### 16.5 总结

模型部署是将机器学习研究成果转化为实际业务价值的关键一步。选择合适的部署模式取决于应用的具体需求(延迟、吞吐量、可用性、成本等)。API 服务部署是最常用和灵活的方式,而批量预测则适用于非实时场景。无论采用哪种方式,都需要考虑可维护性、监控和 MLOps 实践,以确保模型的长期稳定运行和持续改进。

虽然本教程不深入 MLOps 的具体实现细节,但理解部署的基本概念和模式,可以帮助你在设计和构建机器学习模型时,就考虑到后续部署的需求,使整个项目流程更加顺畅。

# 第十七章: 总结与展望

恭喜你!你已经跟随本教程(这本"小书")走过了机器学习从入门到实践的漫长而充实的旅程。我们从最基础的机器学习概念出发,系统性地学习了:

- 核心工具: Python 环境搭建,以及数据处理的基石 NumPy 和 Pandas。
- **Scikit-learn:** 掌握了其一致性的 Estimator API,常用的数据划分、预处理、模型评估方法,以及多种经典的分类算法(逻辑回归、KNN、SVM、决策树、随机森林)。
- XGBoost: 深入理解了梯度提升的思想、XGBoost 的优势、关键参数调优、早停机制、特征重要性评估以及模型微调。
- **高级技术**: 探讨了集成学习的进阶方法 Stacking,处理类别不平衡的各种策略,以及应对大型数据集的内存映射和分块处理技巧。
- **模型选择与部署**: 学习了使用交叉验证、网格搜索、随机搜索进行模型调优和选择的流程,并对模型部署的基本概念和模式有所了解。
- **案例分析**: 通过详细解读天气预测模型的系列脚本,将前面学习到的知识点串联起来,理解它们在实际工作流中的应用。

#### 现在, 你应该具备了:

- 对监督学习(特别是分类问题)的完整工作流有了清晰的认识。
- 能够熟练使用 NumPy 和 Pandas 进行基本的数据操作和准备。
- 能够自信地运用 Scikit-learn 的核心功能来构建、训练和评估常见的机器学习模型。
- 能够深入理解并有效使用强大的 XGBoost 库,包括参数调优和微调。
- 能够理解并分析(甚至尝试实现)像 Stacking 这样的高级集成策略。
- 能够识别并应用处理类别不平衡和大型数据集的基本方法。
- 具备了阅读、理解和分析类似你提供的复杂机器学习脚本的能力。

**回顾你的初始目标:** 达到能够写出那 8 个 Python 脚本的水平。通过本教程的学习,你已经掌握了实现这些脚本所需的所有核心概念、库的使用方法和关键技术点。当然,从"理解"到"熟练编写"还需要大量的实践。

#### 展望未来: 持续学习与实践之路

机器学习是一个日新月异、永无止境的领域。本教程为你打下了坚实的基础,但前方的道路依然广阔:

#### 1. 动手实践:

- 重现案例: 尝试在你自己的机器上运行(或模拟运行)教程中分析的天气预测脚本,修改参数,观察结果变化。
- 。 **应用到新数据**: 寻找你感兴趣的数据集(Kaggle, UCI 机器学习仓库, 政府公开数据等),应用本教程学到的知识,完整地走一遍机器学习流程。
- o 参加竞赛: Kaggle 等平台提供了大量的练手机会和高质量的学习资源。

#### 2. 深化特定领域:

- **回归问题:** 本教程侧重分类,你可以系统学习线性回归、岭回归、Lasso、决策树回归、SVR、XGBoost 回归等模型及其评估指标(MSE, MAE, R²)。
- **无监督学习:** 探索聚类 (K-Means, DBSCAN)、降维 (PCA, t-SNE)、异常检测等技术。
- 自然语言处理 (NLP): 学习文本表示 (词袋、TF-IDF、Word Embeddings like Word2Vec, GloVe)、文本分类、情感分析、序列模型 (RNN, LSTM, Transformer)。
- **计算机视觉 (CV):** 学习图像处理基础、卷积神经网络 (CNN)、目标检测、图像分割。
- o 时间序列分析: 学习 ARIMA, Prophet, 深度学习模型 (RNN, LSTM) 等用于预测时间序列数据。

#### 3. 探索更前沿的模型与技术:

- 。 LightGBM: 另一个高效的梯度提升库,通常比 XGBoost 更快,内存占用更低。
- 。 CatBoost: 对类别特征处理非常友好的梯度提升库。
- 深度学习框架: 学习 TensorFlow/Keras 或 PyTorch, 构建神经网络模型。
- 可解释性机器学习 (XAI): 学习 SHAP, LIME 等技术来理解"黑盒"模型的预测依据。
- AutoML: 了解自动化机器学习工具,它们可以自动完成特征工程、模型选择和超参数调优的部分工作。
- o MLOps: 深入学习模型部署、监控、版本控制、自动化流水线等运维实践。

#### 4. 阅读与交流:

- 阅读经典书籍和论文: 跟讲行业最新讲展。
- o 查阅官方文档: Scikit-learn, Pandas, NumPy, XGBoost 等库的官方文档是最好的学习资源。
- 参与社区: 在 Stack Overflow, GitHub, Reddit (r/MachineLearning), 专业论坛上提问、回答、交流。
- 关注博客和教程: 许多优秀的博客和在线课程持续分享着最新的知识和实践。

#### 结语:

机器学习的旅程充满挑战,但也同样充满乐趣和机遇。最重要的是保持好奇心,勇于实践,不断学习。希望这本"小书"能够成为你机器学习道路上的一个有力的起点和可靠的参考。

 <b>绝知此事要躬行。</b> 上取得丰硕的成果	-	见在,;	是时候动手去创造你自己的机器学习解决方案了!	祝你在这