

Pairs Trading - Cointegration

1. Importing the S&P 500 Stock Data

a. Loading the Necessary Libraries

```
In [1]: import os
import itertools
import numpy as np
import pandas as pd
import statmodels.api as sm
import matplotlib.pyplot as plt
from statmodels.tsa.stattools import coint
import matplotlib.pyplot as plt
from scipy.cluster.hierarchy import linkage, dendrogram, fcluster, incoincident

# Changing Working Directory
os.chdir("C:/Users/Jesul/OneDrive/Desktop/Quantitative Research/Pairs trading - Cointegration")

# Printing my working directory
print("Current Working Directory", os.getcwd())

Current Working Directory: C:/Users/Jesul/OneDrive/Desktop/Quantitative Research/Pairs trading - Cointegration

c. Loading the Data
```

```
In [3]: data = pd.read_csv("SP500_stock_data.csv")

# Previewing the data
print(data.head())

   Date      HMM      AOS      ABT      ABBV      ACN  \
0 19-01-2022  128.609131  74.333565  117.356979  118.657448  324.742422
1 20-01-2022  125.071003  73.081429  116.109398  116.146886  323.752289
2 21-01-2022  124.408150  71.782555  117.152153  115.822830  317.082398
3 22-01-2022  124.516235  73.614052  115.280769  115.980797  320.286144
4 23-01-2022  125.200806  70.455696  114.768692  116.796951  316.440033

   ADBE      AMD      AES      AFL  ...      WMB      WTM  \
0 516.580017  128.270004  19.911016  57.549511  ...  24.374533  215.411179
1 510.850068  121.809999  20.187204  56.769010  ...  24.156227  217.008545
2 499.910004  118.809998  19.807451  55.807239  ...  24.147833  214.334453
3 519.639973  116.529999  19.496744  56.375786  ...  23.736416  219.480560
4 502.720001  111.132997  19.132454  56.769010  ...  24.324316  215.496735

   WDAY      NYM      XEL      YUM      ZBRA  \
0 249.460007  86.744938  60.221244  103.644225  115.870186  313.109985
1 251.770004  84.835892  60.274922  101.743545  115.051178  301.420013
2 249.690002  81.422576  60.398664  102.287987  115.668197  484.179993
3 243.430005  82.595773  59.611261  102.765572  117.582626  490.559998
4 236.110001  84.212524  60.009083  99.212387  113.559996  474.000000

   ZBH      ZTS  \
0 115.336148  193.803757
1 115.720016  195.137146
2 115.249779  193.562210
3 115.293739  193.311020
4 113.754364  188.653809

[5 rows x 14 columns]
```

2. Data Pre-Processing

a. Checking for missing values

```
In [4]: # Count of missing values in each column
null_counts = data.isnull().sum()

# Display columns with the number of missing values
null_counts = null_counts[null_counts > 0].sort_values(ascending=False)

print(null_counts)

SOLV    548
dtype: int64

The stock 'SOLV' has 548 missing data points. This stock has been dropped from the entire dataset as more than half its data points are missing.
```

```
In [5]: data = data.drop(columns=['SOLV'])

print(data.dtypes)

Date      object
HMM      float64
AOS      float64
ABT      float64
ABBV     float64
ACN      float64
ADBE     float64
AMD      float64
AES      float64
AFL      float64
WMB      float64
WTM      float64
ZBH      float64
ZTS      float64
Length: 497, dtype: object
```

b. Formatting 'date' appropriately

```
In [7]: data['Date'] = pd.to_datetime(data['Date'], format='%d-%m-%Y', errors='coerce')
data.set_index('Date', inplace=True) #setting 'date' as an index

print(data.head())

   HMM      AOS      ABT      ABBV      ACN  \
Date
2022-01-19  128.609131  74.333565  117.356979  118.657448  324.742422
2022-01-20  125.071003  73.081429  116.109398  116.146886  323.752289
2022-01-21  124.408150  71.782555  117.152153  115.822830  317.082398
2022-01-24  124.516235  73.614052  115.280769  115.980797  320.286144
2022-01-25  125.200806  70.455696  114.768692  116.796951  316.440033

   ADBE      AMD      AES      AFL      A  ...  \
Date
2022-01-19  516.580017  128.270004  19.911016  57.549511  136.69701  ...
2022-01-20  510.850068  121.809999  20.187204  56.769010  135.771973  ...
2022-01-21  499.910004  118.809998  19.807451  55.807239  133.843435  ...
2022-01-24  519.639973  116.529999  19.496744  56.375786  134.448120  ...
2022-01-25  502.720001  111.129997  19.132454  56.769010  130.992477  ...

   WMB      WTM      NYM      XEL  \
Date
2022-01-19  24.374533  215.411179  249.460007  86.744939  60.221244
2022-01-20  24.156227  217.008545  251.770004  84.835892  60.274922
2022-01-21  24.147833  214.334453  245.690002  81.422576  60.398664
2022-01-24  23.736416  219.480560  247.630005  82.595673  59.611267
2022-01-25  24.324316  215.496735  236.110001  84.212524  60.009083

   XYL      YUM      ZBRA      ZBH      ZTS  \
Date
2022-01-19  103.644325  115.870186  513.109985  115.936348  193.803757
2022-01-20  101.743545  115.051178  501.420013  115.720016  195.137146
2022-01-21  102.287987  115.668197  494.179993  115.249779  193.562210
2022-01-24  102.765572  117.582626  490.559998  116.293739  193.311020
2022-01-25  99.212387  113.559996  474.000000  113.754364  188.653809

[5 rows x 496 columns]
```

c. Converting raw prices to scale-free log-returns

```
In [8]: log_returns = log_returns / data.shift(1) * 100

# Drop the first row as it is null and rounding up to four decimal points
log_returns = log_returns.dropna().round(4)

# printing the dataset of log-returns
print(log_returns.head())

   HMM      AOS      ABT      ABBV      ACN      ADBE      AMD      AES  \
Date
2022-01-20 -2.7896 -1.6988 -1.0688 -1.6255 -3.3059 -1.1154 -5.1018 -1.3776
2022-01-21 -0.5311 -1.7933 0.8941 -0.7924 -2.0817 -2.1648 -2.5559 -1.8991
2022-01-24 0.0868 2.5194 -1.4103 0.1363 0.9991 3.8747 -1.9397 -1.1648
2022-01-25 0.5481 -0.3852 -0.4452 0.7032 -2.2018 -3.3181 -4.7448 -1.8767
2022-01-26 -2.5945 -0.9060 -2.6302 0.8529 -1.5222 -2.3807 -0.3787 1.1213

   AFL      A  ...      WMB      WTM      NYM      XEL  \
Date
2022-01-20 -1.3798 -0.6788 ... -0.8997 0.7388 0.9217 -2.2253 0.5880
2022-01-21 -1.6944 -1.4225 ... -0.0348 -1.1457 -2.4444 -1.6788 0.0987
2022-01-24 1.0136 0.4426 ... -1.7184 2.2784 0.7865 -0.9974 -1.1648
2022-01-25 0.6800 -2.6039 ... 2.4460 -1.8318 -4.7638 1.9386 0.6651
2022-01-26 0.4192 -0.7908 ... 0.8250 -0.4865 -4.8020 -0.3824 -0.3837

   XYL      YUM      ZBRA      ZBH      ZTS  \
Date
2022-01-20 -1.8510 -0.7093 -2.3046 -0.1868 -0.6057
2022-01-21 0.5337 0.5003 -1.4544 -0.4072 -0.8104
2022-01-24 0.4638 1.8761 -0.7352 0.9017 -0.1299
2022-01-25 -3.5188 -1.4546 -3.4340 -2.2078 -2.4387
2022-01-26 -1.2887 -1.5773 0.3580 -0.3810 -2.8099

[5 rows x 496 columns]
```

3. Cluster Analysis

a. Computing correlation based distance matrix

```
In [9]: # Correlation matrix of stocks
corr_matrix = log_returns.corr()

# Convert correlation to distance for clustering
distance_matrix = 1 - corr_matrix

print(distance_matrix.head())

   HMM      AOS      ABT      ABBV      ACN      ADBE      AMD      AES  \
Date
2022-01-20 0.000000 0.534578 0.123233 0.751459 0.618890 0.786337 0.693560
2022-01-21 0.534578 0.000000 0.703753 0.782629 0.570781 0.655509 0.653458
2022-01-24 0.123233 0.703753 0.000000 0.662106 0.613704 0.690275 0.799555
2022-01-25 0.751459 0.782629 0.662106 0.000000 0.814740 0.893400 0.964565
2022-01-26 0.618890 0.570781 0.613704 0.814740 0.000000 0.467642 0.561409

   AFL      A  ...      WMB      WTM      NYM      XEL  \
Date
2022-01-20 0.529711 0.435295 0.604930 0.778335 0.517357 0.607826 0.583906
2022-01-21 0.703753 0.000000 0.805135 0.737037 0.622389 0.687591 0.624678
2022-01-24 0.534578 0.703753 0.782629 0.570781 0.690363 0.687853 0.525939
2022-01-25 0.694936 0.640814 0.513738 0.705764 0.652337 0.727047 0.768720
2022-01-26 0.604545 0.628006 0.523717 0.705678 0.519975 0.629366 0.733607

   XYL      YUM      ZBRA      ZBH      ZTS  \
Date
2022-01-20 0.654230 0.622526 0.603139 0.654574 0.647170 0.765860
2022-01-21 0.653694 0.615478 0.530277 ... 0.771594 0.579120 0.683700
2022-01-24 0.742094 0.637121 0.574443 ... 0.769594 0.589639 0.719513
2022-01-25 0.818710 0.746838 0.702579 ... 0.798528 0.724190 0.890777
2022-01-26 0.707094 0.630794 0.508951 ... 0.764338 0.605173 0.616569

   XYL      YUM      ZBRA      ZBH      ZTS  \
Date
2022-01-20 0.592414 0.554085 0.493720 ... 0.693530 0.557928 0.610643
2022-01-21 0.653565 0.621194 0.602011 ... 0.781592 0.637327 0.699335
2022-01-24 0.620198 0.737160 0.497291 ... 0.690363 0.687853 0.586618
2022-01-25 0.698536 0.668946 0.596548 ... 0.798473 0.657734 0.706693
2022-01-26 0.660590 0.684576 0.474535 ... 0.779625 0.631171 0.615209

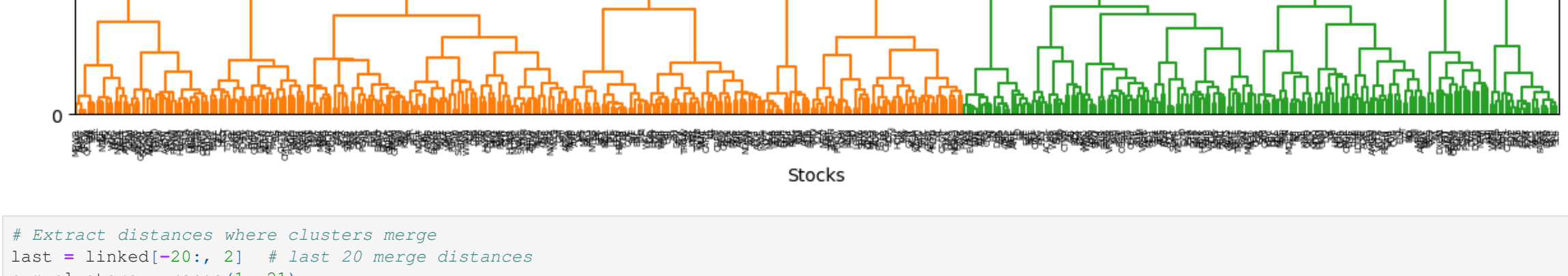
[496 rows x 496 columns]
```

b. Hierarchical cluster analysis

```
In [10]: # Perform hierarchical clustering using Ward's method
linked = linkage(distance_matrix, method='ward')

# Plot dendrogram
plt.figure(figsize=(15, 7))
dendrogram(linked, labels=log_returns.columns, last_rotation=90)
plt.title('Hierarchical Clustering of S&P 500 Stocks')
plt.xlabel('Stocks')
plt.ylabel('Distance')
plt.show()
```

C:/Users/Jesul/VggData/Local/Temp/ipykernel_20648/212663342.py:2: ClusterWarning: scipy.cluster: The symmetric non-negative hollow observation matrix looks suspiciously like an uncondensed distance matrix

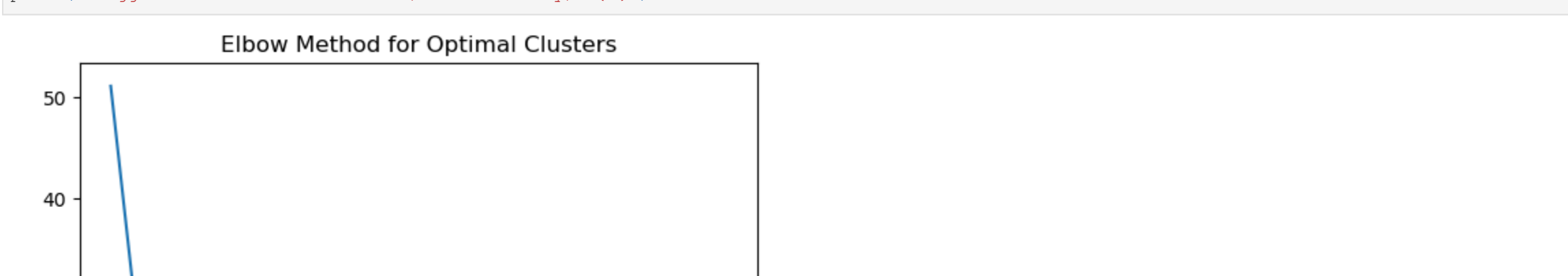


```
In [11]: # Extract distances where clusters merge
last = linked[-20, 2] # Last 20 merge distances
num_clusters = range(1, 21)

# Compute distances between successive merges
acceleration = np.diff(last, 2) # 2nd derivative
k = acceleration.argmax() + 2

plt.plot(num_clusters, last[1:-1])
plt.title('Elbow Method for Optimal Clusters')
plt.xlabel('Number of clusters')
plt.ylabel('Distance')
plt.show()

print(f"Suggested number of clusters (mathematically): {k}")
```



Suggested number of clusters (mathematically): 19

```
In [12]: num_clusters = 6
cluster_labels = fcluster(linked, num_clusters, criterion='maxclust')

# Map stocks to clusters
clusters = pd.DataFrame({'Stock': log_returns.columns, 'Cluster': cluster_labels})
print(clusters.head())

   Stock Cluster
0      HMM      1
1      AMD      1
2      XOM      1
3      KRO      1
4      BAC      1
...      ...
440     TMO      6
441     COP      6
442     CVX      6
443     CRM      6
444     EA      6

[496 rows x 2 columns]
```

3. Cointegration Test

a. Cluster 1

```
In [13]: # Step 1: Select stocks in Cluster 1
clusters['clusters'] = clusters['clusters'].tolist()

cluster_data = log_returns[clusters['clusters'] == 1]

# Step 2: Generate all possible pairs
pairs = list(itertools.combinations(cluster_data['Stocks'], 2))

# Step 3: Test cointegration and store results
coint_results = []

for stock1, stock2 in pairs:
    score, pvalue, _ = coint(cluster_data[stock1], cluster_data[stock2])
    if pvalue < 0.01: # filter for significant pairs
        coint_results.append((stock1, stock2, coint_results[stock2, 'p_value']))

# Step 4: Convert to DataFrame
coint_df_1 = pd.DataFrame(coint_results)
coint_df_1 = coint_df_1.sort_values('p_value')

print(coint_df_1)

   Stock1 Stock2 p_value
0      HMM      ADBE 0.000000
8736     FCX     UPS 0.000000
8737     FCX     GMW 0.000000
8738     FCX     DTS 0.000000
8739     FCX     WBD 0.000000
...      ...
5919     CSX     SBUX 0.000025
5920     CSX     HAZ 0.000025
5921     CSX     SAK 0.000025
13468    CSX     QNSC 0.000040

[18695 rows x 3 columns]
```

```
In [14]: top_50_pairs_1 = coint_df_1.sort_values('p_value').head(50)

print(top_50_pairs_1)

   Stock1 Stock2 p_value
0      HMM      ADBE 0.0
3594     BA      HST 0.0
1584     BA      CRM 0.0
3595     BA      BHM 0.0
3596     BA      HUBB 0.0
3597     BA      FOC 0.0
3598     BA      INTC 0.0
3599     BA      JNJ 0.0
3600     BA      JPM 0.0
3601     BA      JFF 0.0
3602     BA      IP 0.0
3603     BA      JCS 0.0
3604     BA      ZBT 0.0
3605     BA      JEL 0.0
3596     BA      HPQ 0.0
3077     AKOX     CSX 0.0
3078     PMAN     KENV 0.0
6088     DDOG     FINV 0.0
3593     BA      HLT 0.0
2758     APTV     HSH 0.0
2759     BA      HLT 0.0
2751     APTV     HSH 0.0
2752     BA      HLT 0.0
2753     APTV     HSH 0.0
2754     BA      HLT 0.0
2755     APTV     HSH 0.0
2756     BA      HLT 0.0
2757     APTV     HSH 0.0
2758     APTV     HSH 0.0
2759     APTV     HSH 0.0
2760     APTV     HSH 0.0

[496 rows x 3 columns]
```

4. Building a Trading Strategy

a. Plotting our cointegrated pair

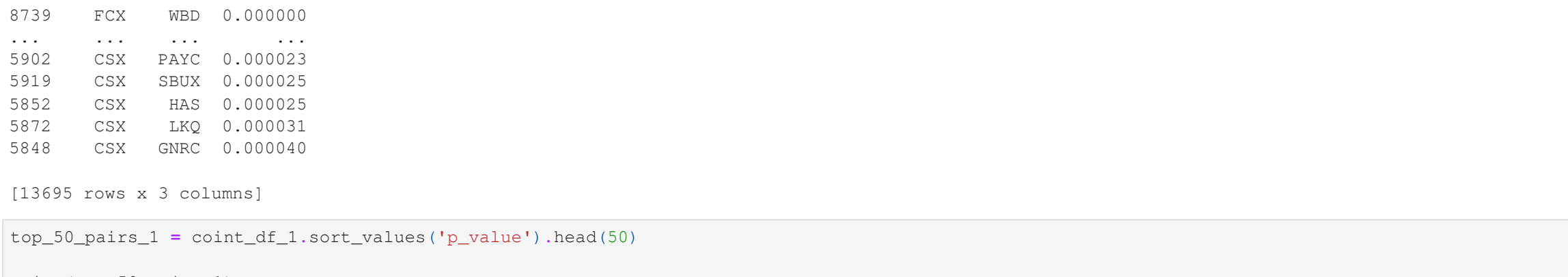
```
In [15]: pair = ('ABT', 'LNT')
pair_data = data[pair]

fig, ax1 = plt.subplots(figsize=(10,4))

# Left y-axis for ABT
color = 'tab:blue'
ax1.set_xlabel('Date')
ax1.set_ylabel('ABT Price', color=color)
ax1.plot(pair_data['ABT'], color=color, label='ABT')

# Right y-axis for LNT
ax2 = ax1.twinx()
color = 'red'
ax2.set_ylabel('LNT Price', color=color)
ax2.plot(pair_data['LNT'], color=color, label='LNT')

plt.title('Closing Prices of ABT and LNT')
fig.tight_layout()
plt.show()
```



```
In [16]: plt.figure(figsize=(10,4))
plt.scatter(pair_data['ABT'], pair_data['LNT'], alpha=0.6)
plt.xlabel('ABT Price')
plt.ylabel('LNT Price')
plt.grid(True)
plt.show()
```



```
In [17]: correlation = pair_data['ABT'].corr(pair_data['LNT'])
print("Correlation between ABT and LNT: ", correlation.round(2))

Correlation between ABT and LNT: 0.82
```

b. Computing the spread: OLS Regression

```
In [18]: pair_data = data[pair].copy()

X = pair_data['ABT']
y = pair_data['LNT']

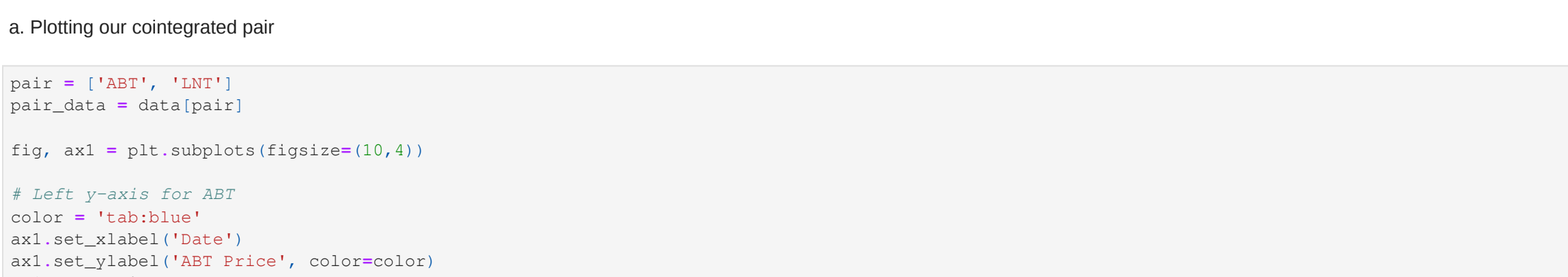
# Add constant for intercept
X_const = sm.add_constant(X)
model = sm.OLS(y, X_const).fit()
beta = model.params['ABT'] # slope
alpha = model.params['const'] # intercept

# Compute spread
pair_data['spread'] = y - (alpha + beta * X)
print(pair_data)

   ABT      LNT      Spread
Date
2022-01-19  117.356979  52.419163 -4.014524
2022-01-20  116.109398  52.692635 -3.240771
2022-01-21  117.152153  52.719101 -3.636664
2022-01-24  115.280769  51.88102 -3.75187
2022-01-25  114.768692  51.439930 -3.958549
...      ...
2025-10-03  133.994156  66.800003 3.712136
2025-10-06  131.149348  68.120003 3.703990
2025-10-07  132.431122  68.070000 5.607282
2025-10-08  133.675583  67.870003 4.909552
2025-10-09  132.719188  67.239998 4.661814

[935 rows x 4 columns]
```

```
In [19]: plt.figure(figsize=(10,4))
plt.plot(pair_data['spread'], label='Spread', color='purple')
plt.axhline(0, color='black', linestyle='--', label='Mean')
plt.title('Spread between ABT and LNT')
plt.xlabel('Date')
plt.ylabel('Spread')
plt.legend()
plt.show()
```



```
In [20]: # Step 1: Lag the spread
spread_lag = spread[1:] # X_t-1
spread_next = spread[1:] # X_t+1
delta_spread = spread_next - spread_lag # X_t - X_{t-1}

# Step 2: Fit AR(2) model: ΔX_t = a + b*X_t + ε_t
X = sm.add_constant(spread_lag)
model = sm.OLS(delta_spread, X).fit()

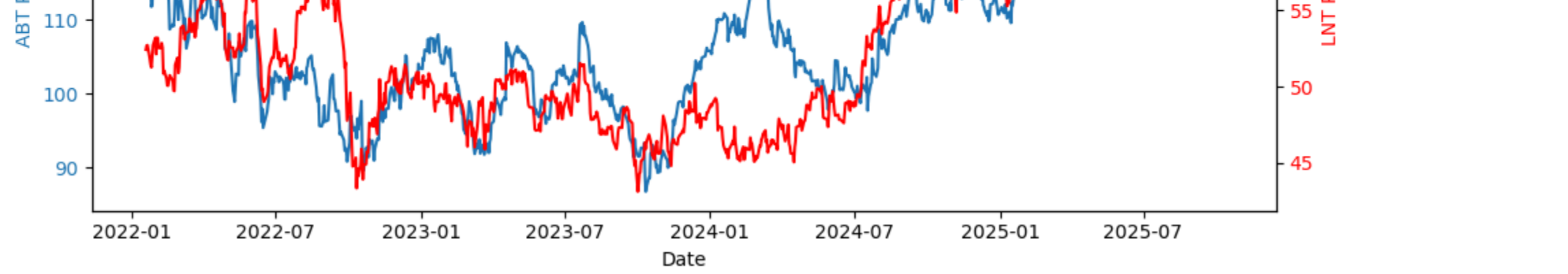
# Step 3: Extract parameters
mu = model.params[0] # intercept
sigma = model.params[1] # slope

# Step 4: Compute OU parameters
theta = -np.log(1 + b) # speed of mean reversion
mu = a / b # long-term mean
sigma = np.std(model.resid) # volatility

# Step 5: Compute half-life
half_life = np.log(2) / theta

# Print the OU model equation
print(f"OU Model Equation: dX_t = {theta} * (X_t - {mu}) dt + {sigma} * dB_t")
half_life (days): 31.57
```

```
In [21]: # Step 6: Plot the spread and OU mean
plt.figure(figsize=(10,4))
plt.plot(pair_data['spread'], label='Spread', color='blue')
plt.axhline(mu, color='red', linestyle='--', label='Long-term mean (mu={mu:.2f})')
plt.title('Spread with OU Model Mean')
plt.xlabel('Date')
plt.ylabel('Spread')
plt.legend()
plt.show()
```

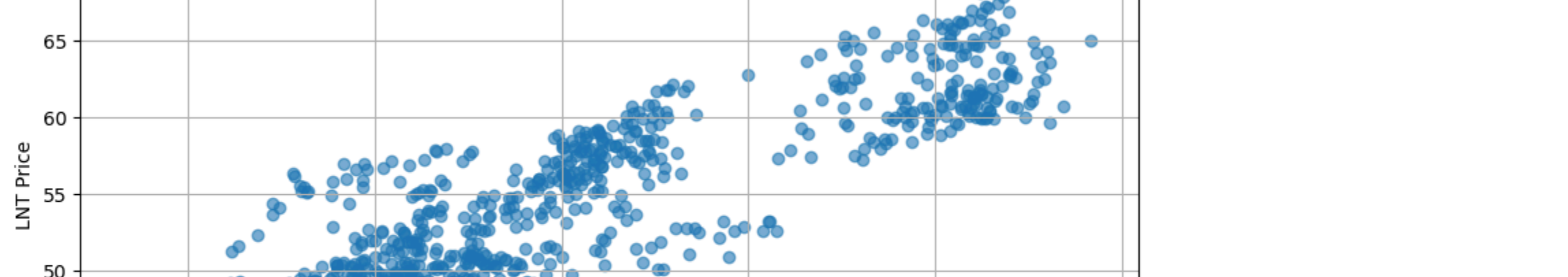


```
In [22]: theta = 0.0243 # speed of mean reversion
mu = -0.4288 # long-term mean
days = 900 # simulation length

# Step 1: Initialize the process
dt = 1 # daily steps
X[0] = mu # start at mean

# Step 2: Simulate OU process
for t in range(1, days):
    X[t] = X[t-1] + theta*(mu - X[t-1])*dt + sigma*np.sqrt(dt)*np.random.normal()

# Step 3: Plot the simulated OU process
plt.figure(figsize=(10,4))
plt.plot(X, label='Simulated Spread', color='blue')
plt.axhline(mu, color='red', linestyle='--', label='Long-term mean (mu={mu:.2f})')
plt.title('OU Process Simulation (900 Days)')
plt.xlabel('Days')
plt.ylabel('Spread')
plt.legend()
plt.show()
```

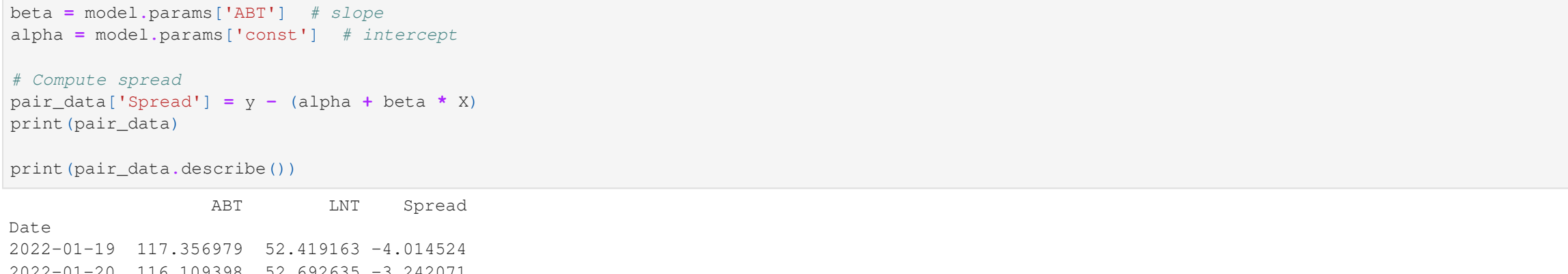


```
In [23]: pair_data['z_score'] = (pair_data['spread'] - pair_data['spread'].mean()) / pair_data['spread'].std()
print(pair_data)

   ABT      LNT      Spread      z_score
Date
2022-01-19  117.356979  52.419163 -4.014524 -1.832934
2022-01-20  116.109398  52.692635 -3.240771 -0.95611
2022-01-21  117.152153  52.719101 -3.636664 -1.070739
2022-01-24  115.280769  51.88102 -3.75187 -1.21149
2022-01-25  114.768692  51.439930 -3.958549 -1.16775
...      ...
2025-10-03  133.994156  66.800003 3.712136 1.094164
2025-10-06  131.149348  68.120003 3.703990 1.582999
2025-10-07  132.431122  68.070000 5.607282 1.625704
2025-10-08  133.675583  67.870003 4.909552 1.467106
2025-10-09  132.719188  67.239998 4.661814 1.374085

[935 rows x 4 columns]
```

```
In [24]: plt.figure(figsize=(10,4))
plt.plot(pair_data['z_score'], label='z-score', color='purple')
plt.axhline(0, color='black', linestyle='--', label='Mean')
plt.axhline(2, color='red', linestyle='--', label='Upper Threshold')
plt.axhline(-2, color='red', linestyle='--', label='Lower Threshold')
plt.title('z-score of ABT and LNT')
plt.xlabel('Date')
plt.ylabel('z-score')
plt.legend()
plt.show()
```



```
In [25]: # Save the dataframe to CSV
file_path = 'pair_data.csv'

pair_data.to_csv(file_path, index=False)

print(f'pair_data has been saved as {file_path}')
```