

检测模型算法

方法类别	方法名称	异常标准简述	备注
基于分布	3σ	值 $> \mu + 3\sigma$ 或 值 $< \mu - 3\sigma$	
	Boxplot	值 $> q3 + 1.5 \text{IQR}$ 或 值 $< q1 - 1.5 \text{IQR}$	
	Grubbs	$z_score > \text{Grubbs临界值}$	
	weighted_avg	值 $> \mu + 3\sigma$ 或 值 $< \mu - 3\sigma$	计算残差的均值和标准差
基于时间序列	ewma	预测值在区间范围之内	
	arima	预测值在区间范围之内	

代码块

```

1  # ThreeSigma
2  # 传入的数据不包括需要检测的数据
3  class OutlierDetectorThreeSigma:
4      def __init__(self, data, threshold=3.0):
5          self.data = np.array(data)
6          self.mean = np.mean(data)
7          self.std_dev = np.std(data)
8          self.threshold = threshold
9
10     def is_outlier(self, value):
11         """
12             判断一个值是否为异常值。
13             如果值落在均值±3*标准差之外，则认为是异常值。
14         """
15         lower_bound = self.mean - self.threshold * self.std_dev
16         upper_bound = self.mean + self.threshold * self.std_dev
17         return value < lower_bound or value > upper_bound
18
19     def get_up_down_level(self):
20         lower_bound = self.mean - self.threshold * self.std_dev

```

```
21         upper_bound = self.mean + self.threshold * self.std_dev
22     return round(lower_bound, 2), round(upper_bound, 2)
23
24 # BoxPlot (箱线图)
25 # 传入的数据不包括需要检测的数据
26 class OutlierDetectorBoxPlot:
27     def __init__(self, data, threshold=1.5):
28         self.data = np.array(data)
29         self.threshold = threshold
30
31     def detect_anomaly_data(self, new_data):
32         """
33             判断一个值是否为异常值。
34             如果值落在下四分位数和上四分位数之外，则认为是异常值。
35         """
36         lower_bound, upper_bound = self.get_up_down_level()
37         if new_data < lower_bound or new_data > upper_bound:
38             return True
39         return False
40
41     def get_up_down_level(self):
42         # 获取箱线图的特征
43         """
44             使用 numpy.percentile 函数计算数据的第25百分位数（下四分位数）、第50百分位数
45             （中位数）和第75百分位数（上四分位数）。
46             结果存储在 quartiles 数组中
47             """
48         quartiles = np.percentile(self.data, [25, 50, 75])
49         """
50             计算四分位距，四分位距 (Interquartile Range, IQR) 是上四分位数和下四分位数之间
51             的差值。
52             它是一个衡量数据分散程度的统计量。
53             """
54         iqr = quartiles[2] - quartiles[0]
55         # 计算上下界
56         lower_bound = quartiles[0] - self.threshold * iqr
57         upper_bound = quartiles[2] + self.threshold * iqr
58         return round(lower_bound, 2), round(upper_bound, 2)
59
60 # Grubbs (格鲁布斯检验)
61 # 基于假设检验的思想，通过计算一个统计量来判断数据集中是否存在显著偏离其他数据点的异常值
62 # 传入的数据包括需要检测的数据
63 class OutlierDetectorGrubbs:
64     def __init__(self, data, alpha=0.05):
65         # 设置显著性水平
66         self.alpha = alpha
67         self.data = data
```

```
66     # 计算均值
67     self.mean = np.mean(self.data)
68     # 计算标准差
69     self.std_dev = np.std(self.data, ddof=1)
70     # 计算最大偏差
71     self.max_deviation = np.max(np.abs(self.data - self.mean))
72     # 计算 Grubbs 统计量
73     self.grubbs_statistic = self.max_deviation / self.std_dev
74     # 计算Grubbs的临界值
75     self.critical_value = self.calculate_critical_value(len(self.data))
76
77 def calculate_critical_value(self, n):
78     # 计算Grubbs' Test的临界值，基于样本大小n和显著性水平alpha
79     """
80         计算 t 值
81         t.ppf 是 SciPy 库中的累积分布函数 (CDF) 的逆函数，用于计算 t 分布的临界值。
82         1 - self.alpha / (2 * n) 是用于计算双尾检验的显著性水平。
83         n - 2 是自由度，因为 Grubbs' Test 使用的是 t 分布，自由度为 n - 2
84     """
85     t_value = t.ppf(1 - self.alpha / (2 * n), n - 2)
86     # 计算 Grubbs 临界值
87     critical_value = ((n - 1) * np.sqrt(t_value**2)) / (np.sqrt(n) *
88     np.sqrt(n - 2 + t_value**2))
89     return critical_value
90
91 def test(self):
92     """
93         如果 Grubbs 统计量大于临界值，则认为数据集中存在异常值。
94         使用 np.argmax(np.abs(self.data - mean)) 找到最大偏差对应的数据点的索引
95         根据索引返回异常值
96     """
97     if self.grubbs_statistic > self.critical_value:
98         outlier_index = np.argmax(np.abs(self.data - self.mean))
99         return True, data[outlier_index]
100    else:
101        return False, None
102
103 def get_up_down_level(self, data):
104     """
105         找出数据中的所有上下区间。
106     """
107     lower_bound = self.mean - self.critical_value * self.std_dev
108     if lower_bound < 0.0:
109         lower_bound = 0.0
110     upper_bound = self.mean + self.critical_value * self.std_dev
111     return round(lower_bound, 2), round(upper_bound, 2)
```

```
112     # WeightedAverage (加权平均)
113     # 传入的数据包括需要检测的数据
114     class OutlierDetectorWeightedAverage:
115         def __init__(self, weights=None):
116             self.weights = weights if weights is not None else None
117             self.data = None
118             self.weighted_averages = None
119
120         def fit(self, data, weights=None):
121             self.data = np.array(data)
122             # 设置每个输入数据的权重，默认都为1
123             self.weights = np.array(weights) if weights is not None else None
124             if self.weights is None:
125                 self.weights = np.ones_like(self.data)
126             # 使用选定的权重计算每个输入数据的加权平均值
127             self.weighted_averages = np.average(self.data, weights=self.weights)
128
129         def detect_anomalies(self, threshold=3):
130             # 计算每个数据的实际值与加权平均值之间的残差
131             deviations = np.abs(self.data - self.weighted_averages)
132             # 计算所有残差的标准差
133             standard_deviation = np.sqrt(np.average((deviations -
134                                         np.average(deviations))**2, weights=self.weights))
135             # 如果某个数据的残差超出定义的阈值，则认为该数据是异常的
136             anomalies = np.where(deviations > threshold * standard_deviation)[0]
137             return anomalies.tolist()
138
139         def get_up_down_level(self, threshold):
140             """
141                 找出数据中的所有上下区间。
142             """
143             lower_bound = self.weighted_averages - threshold *
144             self.standard_deviation
145             if lower_bound < 0.0:
146                 lower_bound = 0.0
147             upper_bound = self.weighted_averages + threshold *
148             self.standard_deviation
149             return round(lower_bound, 2), round(upper_bound, 2)
150
151         # EWMA (Exponentially Weighted Moving Average, 指数加权移动平均)
152         # 传入的数据不包括需要检测的数据
153         class OutlierDetectorEWMA:
154             def __init__(self, alpha=0.2, threshold=3):
155                 self.alpha = alpha
156                 self.threshold = threshold
157                 self.ewma = None
158                 self.anomalies = None
```

```
156         self.deviation_var = 0.0
157
158     def fit(self, data):
159         self.data = np.array(data)
160         self.ewma = np.zeros_like(self.data)
161         # 每个时间点的EWMA值
162         self.ewma[0] = self.data[0]
163         for i in range(1, len(self.data)):
164             # EWMA公式  $St=\alpha xt + (1-\alpha)St-1$ 
165             self.ewma[i] = self.alpha * self.data[i] + (1 - self.alpha) *
166             self.ewma[i-1]
167             deviations = np.abs(self.data - self.ewma)
168             self.deviation_var = deviations[1:].std()
169
170     def is_outlier(self, new_data):
171         """
172             判断一个值是否为异常值。
173         """
174         # 计算新传入的时间点的EWMA值
175         new_data_alpha = self.alpha * new_data + (1 - self.alpha) *
176         self.ewma[-1]
177         # 计算新传入的时间点的实际值与计算得出的EWMA值之间的残差
178         deviation = np.abs(new_data - new_data_alpha)
179         # 如果新传入的时间点的残差超出定义的阈值，则认为新传入的时间点的数据是异常的
180         if deviation <= self.deviation_var * self.threshold:
181             return False
182         return True
183
184     def get_up_down_level(self, new_data):
185         """
186             找出数据中的所有上下区间。
187         """
188         # 计算新传入的时间点的EWMA值
189         new_data_alpha = self.alpha * new_data + (1 - self.alpha) *
190         self.ewma[-1]
191         lower_bound = new_data_alpha - self.threshold * self.deviation_var
192         if lower_bound < 0.0:
193             lower_bound = 0.0
194         upper_bound = self.ewma[-1] + self.threshold * self.deviation_var
195         return round(lower_bound, 2), round(upper_bound, 2)
196
197     # ARIMA (AutoRegressive Integrated Moving Average, 自回归积分滑动平均模型)
198     # 传入的数据不包括需要检测的数据
199     class OutlierDetectorARIMA:
200         def __init__(self, data, order=(1, 1, 1), threshold=3.0):
201             """
202                 初始化ARIMA异常检测器
203             
```

```
200     使用历史数据拟合ARIMA模型。
201     Parameters:
202         order (tuple): ARIMA模型的(p, d, q)参数
203         确定模型的参数 (p, d, q) , 其中:
204             p 是自回归部分的阶数。
205             d 是差分的次数。
206             q 是移动平均部分的阶数。
207         threshold (float): 用于判断异常的阈值
208         """
209         self.order = order
210         self.threshold = threshold
211         self.model = None
212         self.fitted_values = None
213         self.data = data
214         self.std_residuals = 0.0
215
216     # 根据AIC和BIC可以用于比较不同 ARIMA 模型 (即不同的 (p, d, q) 参数组合)
217     def get_best_pdq(self):
218         # 定义 p, d, q 的范围
219         p_values = range(0, 3)
220         d_values = range(0, 2)
221         q_values = range(0, 3)
222
223         best_aic = float("inf")
224         best_bic = float("inf")
225         best_total_aic_bic = float("inf")
226         best_order = None
227         best_order_bic = None
228         best_order_aic_bic = None
229
230         for p in p_values:
231             for d in d_values:
232                 for q in q_values:
233                     try:
234                         # 拟合 ARIMA 模型
235                         model = ARIMA(self.data, order=(p, d, q))
236                         model_fit = model.fit()
237
238                         # 记录 AIC 和 BIC 值
239                         aic = model_fit.aic
240                         bic = model_fit.bic
241
242                         total_aic_bic = model_fit.aic + model_fit.bic
243
244                         if aic < best_aic:
245                             best_aic = aic
246                             best_order = (p, d, q)
```

```
247
248         if bic < best_bic:
249             best_bic = bic
250             best_order_bic = (p, d, q)
251
252         if total_aic_bic < best_total_aic_bic:
253             best_total_aic_bic = total_aic_bic
254             best_order_aic_bic = (p, d, q)
255
256     print(f"Order: ({p},{d},{q}) - AIC: {aic}, BIC: {bic},
257           AIC+BIC: {total_aic_bic}")
258     except:
259         continue
260
261     print(f"Best AIC Order: {best_order}, Best AIC: {best_aic}")
262     print(f"Best BIC Order: {best_order_bic}, Best BIC: {best_bic}")
263     print(f"Best AIC+BIC Order: {best_order_aic_bic}, Best AIC+BIC:
264           {best_order_aic_bic}")
265     # 我们默选Best AIC+BIC
266     return best_order_aic_bic
267
268     def fit(self):
269         """
270             拟合ARIMA模型
271
272         Parameters:
273             data (pd.Series): 输入的时间序列数据
274
275             self.model = ARIMA(self.data, order=self.order)
276             self.model = self.model.fit()
277             # 使用ARIMA模型预测时间序列的未来值
278             self.fitted_values = self.model.fittedvalues
279             residuals = self.data - self.fitted_values
280             # 因为对于处理GMV等数据较大的值,所以需要去掉residuals中的第一个值
281             self.std_residuals = np.std(residuals[1:])
282
283             # 判断一个值是否是异常值
284             def is_outlier(self, test_data):
285                 # 使用训练好的ARIMA模型来预测后下一个值
286                 new_data = self.model.forecast(steps=1)
287                 if abs(test_data - new_data[0]) <= self.threshold * self.std_residuals:
288                     return False
289                 return True
290
291             # 计算上下限
292             def get_up_down_level(self):
293                 """
```

```
292     找出数据中的所有上下区间。  
293     """  
294     # 使用训练好的ARIMA模型来预测后下一个值  
295     new_data = self.model.forecast(steps=1)  
296     lower_bound = new_data[0] - self.threshold * self.std_residuals  
297     if lower_bound < 0.0:  
298         lower_bound = 0.0  
299     upper_bound = new_data[0] + self.threshold * self.std_residuals  
300     return round(lower_bound, 2), round(upper_bound, 2)  
301  
302
```

ARIMA异常检测介绍

ARIMA (AutoRegressive Integrated Moving Average, 自回归积分滑动平均模型) 是一种广泛用于时间序列分析和预测的方法。虽然ARIMA主要用于预测未来值，但它也可以用于检测时间序列中的异常值。通过构建ARIMA模型，我们可以预测时间序列的未来值，并将实际值与预测值进行比较，从而识别出异常值。

基本概念

ARIMA模型的基本概念

1. 自回归 (AR) 部分：

- 使用过去的值来预测未来的值。
- 模型形式: $y_t = c + \phi_1 y_{t-1} + \phi_2 y_{t-2} + \dots + \phi_p y_{t-p} + \epsilon_t$

2. 差分 (I) 部分：

- 对时间序列进行差分操作，使其变得平稳。
- 模型形式: $y'_t = y_t - y_{t-1}$

3. 移动平均 (MA) 部分：

- 使用过去的误差项来预测未来的值。
- 模型形式: $y_t = \mu + \epsilon_t + \theta_1 \epsilon_{t-1} + \theta_2 \epsilon_{t-2} + \dots + \theta_q \epsilon_{t-q}$

使用ARIMA进行异常检测的步骤

1. 数据准备：

- 收集并清洗时间序列数据，确保数据的完整性和一致性。

2. 构建ARIMA模型：

- 使用历史数据拟合ARIMA模型。
- 确定模型的参数 (p, d, q) ，其中：
 - p 是自回归部分的阶数。
 - d 是差分的次数。
 - q 是移动平均部分的阶数。

3. 预测未来值：

- 使用ARIMA模型预测时间序列的未来值。

4. 计算残差：

- 计算实际值与预测值之间的残差。

5. 计算残差的统计特性：

- 计算残差的均值和标准差。

6. 定义异常检测阈值：

- 通常使用3倍标准差作为阈值。

7. 检测异常值：

- 如果某个时间点的残差超出定义的阈值，则认为该时间点的数据是异常的。

检测效果展示

异常数据检测源数据-核心指标-在营门店

EWMA异常检测介绍

EWMA (Exponentially Weighted Moving Average, 指数加权移动平均) 是一种常用的时间序列分析方法，特别适用于检测数据中的异常值。EWMA通过赋予较新的数据更高的权重，而较旧的数据权重逐渐减小，使得模型对数据的最新变化更加敏感。以下是使用EWMA进行异常检测的具体步骤和方法。

基本概念

1. EWMA公式：

$$S_t = \alpha x_t + (1 - \alpha) S_{t-1}$$

其中：

- S_t 是时间 t 的EWMA值。
- x_t 是时间 t 的观测值。
- S_{t-1} 是时间 $t - 1$ 的EWMA值。
- α 是平滑因子 (smoothing factor) , 取值范围为 $0 < \alpha \leq 1$ 。

2. 初始值：

- 通常, 初始值 S_0 可以设置为第一个观测值 x_1 , 或者设置为一个合理的初始估计值。

异常检测步骤

1. 数据准备：

- 收集时间序列数据, 确保数据的时间顺序性和完整性。

2. 计算EWMA值：

- 使用上述公式计算每个时间点的EWMA值。

3. 计算残差：

- 计算每个时间点的实际值与EWMA值之间的残差。

4. 计算残差的统计特性：

- 计算残差的均值和标准差。

5. 定义异常检测阈值：

- 通常使用3倍标准差作为阈值。

6. 检测异常值：

- 如果某个时间点的残差超出定义的阈值, 则认为该时间点的数据是异常的。

检测效果展示

异常数据检测源数据-核心指标-在营门店