Формирование обучающей выборки для predictive clustering import json import numpy as np import pandas as pd import os import time import glob import datetime import itertools from tqdm import tqdm from datetime import datetime import matplotlib.pyplot as plt from itertools import product from scipy.special import gamma from sklearn.neighbors import KDTree from scipy.spatial import cKDTree from collections import defaultdict Формирование обучающей выборки для кластеризации по Уишарту In [4]: new train tot3 = np.load('new train tot3.npy') print(new train tot3.shape) new train tot3 # список наблюдений для обучающей выборки (3912, 5)Out[4]: array([['SPI', '2020-09-23', '720', '34', '1'], ['GLSI', '2020-12-09', '666', '54', '1'], ['AACG', '2021-02-04', '627', '55', '1'], ['EVH', '2020-08-20', '486', '28', '0'], ['MXIM', '2019-10-29', '488', '30', '0'], ['PEI', '2020-10-08', '519', '28', '0']], dtype='<U11') Нельзя, чтобы обучающие выборки по каждому наблюдению слиплись, иначе шаблоны пройдутся не только по своим наблюдениям но и по чужим, выходя за границы, поэтому будем резать и сохранять в файлы (опер памяти, чтобы все это делать сразу не хватит). from itertools import product patterns1=np.array(list(product(np.arange(1, 11, 1), repeat=4)), dtype=np.uint8) np.random.seed(123) # возьмем 1000 шаблонов, слишком долго будет вычисляться, если возьмем все 10 тыс. patterns = patterns1[np.random.choice(patterns1.shape[0], 1000, replace=False)] patterns Out[5]: array([[3, 7, 6, 7], [1, 5, 5, 6], [10, 6, 1, 6], [3, 7, 5, 8], 4, 10], [1, 5, [9, 5, 1, 5]], dtype=uint8) #np.save('patterns 1000 patterns for wishart.npy', patterns) # отдельные суммы расстояний для каждого шаблона q = np.sum(patterns[:,:0],axis=1).reshape((patterns.shape[0],1)) a = np.sum(patterns[:,:1],axis=1).reshape((patterns.shape[0],1)) b = np.sum(patterns[:,:2],axis=1).reshape((patterns.shape[0],1)) c = np.sum(patterns[:,:3],axis=1).reshape((patterns.shape[0],1)) d = np.sum(patterns[:,:4],axis=1).reshape((patterns.shape[0],1)) np.max(d) # макс обрезка по шаблону с конца ряда Out[8]: 39 # обучающая выборка для каждого дневного наблюдения из списка класса min ts = -0.9274193548387099 # для сдвига всей выборки вверх, отриц значения здесь недопустимы for i in tqdm(range(len(new train tot3))): # откроем каждое наблюдение из обуч выборки ticker name = new train tot3[i][0] day = new train tot3[i][1] parent dir = 'C:/Users/Kuanysh/Downloads/pump and dump/train3' file_min = os.path.join(parent_dir, ticker_name + '_'+ day +'.csv') train data = pd.read csv(file min).to numpy().ravel() # сделаем ряд стационарным через np.diff, очень важно ее сделать ts diff = np.diff(train data) ts diff = ts diff - np.full((ts diff.shape[0],), min ts) # но также важно сдвинуть все ряды вверх # запишем наблюдения по каждому паттерну в тензор и правильно переформатируем (спец трюк) sds=[] for i in range(len(ts diff)-np.max(d)): sd=np.hstack((ts_diff[q+i], ts_diff[a+i], ts_diff[b+i], ts_diff[c+i], ts_diff[d+i])) sds.append(sd) sds=np.array(sds) sds=np.stack(sds,1) # запишем эти тензоры в файлы parent dir2='C:/Users/Kuanysh/Downloads/pump and dump/train npys' np dir = os.path.join(parent dir2, ticker name + ' '+ day +'.npy') with open(np dir, 'wb') as f: np.save(f, sds) 3912/3912 [02:52<00:00, 2 100%| 2.65it/s] # создадим подпапки для кадого из 1000 шаблонов, чтобы туда собирать срезанные по каждому шаблону обуч-ие выбок for n in range(1000): # 1000 шаблонов parent dir = 'C:/Users/Kuanysh/Downloads/pump and dump/1000patterns' path = os.path.join(parent dir, str(n)) os.makedirs(path) for i in tqdm(range(len(new_train_tot3))): # вскроем обуч выборки по каждому наблюдению ticker name = new train tot3[i][0] day = new train tot3[i][1] parent dir2='C:/Users/Kuanysh/Downloads/pump and dump/train npys' np dir = os.path.join(parent dir2, ticker name + ' '+ day +'.npy') train ts = np.load(np dir, mmap mode='r+') sh = train ts.shape[0] # 1000 шаблонов # теперь поместим срезанные по каждому шаблону обуч-ие выборки в их подпапки parent dir3 = 'C:/Users/Kuanysh/Downloads/pump and dump/1000patterns' for k in range(sh): path1 = os.path.join(parent dir3, str(k)) path2 = os.path.join(path1, ticker name + ' ' + day + ' ' + str(k) +'.npy') with open(path2, 'wb') as f: np.save(f, train ts[k]) 100%| | 3912/3912 [28:59<00:00, 2.25it/s] for i in tqdm(range(1000)): # 1000 шаблонов # вскроем каждую подпаку и соберем все обучающие выборки и объеденим наблюдения по каждому шаблону parent dir = 'C:/Users/Kuanysh/Downloads/pump and dump/1000patterns' path = os.path.join(parent dir, str(i)) files = glob.glob(path + '/*') # обнаружим все файлы list pats = [] for file in files: one train pat = np.load(file, mmap mode='r+').tolist() list_pats.extend(one_train_pat) # объединение по каждому шаблону # объединенные наблюдения поместим в отдельную папку directory = 'C:/Users/Kuanysh/Downloads/pump and dump/total patterns' with open(os.path.join(directory, 'tot'+ str(i)), "w") as fp: json.dump(list pats, fp) 1000/1000 [4:21:50<00:00, 1 5.71s/it] np.random.seed(12345) choice = np.random.choice(1473014, 400000, replace=False) #np.save('choice 500000.npy', choice) choice Out[2]: array([727536, 633366, 16414, ..., 940512, 1279791, 252998]) # теперь наконецто объединим все шаблоны воедино, и таким образом, вся обучающая выборка будет правильно сформи list tot = [] for i in tqdm(range(1000)): # 1000 шаблонов parent dir = "C:/Users/Kuanysh/Downloads/pump and dump/total patterns" path = os.path.join(parent dir, 'tot' + str(i)) with open(path, "r") as fp: tot = json.load(fp) tot = np.array(tot) rnd tot = tot [choice] list tot.append(rnd tot) 1000/1000 [1:11:52<00:00, 100%| 4.31s/it] In [4]: list_tot_=np.array(list_tot, dtype=np.float32) #np.save('tot all.npy', list tot) list tot all = np.load("tot all.npy", mmap mode='r+') Кластеризация по Уишарту: # кластеризуем обучающую выборку по Уишарту: class Wishart: def init (self, wishart neighbors, significance level): self.wishart neighbors = wishart neighbors # Number of neighbors self.significance level = significance level # Significance level def fit(self, X): #kdt = KDTree(X, metric='euclidean') kdt = cKDTree(X, leafsize=40) #add one because you are your neighb. distances, neighbors = kdt.query(X, k = self.wishart_neighbors + 1, p=2, n_jobs=-1) #distances, neighbors = kdt.query(X, k = self.wishart_neighbors + 1, return_distance = True) neighbors = neighbors[:, 1:] distances = distances[:, -1] indexes = np.argsort(distances) size, dim = X.shapeself.object_labels = np.zeros(size, dtype = int) - 1 #index in tuple #min_dist, max_dist, flag_to_significant self.clusters = np.array([(1., 1., 0)])self.clusters_to_objects = defaultdict(list) #print('Start clustering') for index in indexes: neighbors clusters =\ np.concatenate([self.object_labels[neighbors[index]], self.object_labels[neighbors[index]]]) unique_clusters = np.unique(neighbors_clusters).astype(int) unique_clusters = unique_clusters[unique_clusters != -1] if len(unique_clusters) == 0: self._create_new_cluster(index, distances[index]) max cluster = unique clusters[-1] min_cluster = unique_clusters[0] if max_cluster == min_cluster: if self.clusters[max_cluster][-1] < 0.5:</pre> self._add_elem_to_exist_cluster(index, distances[index], max_cluster) self._add_elem_to_noise(index) else: my clusters = self.clusters[unique_clusters] flags = my_clusters[:, -1] if np.min(flags) > 0.5: self._add_elem_to_noise(index) else: significan = np.power(my_clusters[:, 0], -dim) - np.power(my_clusters[:, 1], -dim) significan *= self.wishart_neighbors significan /= size significan /= np.power(np.pi, dim / 2) significan *= gamma(dim / 2 + 1) significan_index = significan >= self.significance_level significan_clusters = unique clusters[significan index] not_significan_clusters = unique_clusters[~significan_index] significan_clusters_count = len(significan_clusters) if significan_clusters_count > 1 or min_cluster == 0: self._add_elem_to_noise(index) self.clusters[significan clusters, -1] = 1 for not_sig_cluster in not_significan_clusters: if not sig cluster == 0: continue for bad index in self.clusters to objects[not sig cluster]: self. add elem to noise(bad index) self.clusters_to_objects[not_sig_cluster].clear() else: for cur_cluster in unique_clusters: if cur_cluster == min_cluster: continue for bad index in self.clusters to objects[cur cluster]: self._add_elem_to_exist_cluster(bad_index, distances[bad_index], min_cluste self.clusters to objects[cur cluster].clear() self. add_elem_to_exist_cluster(index, distances[index], min_cluster) return self.clean_data() def clean data(self): unique = np.unique(self.object_labels) index = np.argsort(unique) **if** unique[0] != 0: index += 1 true cluster = {unq : index for unq, index in zip(unique, index)} result = np.zeros(len(self.object_labels), dtype = int) for index, unq in enumerate(self.object_labels): result[index] = true cluster[unq] return result def _add_elem_to_noise(self, index): self.object labels[index] = 0 self.clusters_to_objects[0].append(index) def create new cluster(self, index, dist): self.object_labels[index] = len(self.clusters) self.clusters to objects[len(self.clusters)].append(index) self.clusters = np.append(self.clusters, [(dist, dist, 0)], axis = 0) def add elem to exist cluster(self, index, dist, cluster label): self.object_labels[index] = cluster_label self.clusters to objects[cluster label].append(index) self.clusters[cluster_label][0] = min(self.clusters[cluster_label][0], dist) self.clusters[cluster_label][1] = max(self.clusters[cluster_label][1], dist) claster labs Wish2 = [Wishart(22, 0.2).fit(list tot all[i]) for i in tqdm(range(len(patterns)))] 0%| | 0/1000 [00:00 <?, ?it/s]<ipython-input-3-59bae2c59e79>:12: DeprecationWarning: The n_jobs argument has been renamed "worker s". The old name "n_jobs" will stop working in SciPy 1.8.0. distances, neighbors = kdt.query(X, k = self.wishart_neighbors + 1, p=2, n_jobs=-1) <ipython-input-3-59bae2c59e79>:52: RuntimeWarning: divide by zero encountered in power significan = np.power(my_clusters[:, 0], -dim) - np.power(my_clusters[:, 1], -dim) <ipython-input-3-59bae2c59e79>:52: RuntimeWarning: invalid value encountered in subtract significan = np.power(my clusters[:, 0], -dim) - np.power(my clusters[:, 1], -dim) 100%| | 1000/1000 [2:19:03<00:00, 8.34s/it] np.save('claster_labs_Wish2.npy', claster_labs_Wish2) # получим центры кластеров sds = list_tot_all sds_Wishart = [] for i in range(len(patterns)): $sds_i_W = []$ for cl in set(claster_labs_Wish2[i]): **if** cl != 0: sds_i_W.append(np.mean(sds[i][claster_labs_Wish2[i] == cl], axis=0)) sds_i_W = np.array(sds_i_W) sds_Wishart.append(sds_i_W) # выровняем ряды по одной длине и заполним ненулевыми значениями from tslearn.utils import to time series dataset N = np.array([0.6666], dtype=np.float32)train_pat = to_time_series_dataset(sds_Wishart) train_pat[np.isnan(train_pat)]=N[0] train pat.shape Out[11]: (1000, 2441, 5) train_pat[200][200:500] Out[16]: array([[0.88575757, 0.9274323 , 0.92741054, 0.92741394, 0.92739296], [0.92738473, 0.93124068, 0.93088257, 0.92741936, 0.92741936],[0.92213601, 0.92744416, 0.92741948, 0.92741531, 0.93233728],[0.99365377, 0.91254997, 0.92790794, 0.92684883, 0.92741936], [0.92740697, 0.92776537, 0.92739129, 1.03847671, 0.92732012],[0.89711171, 0.98803461, 0.92741936, 0.92741936, 0.92726785]])np.save('train_pat_1000_wishart_totlist3.npy', train_pat)