

一、Granger因果关系分析

Granger 因果关系分析是一种统计方法，用于确定一个时间序列是否可以预测另一个时间序列。

- 基本思想是，如果一个时间序列 X 中的过去值能够提供 Y 未来值的显著信息，那么可以说 X Granger 引起了 Y。
- 重要的是，这种因果关系并不意味着真实的因果关系，而是统计学意义上的预测能力。

```
131 class Aranger:
132     def granger_analysis(self, d1, d2):
133         return len(significant_lags) > 0
134
135     1 usage
136     def analysis_all(self):
137         """对所有列，两两之间分析"""
138         self.set_origin_data()
139         col_names = self.df.columns.tolist()
140         # 初始化相关性矩阵
141         matrix = pd.DataFrame( data=0, index=col_names, columns=col_names)
142         # 使用嵌套循环遍历所有可能的变量组合
143         for i in range(len(col_names)):
144             for j in range(i + 1, len(col_names)):
145                 var1 = col_names[i]
146                 var2 = col_names[j]
147                 print(f"分析 {var1} 与 {var2} 之间的 Granger 因果关系:")
148                 # 方向 var1 -> var2
149                 has_causality_1 = self.granger_analysis(var1, var2)
150                 # 方向 var2 -> var1
151                 has_causality_2 = self.granger_analysis(var2, var1)
152                 if has_causality_1:
153                     matrix.loc[var1, var2] = 1
154                     # matrix.loc[var2, var1] = 1
155                 if has_causality_2:
156                     matrix.loc[var2, var1] = 1
157                 print('-' * 50) # 分隔不同变量组合的输出
158         print("Granger 因果关系矩阵: (行是因，列是果，对称位置都为1，则互为因果)")
159         print(matrix)
160
161 if __name__ == '__main__':
162     # 读取数据，第二行开始
163     df = pd.read_csv( filepath_or_buffer: '000001.txt', encoding='ansi', skip_rows=1, sep=',', engine='python')
164     analysis_obj = Aranger(df)
165     analysis_obj.analysis_all()
```

运行analysis_1

选择的最小滞后阶数: 7:
成交额的ADF检验结果:
ADF统计量: -3.3825
p值: 0.0116
临界值: {'1%': '-3.4312', '5%': '-2.8619', '10%': '-2.5670'}
结论: 成交额在5%显著性水平下平稳

成交量的ADF检验结果:
ADF统计量: -3.0321
p值: 0.0320
临界值: {'1%': '-3.4312', '5%': '-2.8619', '10%': '-2.5670'}
结论: 成交量在5%显著性水平下平稳

选择的最小滞后阶数: 7:

Granger 因果关系矩阵: (行是因，列是果，对称位置都为1，则互为因果)

	开盘	最高	最低	收盘	成交量	成交额
开盘	0	1	1	1	1	1
最高	1	0	1	1	1	1
最低	1	1	0	1	1	1
收盘	1	1	1	0	1	1
成交量	1	1	1	1	0	1
成交额	1	1	1	1	1	0

进程已结束，退出代码为 0

二、方向性条件概率

方向性条件概率（Directional Conditional Probability）通常是指在给定方向上的某个条件下事件发生的概率。

它涉及概率分布的方向性或时间顺序，这在时间序列分析、因果关系分析或贝叶斯网络中可能会遇到。

- 在时间序列分析中，方向性条件概率考虑时间顺序。例如，给定过去的观测值，预测未来观测值的条件概率：

$$P(Y_{t+1}|Y_t, Y_{t-1}, \dots, Y_{t-k})$$

条件概率是基于之前的 k个观测值，预测时间序列在 t+1时刻的值。

- 在因果关系分析中，方向性条件概率用于表示因果关系的方向。例如，假设我们有两个变量 X 和 Y，我们可能会计算：

$$P(Y|X)$$

表示在 X发生的条件下 Y发生的概率。如果 P(Y|X)和 P(X|Y)显著不同，这可能暗示一个方向性的因果关系。

```
def directional_conditional_probability(series_x, series_y):
    # 计算series_x和series_y中大于0的逻辑数组
    positive_x = series_x > 0
    positive_y = series_y > 0

    # 计算在series_x上升时series_y也上升的概率
    p_y_given_x = sum(positive_x & positive_y) / sum(positive_x)
    # 计算在series_x下降时series_y上升的概率
    p_y_given_not_x = sum(~positive_x & positive_y) / sum(~positive_x)

    return p_y_given_x, p_y_given_not_x

# 准备需要分析的列，这里包括了开盘、最高、最低、收盘、成交量和成交额的差分
columns = ['开盘_diff', '最高_diff', '最低_diff', '收盘_diff', '成交量_diff', '成交额_diff']
# 生成所有可能的列组合，用于分析每一列与其他列之间的方向性条件概率
combinations = [(x, y) for idx, x in enumerate(columns) for y in columns[idx + 1:]]

# 对所有组合进行方向性条件概率分析，并打印结果
for x, y in combinations:
    # 调用directional_conditional_probability函数计算条件概率
    p_y_given_x, p_y_given_not_x = directional_conditional_probability(data[x], data[y])

    # 从列名中提取主要的股票指标名称，例如从'开盘_diff'提取'开盘'
    x_label = x.split('_')[0]
    y_label = y.split('_')[0]

    # 打印结果，使用GREEN颜色代码高亮显示上升的概率
    print(f'{GREEN}当 {x_label} 上升时, {y_label} 也上升的概率为: {p_y_given_x:.4f}:{RESET}')
    # 打印结果，显示下降时另一个序列上升的概率
    print(f'当 {x_label} 下降时, {y_label} 反而上升的概率为: {p_y_given_not_x:.4f}')
    print() # 打印空行，增加输出的可读性
```

运行 analysis_2

↑

↓

≡

⌵

⌶

🗑

当 最高 上升时, 成交量 也上升的概率为: 0.6287:

当 最高 下降时, 成交量 反而上升的概率为: 0.3106

当 最高 上升时, 成交额 也上升的概率为: 0.6486:

当 最高 下降时, 成交额 反而上升的概率为: 0.2896

当 最低 上升时, 收盘 也上升的概率为: 0.7222:

当 最低 下降时, 收盘 反而上升的概率为: 0.3217

当 最低 上升时, 成交量 也上升的概率为: 0.5313:

当 最低 下降时, 成交量 反而上升的概率为: 0.4268

当 收盘 上升时, 成交量 也上升的概率为: 0.5312:

当 收盘 下降时, 成交量 反而上升的概率为: 0.4222

当 收盘 上升时, 成交额 也上升的概率为: 0.5452:

当 收盘 下降时, 成交额 反而上升的概率为: 0.4082

当 成交量 上升时, 成交额 也上升的概率为: 0.9048:

当 成交量 下降时, 成交额 反而上升的概率为: 0.0891

进程已结束，退出代码为 0

三、交互延迟互信息

交互延迟互信息（Interaction Delay Mutual Information，**IDMI**）是一种用于衡量两个时间序列之间的相互作用和延迟关系的统计方法。

它结合了**互信息**和**时间延迟分析**，用于识别和量化在时间序列数据中存在的非线性依赖关系。

- 互信息（Mutual Information）是用来衡量两个随机变量之间的非线性依赖关系的量度。它表示两个变量之间共享的信息量。互信息越大，说明两个变量之间的关系越强。互信息可以计算为：

$$I(X;Y) = \sum_{x \in X} \sum_{y \in Y} p(x,y) \log \frac{p(x,y)}{p(x)p(y)} = P(X) + P(Y) - P(X,Y)$$

其中：

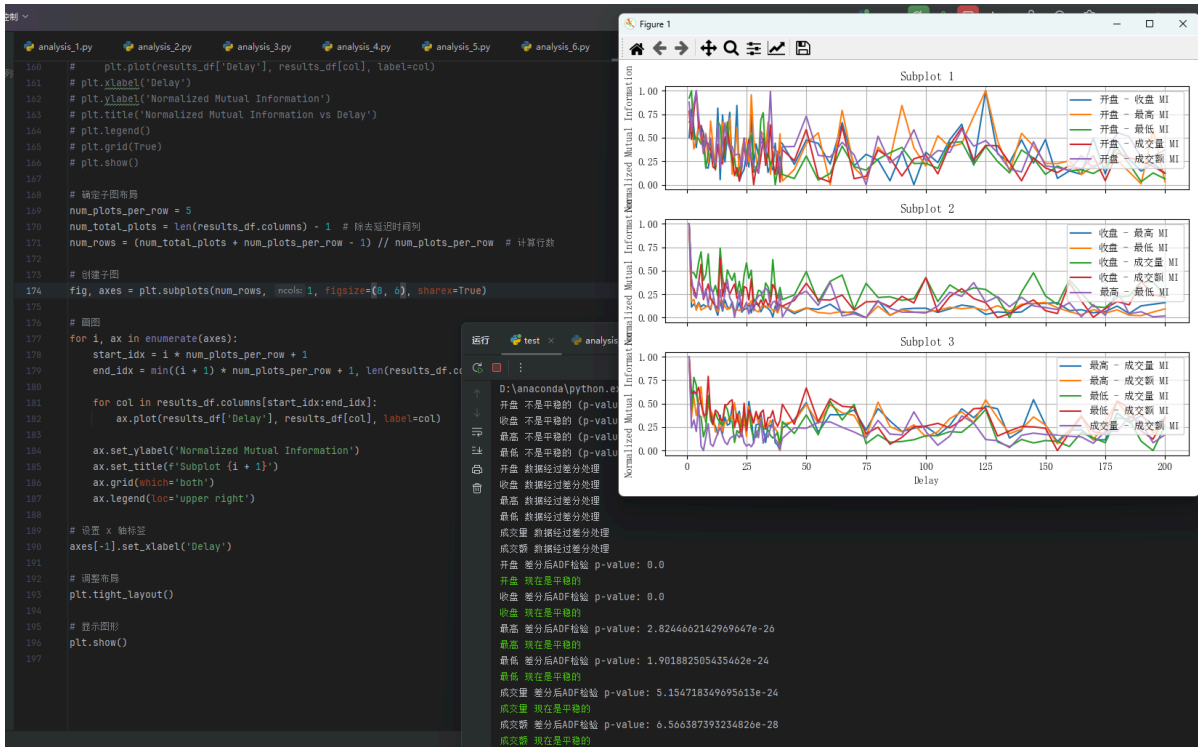
- p(x, y) 是 X 和 Y 的**联合概率分布**。
- p(x) 和 p(y) 是 X 和 Y 的**边际概率分布**。

- 交互延迟互信息通过将 **时间序列的延迟（滞后）** 与互信息结合，用于衡量一个时间序列的过去状态对另一个时间序列当前状态的影响。其步骤通常包括：

计算互信息：计算两个时间序列在不同延迟下的互信息，以量化它们的相互依赖。

引入延迟：考虑时间序列中的滞后或延迟。例如，对于时间序列 X 和 Y，计算X(t-k) 对 Y(t) 的互信息，其中 k 是滞后期。

计算交互延迟互信息：在不同的**延迟期**上计算互信息，得到关于时间序列相互作用的延迟信息。这可以揭示出时间序列之间的复杂非线性关系和交互模式。



四、时间滞后相干性分析

时间滞后相干性分析（Time-Lagged Coherence Analysis）是一种用于研究时间序列之间在不同时间滞后下的线性关系和相干性的技术。

它特别适用于探究两个时间序列之间在不同时间延迟上的相互依赖关系。

1. 相干性（Coherence）：

相干性是衡量两个时间序列之间在频域上的线性相关性。相干性值范围从0到1，0表示无相干性，1表示完全的相干性。高相干性表明两个时间序列在相同频率下有强的线性关系。

2. 时间滞后（Time Lag）：

时间滞后是指在一个时间序列的观测值与另一个时间序列的观测值之间的时间差。在分析中，我们可能会考虑不同的滞后期，以了解两个时间序列在这些滞后期下的关系如何变化。

分析过程

1. 计算相干性：

相干性分析通常在频域中进行。首先将时间序列转换到频域，这可以通过傅里叶变换实现。然后，计算两个时间序列的功率谱密度（Power Spectral Density, PSD）以及它们的互谱（Cross-Spectrum）。

相干性可以通过以下公式计算：

$$C_{XY}(f) = \frac{|S_{XY}(f)|^2}{S_{XX}(f) \cdot S_{YY}(f)}$$

其中：

$S_{XY}(f)$ 是 (X) 和 (Y) 的互谱。

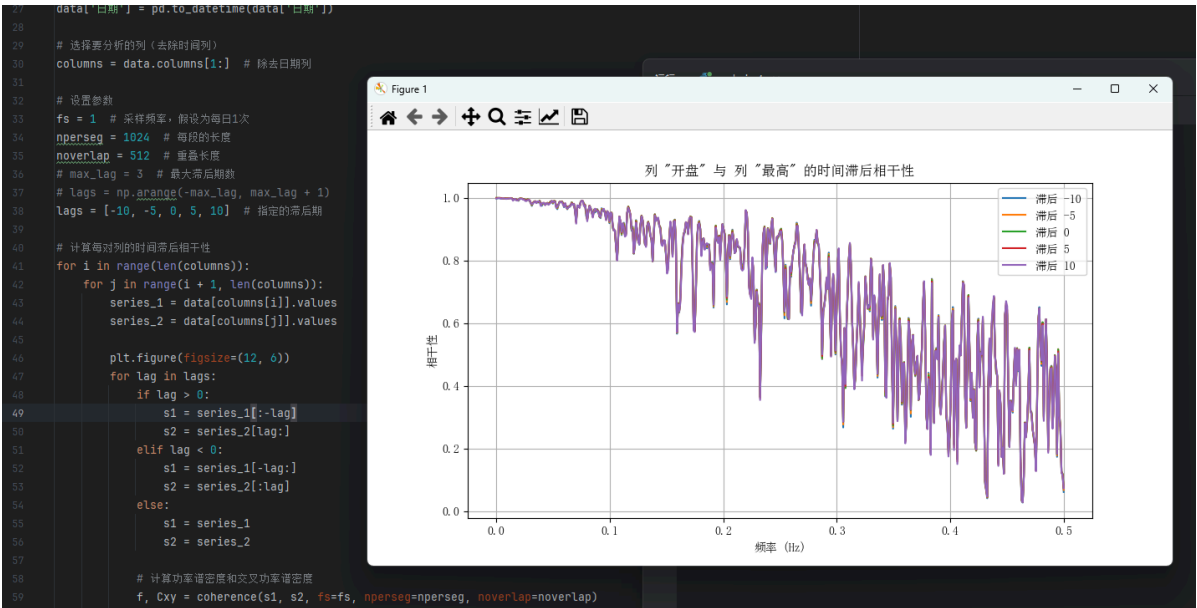
$S_{XX}(f)$ 和 $S_{YY}(f)$ 是 (X) 和 (Y) 的功率谱密度。

2. 考虑不同时间滞后：

在时间滞后相干性分析中，分析不同时间滞后下的相干性可以揭示时间序列之间的复杂依赖关系。这通常涉及到将时间序列在不同滞后后期进行处理，并计算相应的相干性。

3. 分析结果：

分析结果通常以图形方式展示，其中相干性随时间滞后变化的图表可以揭示不同滞后期下的相干性模式。高相干性值表示两个时间序列在特定频率和滞后下有强的线性关系，而低相干性值则表示较弱的关系。



五、向量自回归条件异方差

向量自回归条件异方差模型（Vector Autoregressive Conditional Heteroskedasticity, VARCH）是一种扩展的模型，用于处理多变量时间序列数据中存在的条件异方差性。

它结合了向量自回归（VAR）模型和条件异方差（ARCH/GARCH）模型的特点，适用于分析时间序列中多个变量的相互影响及其波动性。

向量自回归（VAR）模型：

向量自回归（VAR）模型是一种用于分析多个时间序列变量之间的线性动态关系的模型。

在 VAR 模型中，每个时间序列的当前值都是其过去值的线性函数，同时也可能受到其他时间序列过去值的影响。V

条件异方差（ARCH/GARCH）模型：

条件异方差（ARCH）模型用于处理时间序列数据中的波动性聚集现象。于1986年扩展为 GARCH（Generalized ARCH）模型。

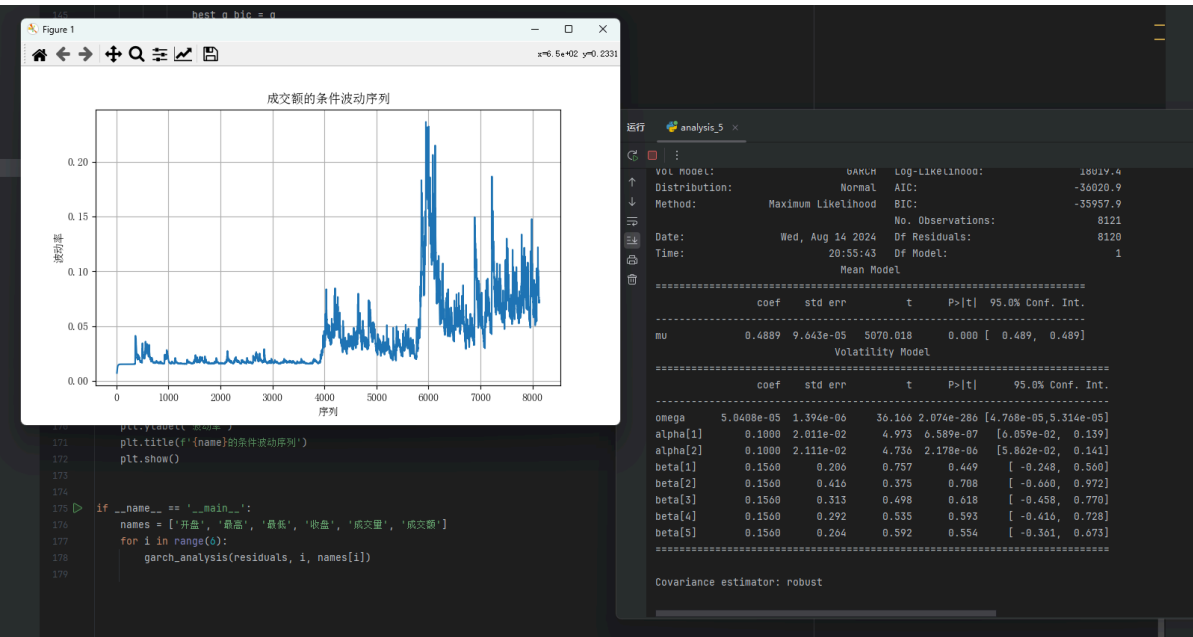
- **ARCH 模型：**假设时间序列的条件方差依赖于过去误差的平方。
- **GARCH 模型：**在 ARCH 模型的基础上，引入了过去方差的滞后项，使得条件方差依赖于过去误差和过去方差。

向量自回归条件异方差（VARCH）模型：

VARCH 模型结合了 VAR 模型和 ARCH/GARCH 模型，处理多变量时间序列中的波动性和依赖关系。VARCH 模型的主要特点包括：

- 1. **向量自回归**：模型的基础是 VAR 模型，用于描述时间序列变量之间的线性动态关系。
- 2. **条件异方差**：在 VAR 模型的基础上，引入条件异方差模型（如 GARCH），用于建模误差项的条件方差。

确定滞后阶数常用的方法包括 **信息准则**（AIC、BIC、HQIC、FPE）和**交叉相关函数**（CCF）。



六、递归量化分析

第1步：相空间重构

- 1. **确定嵌入维度**：
 - 方法：使用FNN（False Nearest Neighbors）方法来确定嵌入维度。
 - 步骤：
 1. 计算不同嵌入维度下的假最近邻数量。
 2. 选择假最近邻数量最小的维度作为最佳嵌入维度。
- 2. **确定延迟时间**：
 - 方法：使用最小二乘法（或自相关函数）来确定延迟时间。
 - 步骤：
 1. 计算自相关函数或互相关函数。
 2. 找到自相关函数首次显著下降的位置，作为延迟时间。

第2步：距离矩阵确定

- 1. **计算距离矩阵**：
 - 互信息距离：
 1. 计算每对点的互信息量。

- 2. 使用互信息量作为距离度量。
 - **曼哈顿距离：**
 - 1. 对于每对点，计算曼哈顿距离。
 - 2. 曼哈顿距离 = $\sum |x_i - y_i|$ 。
 - **余弦相似度：**
 - 1. 计算每对点的余弦相似度。
 - 2. 余弦相似度 = $(x \cdot y) / (||x|| * ||y||)$ 。
 - **同时运用：**
 - 1. 将上述三种距离计算结果结合，形成一个综合距离矩阵。
2. **阈值确定：**
- **K-最近邻阈值法：**
 - 1. 计算不同K值的K-最近邻距离。
 - 2. 选择合适的K值，使得距离矩阵能有效反映时间序列中的动态结构。
 - **确定K值的方法：**
 - **MSE (均方误差)：**
 - 1. 对不同K值，计算训练集和测试集上的均方误差。
 - 2. 选择MSE最小的K值。
 - **邻近误差：**
 - 1. 计算每个K值下的邻近误差。
 - 2. 选择邻近误差最小的K值。
 - **准确率-距离平衡：**
 - 1. 计算不同K值下的准确率和距离平衡。
 - 2. 选择准确率和距离平衡最优的K值。
 - **AIC和BIC (赤池信息准则和贝叶斯信息准则)：**
 - 1. 计算不同K值下的AIC和BIC值。
 - 2. 选择AIC和BIC值最小的K值。

后续步骤：RQA分析

- 1. **计算递归图：**
 - 使用前述的距离矩阵，生成递归图 (recurrence plot) 。
- 2. **量化指标：**
 - 计算以下RQA量化指标：
 - **Recurrence Rate (RR)：**递归图中点的比例。
 - **Determinism (DET)：**测量递归图中的确定性。
 - **Average Diagonal Line Length (L)：**测量递归图中的平均对角线长度。
 - **Entropy (ENT)：**递归图的熵值。
- 3. **分析结果：**
 - 通过量化指标分析时间序列的动态相关性。
 - 比较不同时间序列的RQA指标，以检测它们之间的相关性。

```

324 def start_analysis(filename, select_col):
325     df = df[-365:]
326
327     result = adfuller(df[select_col].values)
328     print(f'ADF 统计量: {result[0]:.4f}')
329     print(f'p 值: {result[1]:.4f}')
330     if result[1] > 0.05:
331         print(f'{RED}序列不平稳, 将进行差分。{RESET}')
332         df[select_col] = df[select_col].diff().dropna()
333     else:
334         print(f'{GREEN}序列平稳{RESET}')
335
336     df = df.dropna(subset=[select_col])
337
338     # 数据标准化处理
339     scaler = MinMaxScaler()
340     # 对df进行拟合和转换
341     df_scaled = scaler.fit_transform(df)
342     df = pd.DataFrame(df_scaled, columns=df.columns)
343
344     # 最佳延迟时间
345     delay_time = delay_time_selection_mi(df[select_col].values)
346     # 获取最佳嵌入维度
347     embedding_dim = fnn_method(df[select_col].values, max_delay=delay_time)
348     # 执行RQA分析
349     rqa_analysis(df[select_col].values, embedding_dim, delay_time)
350
351     if __name__ == '__main__':
352         file = '000001.txt'
353         col = '成交量'
354
355         print(f'{GREEN}=====分析{col}===== {RESET} ')
356         start_analysis(file, col)
357         print(f'{GREEN}===== {RESET} ')
358
359

```

```

0:\anaconda\python.exe E:\Analysis\序列相关性检测\analysis_6.py
=====分析成交量=====
ADF 统计量: -3.9457
p 值: 0.0017
序列平稳
最佳延迟时间为: 4
最佳嵌入维度为: 4
最佳K值: 10
阈值为: 0.2469
递归率RR: 0.2476
平均对角线长度: 3.6194
最大对角线长度: 350.0000
平均垂直线长度: 4.9929
最大垂直线长度: 332.0000
递归点百分比: 24.76%

```

