import numpy as np import pandas as pd import json import os import time import glob

from tqdm import tqdm

len(train\_list)

time series = []

100%|

In [4]:

Out[5]: 9304

1.42it/s]

len(ts\_list)

Out[3]: 9304

from datetime import datetime from dtaidistance import dtw

with open ("train list 2.txt", "r") as fp:

Соберем датаесет из рядов по этой выборке:

for i in tqdm(range(len(train list))):

ticker name = train list[i][0]

with open('ts list.npy', 'wb') as fp: np.save(fp, time series)

return array(a, dtype, copy=False, order=order, subok=True)

ts list = np.load("ts list.npy", allow pickle=True)

определения их схожести - Dynamic Time Warping

(cs.ucr.edu/~eamonn/KAIS\_2004\_warping.pdf)

from IPython.display import Image

Image('dtw eucl.png')

Euclidean

**def** dtw1(s1, s2):  $DTW = \{ \}$ 

for i in range(len(s1)):

for i in range(len(s2)):

for i in range(len(s1)):

DTW[(-1, -1)] = 0

def dtw2(s1, s2, window):

dtw matrix[0, 0] = 0.

подсчет матрицы:

ts dist matrix

Out[10]: array([[0.

In [9]:

for i in range(1, n+1):

dist = squareform(dist matrix)

4.96603827],

8.484980691,

7.33282846],

7.27461574],

Найдем linkage matrix:

#pdist(ts dist matrix)

# построим дендрограмму:

print(linkage matrix.shape)

# Обрежем дерево на значении 500

2,

66, 67, 79, 80,

196,

cluster\_labels.shape

new\_list\_arr = np.array(train list)

Out[26]: (9304,)

new\_arr

#

Out[6]: array([['SPI', '2020-09-23', '60'],

#for ticker in unq[0].tolist():

for i in tqdm(range(len(new arr))): ticker name = new arr[i][0]

df = pd.read csv(file)

print(new\_arr[i])

plt.close() except IndexError:

pass

кластеров. Останется 2729 наблюдений.

print("error")

try:

os.makedirs(path)

78, 21,

3,

from scipy.cluster.hierarchy import fcluster

# hand-select an appropriate cut-off on the dendrogram

5,

14, 15, 16, 17, 18, 19, 20, 21, 22,

70,

83,

96,

array([108, 23, 78, 111, 119, 25, 62, 61, 59, 46,

unq = np.unique(cluster\_labels, return\_counts=True)

4,

82,

95,

27, 28, 29, 30, 31, 32,

53, 54, 55, 56, 57, 58,

68, 69,

81,

94,

56, 183, 121, 75, 93,

64, 121,

42, 24,

35, 30, 14, 14, 35, 50,

67, 67, 29, 41, 124, 19,

plt.show()

linkage\_matrix

(9303, 4)

Out[23]: (array([ 1,

In [14]:

from warnings import simplefilter

simplefilter("ignore", ClusterWarning)

linkage\_matrix = ward(ts\_dist\_matrix)

fig = plt.figure(figsize=(20, 100))

[2.16297457, 0. 2.30455842],

#np.save('ts dist matrix.npy', dist)

ts\_dist\_matrix = np.load("ts\_dist\_matrix.npy")

[6.21311378, 7.28620543, 0.

DTW[(i, -1)] = float('inf')

DTW[(-1, i)] = float('inf')

for j in range(len(s2)):

n, m = s1.shape[0], s2.shape[0]

return np.sqrt(dtw\_matrix[-1, -1])

dtw.try import c() # проверим есть ли ОМР для Си

# используем squareform для перехода на обычную матрицу from scipy.spatial.distance import pdist, squareform

[3.20641225, 6.53808198, 2.307163 , ..., 0.

return linkage(y, method='ward', metric='euclidean')

#np.save('linkage\_matrix.npy', linkage\_matrix)

linkage\_matrix = np.load("linkage\_matrix.npy")

Out[14]: array([[3.78000000e+02, 5.90200000e+03, 1.48850550e+01, 2.00000000e+00],

[8.560000000e+02, 1.29800000e+03, 1.83981784e+01, 2.00000000e+00],[1.611000000e+03, 7.658000000e+03, 1.87199933e+01, 2.000000000e+00],

[1.86020000e+04, 1.86030000e+04, 8.90605951e+03, 5.71000000e+03], [1.86000000e+04, 1.86010000e+04, 1.33028215e+04, 3.59400000e+03], [1.86040000e+04, 1.86050000e+04, 2.78567638e+04, 9.30400000e+03]])

По linkage matrix сделаем иерархическую кластеризацию:

cluster labels = fcluster(linkage matrix, 500, criterion='distance')

6,

71,

84,

97,

9,

79, 136,

34, 27,

cluster\_labels = cluster\_labels.reshape(len(cluster\_labels),1)

new\_arr = np.hstack((new\_list\_arr, cluster\_labels)) #np.save('new\_arr\_after\_cluster.npy', new\_arr2)

new\_arr = np.load("new\_arr\_after\_cluster.npy")

['GLSI', '2020-12-09', '94'], ['APM', '2020-09-29', '104'],

['WTRH', '2020-03-11', '129'], ['ACRE', '2020-04-09', '97'],

['MVIS', '2020-10-08', '18']], dtype='<U11')

path = os.path.join(parent dir, str(ticker))

# parent\_dir = 'C:/Users/Kuanysh/Downloads/pump\_and\_dump/clusters\_1'

parent dir = 'C:/Users/Kuanysh/Downloads/pump and dump/agg tickers 1m'

warn too much data=1000, title=str(new\_arr[i]), \

pictures dir = 'C:/Users/Kuanysh/Downloads/pump and dump/clusters 1' + '/' + new arr[i][2]

mpf.plot(df.iloc[:,1:6], figsize =(2048/100,1080/100), type='candle', volume=True, \

Далее рассмотрим все подпапки и отберем нужные кластеры из 166 кластеров вручную, чтобы по ним можно было найти центры

file2 = os.path.join(pictures\_dir, ticker\_name + '\_'+ str(new\_arr[i][1]) +'.png')

Создадим графики всех этих наблюдений в кластерах:

file = os.path.join(parent dir, ticker name + '.csv')

scale padding=0.2, savefig=file2)

df['date'] = pd.to\_datetime(df['time']).dt.date df['times'] = pd.to datetime(df['time']).dt.time

df['date'] = pd.to\_datetime(df['date']) df = df[df['date'] == str(new\_arr[i][1])] df=df.sort\_values(by='times', ascending=True)

df.index = pd.DatetimeIndex(df['time'])

7,

33,

40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52,

72,

85,

98,

26,

32,

53,

39,

82,

121, 59, 105, 27, 26, 35, 52, 49, 40, 50, 78, 59, 53, 25, 68, 57, 37, 70, 73, 111, 41, 56, 48, 77, 42, 85,

23, 66, 37, 37, 70, 73, 111, 41, 30, 40, 77, 42, 63, 67, 71, 69, 77, 54, 84, 85, 35, 64, 67, 25, 56, 22, 57, 14, 29, 31, 23, 28, 22, 19, 25, 17, 44, 61, 48, 37, 50, 50, 79, 45, 64, 77, 44, 44, 15, 79, 30, 33, 49, 70, 43, 72, 38, 82, 40, 46, 74, 20, 94, 31, 66, 100, 111, 14, 19, 34, 14, 10, 13, 31, 79], dtype=int64))

Создадим список из наблюдений и соответствующих им кластеров:

105, 106, 107, 108, 109, 110, 111, 112, 113, 114, 115, 116, 117, 118, 119, 120, 121, 122, 123, 124, 125, 126, 127, 128, 129, 130, 131, 132, 133, 134, 135, 136, 137, 138, 139, 140, 141, 142, 143, 144, 145, 146, 147, 148, 149, 150, 151, 152, 153, 154, 155, 156, 157, 158, 159, 160, 161, 162, 163, 164, 165, 166], dtype=int32),

59, 60,

8,

34,

73,

86,

51,

63,

35**,** 

13,

68,

9, 10,

61, 62,

35,

74,

87,

86,

57,

71,

25,

60,

23,

36,

75,

88,

73,

47,

47,

36,

69,

11,

24,

37,

63**,** 

76,

89,

76,

58,

63,

58,

58,

63,

99, 100, 101, 102, 103, 104,

12, 13,

80, 103,

68, 46,

75, 109,

38,

36,

33,

65**,** 

78,

64**,** 

77,

90,

[2.97714582, 7.49029629, 1.77260901, ..., 1.19643093, 0.

from scipy.cluster.hierarchy import ward, dendrogram, ClusterWarning

[4.96603827, 2.30455842, 8.48498069, ..., 7.33282846, 7.27461574,

dist= (s1[i]-s2[j])\*\*2

return sqrt(DTW[len(s1)-1, len(s2)-1])

w = np.max(np.array([window, np.absolute(n-m)]))

dtw matrix = np.full((n + 1, m + 1), np.inf)

cost = (s1[i-1] - s2[j-1])\*\*2

day = train list[i][1]

time series.append(ts)

train list = json.load(fp)

import datetime

Кластеризация малой выборки по DTW

import matplotlib.pyplot as plt

parent dir train = 'C:/Users/Kuanysh/Downloads/pump and dump/train2'

Нам нужна относительно маленькая выборка, чтобы по ней можно было найти центры кластеров и сравнивать их уже со всей

ts = pd.read csv(os.path.join(parent dir train, ticker name + ' '+ day +'.csv')).to numpy().ravel()

C:\ProgramData\Anaconda3\lib\site-packages\numpy\core\\_asarray.py:171: VisibleDeprecationWarning: Creating an n darray from ragged nested sequences (which is a list-or-tuple of lists-or-tuples-or ndarrays with different len gths or shapes) is deprecated. If you meant to do this, you must specify 'dtype=object' when creating the ndarr

Теперь посчитаем матрицу попарных евклидовых расстояний между рядами отного датасета для

Есть простой алгоритм, который считает все расстояния между рядами, сложности O(nm) где n и m длины рядов:

DTW[(i, j)] = dist + min(DTW[(i-1, j)], DTW[(i, j-1)], DTW[(i-1, j-1)])

Есть и более облегченная его версия с окном, которая позволяет значительно уменьшить вычислительную сложность:

for j in range (np.max(np.array([1, i-w])), np.min(np.array([m, i+w]))+1):

, 2.16297457, 6.21311378, ..., 3.20641225, 2.97714582,

, 7.28620543, ..., 6.53808198, 7.49029629,

*,* ..., 2.307163 *,* 1.77260901*,* 

C:\ProgramData\Anaconda3\lib\site-packages\scipy\cluster\hierarchy.py:834: ClusterWarning: scipy.cluster: The s

dn = dendrogram(linkage\_matrix, orientation='left', color\_threshold=4, truncate\_mode='level', p=500, labels=transfer | truncat

ymmetric non-negative hollow observation matrix looks suspiciously like an uncondensed distance matrix

, 1.19643093,

dist\_matrix = dtw.distance\_matrix\_fast(time\_series, window=60, compact=True)

DTW

dtw\_matrix[i, j] = cost + np.min(np.array([dtw\_matrix[i-1, j], dtw\_matrix[i, j-1], dtw\_matrix[i-1,

Применим библиотеку dtaidistance, которая будет использовать код на Си и паралеллить

| 9304/9304 [00:06<00:00, 136

большой выборкой. Отберем наблюдения с дневным диапазоном >50% - более яркие представители пампов и дампов.