

核心课 ③ · 2020. 11. 07.

65

作业回顾: XGBoost 适合特征多, 特征2程比较差的情况下, 其敏感度高。

LightGBM ←

Cat Boost: 如分类特征多, 适合。

★ 模型融合: Stacking bagging, boosting.

训练集 → 训练, 同域。

Voting / Averaging <sup>平均</sup>

LightGBM 进阶, 需调整参数:  $\text{max-depth} = -1$   
 $\text{num-leaves} = 30$   
 $\text{learning-rate} = 0.2$

LightGBM stacking  $\text{model.fit}(\dots, \text{early-stopping-rounds} = 10)$  <sup>学习率大, 收敛快</sup>

K-Fold

1-层模型 → 训练 → 验证 → 测试。

此时用 "Averaging" 方法, 会读每个测试集等效。

$\Rightarrow w_1 \cdot \text{out}_1 + w_2 \cdot \text{out}_2 + \dots + w_n \cdot \text{out}_n$  ← 线性加权方法  
权重  $w_1, w_2, \dots, w_n$

用固定框架解决问题:

读取数据 → 数据预处理 →

特征工程 → 筛选特征 → 固定特征 →

~~固定特征~~ 调参 → 运行模型 → 模型融合

例子: 如基于特征工程-筛选-固定特征-调参-运行模型。

↑  
特征工程  
2.4. StandardScaler

~~特征工程~~  
cat-col



类别 ysv → label encode

xgboost

直接指定到类别特征

lgb/catboost → data[i].astype('category')

213.103 3个模型

样本量

subsample=0.75

注意!

learning-rate=0.1

xgb-model: objective='binary:logistic'  
tree-method='gpu\_hist' max-depth=6  
lgb-model: objective='binary' metric='auc' max-length=-1  
cat-model: task-type='gpu' eval-metric='auc' max-length=7  
iterations=1000 subsample=0.75

回归评价指标:

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

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

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

$$R^2 = 1 - \frac{\sum (y - \hat{y})^2}{\sum (y - \bar{y})^2} = 1 - \frac{SS_{residual}}{SS_{total}}$$

0: 没学好  
1: 学得好

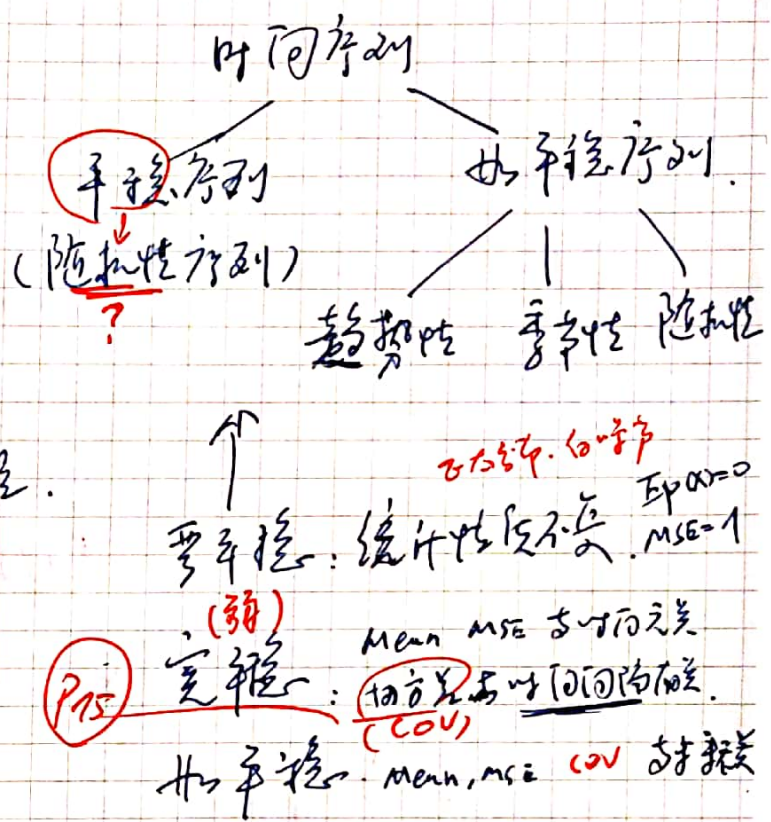


时间序列的概述：根据已有的时间序列数据预测未来值。

- 趋势的延续性。
- 未来的不确定性。(数据不规则)
- 不考虑因果关系。

预测步骤：

1. 数据平稳性检测。
2. 若不平稳，则差分处理。<sup>diff</sup>
3. 确定最佳模型。
4. 模型预测与检验。



ARIMA 模型：描述历史值与当前值之间的关系。同变量自身的历史时间数据对自身进行预测。

差分 integrated

(AR) I MA  
自动回归模型  
移动平均

autoregression

滞后项

$$y_t = \mu + \sum_{i=1}^p \gamma_i y_{t-i} + \epsilon_t$$

$\gamma_i$  自相关系数  $\epsilon_t$  误差

$$y_t = \mu + \epsilon_t + \sum_{i=1}^q \theta_i \epsilon_{t-i}$$

Moving Average, 是误差项的累加。消除随机波动。

ARIMA

$$y_t = \mu + \sum_{i=1}^p \gamma_i y_{t-i} + \epsilon_t + \sum_{i=1}^q \theta_i \epsilon_{t-i}$$

检验自相关性：统计特征 (MSE, Mean) 确定  $(\rho_i \geq 0.5)$   
与自身相关

P, Q 自定  
 $\gamma_i, \theta_i$  系数  
 $\mu, \epsilon_t$  参数

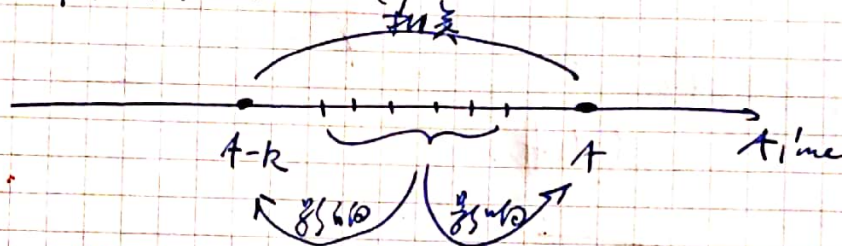


ACF: 自相关函数 (Autocorrelation function)

$$ACF(k) = \rho_k = \frac{Cov(y_t, y_{t-k})}{Var(y_t)}$$

↑  
[-1 ~ 1]

PACF: 偏自相关函数 (partial autocorrelation function)



定义

剔除P 得到真正的 "ACF"

	ACF	PACF
AR(p)	~0	P后截尾
MA(p)	P后截尾	~0
ARMA(p, q)	P后截尾	P后截尾

→ 确定 I, (0, 1, 2)

什么是差分?

模型残差检测: ① ARIMA模型 ② 是否为常数的正态分布

② QQ图: 线性 (即正态分布)

统计子模型

from Statsmodels. tsa. arima-model import ARIMA.

from " " . graphics. Asplots import plot-acf, plot-pacf.