

Part1

```
1. def pq(data, P, init_centroids, max_iter):
2.     data = np.array(data, dtype='float32')
3.     init_centroids = np.array(init_centroids, dtype='float32')
4.     N, M = data.shape
5.     K = 256
6.     MP = int(M / P)
7.     result_centroids = np.array([])
8.
9.     for i in range(P):
10.         sliced_data = data[:, i * MP:(i + 1) * MP] # sliced data by P
11.         centroids = init_centroids[i,].copy() #copy another centroids
12.
13.         for _ in range(max_iter):
14.             ### change to L1 distance by using distance_matrix methol with p
               arameter 1
15.             row_dis = distance_matrix(sliced_data, centroids, 1)
16.             row_label = np.argmin(row_dis, axis=1) # then sort the distance
17.
18.             for j in range(K):
19.                 index_list = np.where(row_label == j)[0] # search the point
                   that belong to centroid
20.
21.                 temp_data = [sliced_data[index] for index in index_list]
22.                 if len(temp_data) != 0: # if no point belongs to centroid, c
                   entroid should stay the same
23.                     ### because change to L1 distance, kmedian should be the
                       beat methol instead of kmeans
24.                     centroids[j] = np.median(temp_data, axis=0) #updata the
                           centroids
25.
26.
27.             row_dis = distance_matrix(sliced_data, centroids, 1)
28.             row_label = np.argmin(row_dis, axis=1)
29.
30.             if i == 0:
31.                 result_label = np.array([row_label]).T.copy()
32.                 result_centroids = np.array([centroids.copy()])
33.             else:
34.                 temp = np.append(result_centroids, centroids)
35.                 dim = result_centroids.shape
36.                 result_centroids = temp.reshape(dim[0] + 1, dim[1], dim[2])
37.                 result_label = np.column_stack((result_label, row_label)).astype
                   (np.uint8)
38.
39.     return result_centroids, result_label
```

The implementation of part 1 already written on the annotations.

About change to L1 distance, there are two point:

1. use `scipy.spatial.distance_matrix` to calculate L1 distance
2. use k-median instead of k-means

Part2

```
1. def query(queries, codebooks, codes, T):
2.     # print(codes)
3.     Q, M = queries.shape
4.     P, K, _ = codebooks.shape
5.     N, P = codes.shape
6.     MP = int(M / P)
7.     candidates = []
8.
9.     inverted_index = {}
10.    for i in range(N): # use a dict with tuple as key to search index instead of np.where
11.        if tuple(codes[i]) not in inverted_index.keys():
12.            inverted_index[tuple(codes[i])] = [i]
13.        else:
14.            inverted_index[tuple(codes[i])].append(i)
15.
16.    one_list = np.array([np.zeros((P)).astype(np.int32)] * P)
17.    for i in range(P): # one_list = [[0,1],[1,0]] if dim == 2, [[1,0,0],[0,1,0],[0,0,1]] if dim==3.....dim==4
18.        one_list[i][i] = 1
19.
20.    for q in range(Q):
21.        for i in range(P):
22.            sliced_data = queries[q, i * MP:(i + 1) * MP]
23.            row_dis = distance_matrix([sliced_data], codebooks[i], 1)
24.
25.            ### sort distance
26.            sort_index = row_dis[0].argsort()
27.            if i == 0:
28.                sorted_matrix = np.array([sort_index])
29.                dis_matrix = np.array(row_dis)
30.            else:
31.                sorted_matrix = np.append(sorted_matrix, [sort_index], axis=0)
32.                dis_matrix = np.append(dis_matrix, row_dis, axis=0)
33.
34.            distance_dict = {(0,) * P: distance(index_conv([0] * P, sorted_matrix), dis_matrix)} # add(0..0) to dis_dict
35.            used_index = {} # store used_index in dict.keys(), because tuple can be hashed and it's fast
36.            temp_candidate = set()
37.
38.            while len(temp_candidate) < T:
39.                minimal_index = min(distance_dict, key=distance_dict.get) # find the minimal distance
40.                ###calculate the invert index
```

```
41.         minimal_invert_index = tuple([sorted_matrix[i][minimal_index[i]]
42.         for i in range(len(minimal_index))])
43.         if minimal_invert_index in inverted_index.keys(): # use dict.ke
44.             y to check if used because it's fast
45.             add_set = tuple(inverted_index[minimal_invert_index])
46.             temp_candidate = temp_candidate.union(add_set) # add to tem
47.             p_candidate
48.             distance_dict.pop(minimal_index) # delete this cell after used
49.             used_index[minimal_index] = True # add the used cell into used_
50.             index
51.             for one in one_list:
52.                 new_index = tuple(one + list(minimal_index)) # get the neigh
53.                 borhood cell index
54.                 if new_index not in used_index.keys() and max(new_index) < 2
55.                 56: # check if used or out of index
56.                 distance_dict[new_index] = distance(index_conv(new_index
57.                 , sorted_matrix), dis_matrix) # if not, add it
58.                 candidates.append(temp_candidate)
59.             return candidates
```

The implementation of part 1 already written on the annotations.

About extended the algorithm 3.1 to a more general case with $P > 2$:

Just extend the index length, if $p == 4$, the first index ==
[0,0,0,0], and the other work will check every dim in thid index.

About how you efficiently retrieve the candidates:

Store most of data in tuple as key in dict because it can be
hashed and search it in dict will be fast and efficiently in $O(1)$
time