

时间序列监督学习

汪小圈

2025-04-27

时间序列监督学习的基本概念

- 时间序列数据是按时间顺序收集的一系列数据点
 - 股票价格、气象数据、用户行为记录等
- 时间序列监督学习是利用历史时间序列数据预测未来值的机器学习任务
- 与传统监督学习不同，时间顺序和时间依赖性成为关键考量因素

时间序列监督学习的步骤：数据收集与预处理

- **数据收集**：获取包含时间戳的序列数据
- **数据清洗**：处理缺失值、异常值和噪声
- **数据标准化**：对数据进行归一化或标准化处理
- **特征工程**：
 - 滞后特征 (Lag Features)：使用过去的观测值作为特征
 - 窗口统计特征：计算滑动窗口内的统计量（均值、方差等）
 - 时间特征：提取时间戳中的周期性信息（小时、星期几、月份等）
 - 差分特征：计算数据点之间的差值，消除趋势

特殊的样本分割方法

时间序列数据的样本分割与传统监督学习有显著区别：

- **时间序列分割原则：** 必须严格按照时间顺序分割数据
 - 确保训练数据在测试数据之前
 - 未来的数据在实际预测时是不可获取的
- **传统随机分割的问题：**
 - 可能导致数据泄露
 - 模型可能利用未来信息预测过去
 - 导致过于乐观的评估结果

时间序列数据分割方法：扩展窗口法

扩展窗口法 (Expanding Window): 使用所有历史数据进行训练，预测窗口不断向前移动

第 1 轮:

- 训练集: [月份 1-3]
- 测试集: [月份 4]

第 2 轮:

- 训练集: [月份 1-4]
- 测试集: [月份 5]

第 3 轮:

- 训练集: [月份 1-5]
- 测试集: [月份 6]

优势:

- 充分利用所有历史数据
- 随着时间推移，模型可获得更多训练数据
- 适合长期趋势预测和季节性数据

时间序列数据分割方法：滑动窗口法

滑动窗口法 (Sliding Window): 使用固定长度的数据窗口进行训练和测试

第 1 轮:

- 训练集: [月份 1-3]
- 测试集: [月份 4]

第 2 轮:

- 训练集: [月份 2-4]
- 测试集: [月份 5]

第 3 轮:

- 训练集: [月份 3-5]
- 测试集: [月份 6]

优势:

- 适合捕捉近期模式和趋势
- 减少较旧数据的影响
- 计算效率高, 训练集大小恒定

时间序列数据分割方法：多步预测分割

多步预测分割：专注于一次性预测多个未来时间点

- 训练集: [月份 1-6]
- 测试集: [月份 7, 8, 9] // 一次性预测多个未来月份

实现方式：

- ① **直接法：**构建单一模型直接预测多个时间点的值
- ② **递归法：**使用单步预测模型，将前一步的预测结果作为下一步的输入特征

优势：

- 适合需要长期规划的场景
- 可以捕捉时间序列的长期依赖性
- 评估模型在不同预测步长上的性能

各分割方法的对比与选择

分割方法	训练数据量	计算成本	适用场景	主要优势
扩展窗口法	不断增加	高	长期稳定数据、季节性数据	充分利用所有历史数据
滑动窗口法	固定	低	快速变化的时间序列	专注于最新趋势，减少旧数据影响
多步预测分割	固定	中等-高	需要长期预测的应用	评估模型在多个未来时间点的预测能力

- 适合时间序列的模型：
 - 传统统计模型：ARIMA、指数平滑法
 - 机器学习模型：随机森林、XGBoost、LSTM、Transformer
 - 混合模型：结合统计模型和机器学习模型的优点
- 训练注意事项：
 - 考虑时间依赖性
 - 避免数据泄露（不使用未来信息）
 - 合理设置预测步长（短期或长期预测）

- 评估指标：
 - MSE、RMSE、MAE：衡量预测误差
 - MAPE：相对误差百分比
 - 方向准确率：预测趋势变化的正确率
- 超参数调优：
 - 使用时间序列交叉验证而非随机交叉验证
 - 考虑时间窗口大小、滞后阶数等特有超参数

- 滚动预测：定期使用新数据更新模型
- 模型监控：检测数据分布变化、模型漂移
- 自适应更新：根据预测表现动态调整模型

时间序列监督学习与传统监督学习的区别

方面	时间序列监督学习	传统监督学习
数据依赖性	数据点间存在时间依赖关系	假设数据点独立同分布
样本分割	严格按时间顺序分割，不可随机打乱	通常随机分割，可以打乱顺序
交叉验证	使用时间序列交叉验证方法	使用 K 折交叉验证等随机方法
特征工程	重视滞后特征、时间特征和窗口特征	关注静态特征和特征间关系
数据泄露	易发生时序泄露（使用未来信息）	主要关注特征泄露

挑战与解决方案

- **非平稳性**：数据分布随时间变化
 - 解决方案：差分、移动平均、时间分解
- **季节性**：数据存在周期性模式
 - 解决方案：季节性分解、季节性特征
- **长期依赖**：当前预测可能依赖于很久之前的数据
 - 解决方案：使用 LSTM、注意力机制等处理长期依赖
- **多变量预测**：多个相关时间序列的联合预测
 - 解决方案：多变量模型、变量间关系建模

- **股票价格预测**：使用历史价格、交易量和技术指标预测未来价格走势
- **电力负荷预测**：基于历史用电量和天气数据预测未来用电需求
- **疾病传播预测**：根据历史疫情数据预测未来感染人数
- **产品需求预测**：利用历史销售数据预测未来产品需求

选择合适的时间序列分割方法

① 考虑数据特性：

- 快速变化的数据（如股票价格）→ 滑动窗口
- 长期稳定的数据（如宏观经济指标）→ 扩展窗口

② 考虑预测目标：

- 短期预测（1-2 步）→ 滑动窗口
- 长期预测（多步）→ 多步预测分割

③ 考虑计算资源：

- 资源有限 → 滑动窗口（训练集大小固定）
- 资源充足 → 扩展窗口（可利用更多历史数据）

- 时间序列监督学习与传统监督学习最大的区别在于对数据时间依赖性的处理
- 在样本分割、特征工程、模型选择和评估方面都需要考虑时间顺序的影响
- 避免数据泄露，捕捉时间依赖模式是关键
- 根据具体应用场景选择适当的分割方法：
 - 扩展窗口法：长期稳定数据
 - 滑动窗口法：快速变化的数据
 - 多步预测分割：长期预测需求
- 掌握合适的时间序列分析技术对于处理实际业务中的预测问题至关重要