

```

1 # -*- coding: utf-8 -*-
2 '''Created by Tai Dinh
3 This file is used for k-CMM algorithm: clustering mixed
4 numeric and categorical data with missing values
5 '''
6 from __future__ import division
7 import sys
8 import numpy as np
9 import evaluation
10 from collections import defaultdict
11 from scipy import *
12 import c45
13 import math
14 import matplotlib.pyplot as plt
15 import pandas as pd
16 import itertools
17 # For measuring the time for running program
18 # source: http://stackoverflow.com/a/1557906/6009280
19 # or https://www.w3resource.com/python-exercises/python-
20 # basic-exercise-57.php
21 # import atexit
22 from time import time, strftime, localtime
23 from datetime import timedelta
24 # For measuring the memory usage
25 import tracemalloc
26
27
28 def get_max_value_key(dic):
29     '''Fast method to get key for maximum value in dict'''
30     v = list(dic.values())
31     k = list(dic.keys())
32     return k[v.index(max(v))]
33
34 ...
35 This function is used to calculate the dissimilarity
36 between 2 attributes x and y at the iattr(d)
37 using the global_attr_freq in which the global_axttr_freq[i]
38 ] [x] is the frequency of the attribute at
39 i with value x in the whole samples:
40
41
42 def attr_dissim(x, y, iattr, global_attr_freq):
43     ...
44     Dissimilarity between 2 categorical attributes x and y
45     at the attribute iattr, i.e
46         dis(x, y) = 1 - 2 * log(P{x, y}) / (log(P{x}) + log(P{y})

```

```

45 })) ...
46     ...
47     if (global_attr_freq[iattr][x] == 1.0) and (
48         global_attr_freq[iattr][y] == 1.0):
49         return 0
50     if x == y:
51         numerator = 2 * math.log(global_attr_freq[iattr][x])
52     else:
53         numerator = 2 * math.log((global_attr_freq[iattr][x]
54 ] + global_attr_freq[iattr][y]))
55         denominator = math.log(global_attr_freq[iattr][x]) +
56         math.log(global_attr_freq[iattr][y]) #Noted by Tai Dinh,
57         Equation 21, page 124
58     return 1 - numerator / denominator
59
60 """information theoretic based similarity measure
61 This function is used to calculate the dissimilarity
62 between a centroid and a vector a
63 """
64
65 def vector_matching_dissim(centroid, categorical, a,
66 global_attr_freq):
67     # Get distance between a centroid and a
68     distance = 0.
69     attrIndex = 0
70     for ic, curc in enumerate(centroid):
71         if ic in categorical:
72             keys = curc.keys()
73             for key in keys:
74                 distance += curc[key] * attr_dissim(key, a[ic],
75 attrIndex, global_attr_freq)
76             attrIndex += 1
77         else:
78             tmp = float(circ)-float(a[ic])
79             distance += pow(tmp,2)
80     return distance
81
82 """
83 This function is used to calculate the distances between
84 centroid clusters and a data point, using the
85 global_attr_freq.
86 global_axttr_freq[i][x] is the probability of the attribute
87 at position i and value x on the whole samples.
88 categorical is the set of categorical attributes in the
89 data point.
90 """
91
92 def vectors_matching_dissim(vectors, categorical, a,
93 global_attr_freq):

```

```

82     '''Get nearest vector in vectors to a'''
83     min = np.Inf
84     min_clust = -1
85     for clust in range(len(vectors)):
86         distance = vector_matching_dissim(vectors[clust],
87             categorical, a, global_attr_freq)
88         if distance < min:
89             min = distance
90             min_clust = clust
91     return min_clust, min
92 ...
93 This function is used to transfer a vector point from this
94 cluster (from_clust) to another cluster (to_clust)
95 ipoint is the index of vector in the samples.
96 membership[cluster_index, ipoint] = 1 means vector with
97 the index ipoint belongs to the cluster_index.
98
99 cl_attr_freq[cluster_index][iattr][curattr] is the
100 frequency of the attribute having value curattr at iattr
101 in cluster cluster_index.
102 In fact, k are kept in the cl_attr_freq instead of k/N
103 such that k is the number of appearance of attribute at
104 the iattr with value curattr,
105 N is the number of data objects in the cluster. The reason
106 is k and N are probably change, in this case
107 recalculating k/N is more complex than k.
108
109 Note that global_attr_freq stores frequency (k/N) because
110 it only needs one time to calculate and values keep
111 permanently.
112 ...
113 def move_point_between_clusters(point, ipoint, to_clust,
114     from_clust,
115     cl_attr_freq, membership, categorical):
116     '''Move point between clusters, categorical attributes
117     ...
118     membership[to_clust, ipoint] = 1
119     membership[from_clust, ipoint] = 0
120     # Update frequencies of attributes in clusters
121     attrIndex = 0
122     for iattr, curattr in enumerate(point):
123         if iattr in categorical:
124             cl_attr_freq[to_clust][attrIndex][curattr] +=
125             1
126             cl_attr_freq[from_clust][attrIndex][curattr]
127             -= 1
128             attrIndex +=1
129     return cl_attr_freq, membership

```

```

117
118 def matching_dissim(a, b):
119     '''Simple matching dissimilarity function'''
120     return np.sum(a != b, axis=1)
121
122 ...
123 Randomly initialize vectors in X into clusters
124 ...
125
126 def _init_clusters(X, centroids, n_clusters, nattrs,
127     npoints, verbose, categorical):
128     # __INIT_CLUSTER__
129     # if verbose:
130     #     print("Init: Initialzing clusters")
131     membership = np.zeros((n_clusters, npoints), dtype='int64')
132     # cl_attr_freq is a list of lists with dictionaries
133     # that contain the
134     # frequencies of values per cluster and attribute.
135     cl_attr_freq = [[defaultdict(int) for i in categorical]
136                     for _ in range(n_clusters)]
137     for ipoint, curpoint in enumerate(X):
138         # Initial assignment to clusters
139         clust = np.argmin(matching_dissim(centroids,
140                                         curpoint))
141         membership[clust, ipoint] = 1
142         # Count attribute values per cluster
143         attrIndex = 0
144         for iattr, curattr in enumerate(curpoint):
145             if iattr in categorical:
146                 cl_attr_freq[clust][attrIndex][curattr] +=
147                     1
148                 attrIndex += 1
149
150     # Move random selected point from largest cluster to
151     # empty cluster if exists
152     for ik in range(n_clusters):
153         if sum(membership[ik, :]) == 0:
154             from_clust = membership.sum(axis=1).argmax()
155             choices = \
156                 [ii for ii, ch in enumerate(membership[
157                     from_clust, :]) if ch]
158             rindex = np.random.choice(choices)
159             # Move random selected point to empty cluster
160             cl_attr_freq, membership =
161             move_point_between_clusters(
162                 X[rindex], rindex, ik, from_clust,
163                 cl_attr_freq, membership, categorical)
164
165
166

```

```

157     return cl_attr_freq, membership
158
159 ...
160 This function is used to calculate the lambda as the
161 formula in slide page 15
162 cl_attr_freq[iattr][curattr] is the probability of
163 attribute at the iattr having value curattr in cluster
164 clust_members is the number of data objects in the cluster
165 ...
166
167 def cal_lambda(cl_attr_freq, clust_members):
168     '''Re-calculate optimal bandwidth for each cluster'''
169     if clust_members <= 1:
170         return 0.
171
172     numerator = 0.
173     denominator = 0.
174
175     for iattr, curattr in enumerate(cl_attr_freq):
176         n_ = 0.
177         d_ = 0.
178         keys = curattr.keys()
179         for key in keys:
180             n_ += math.pow(1.0 * curattr[key] /
181                           clust_members, 2)
182             d_ += math.pow(1.0 * curattr[key] /
183                           clust_members, 2)
184             numerator += (1 - n_)
185             denominator += (d_ - 1.0 / (len(keys)))
186
187             # print denominator
188             # assert denominator != 0, "How can denominator equal
189             # to 0?"
190             if clust_members == 1 or denominator == 0:
191                 return 0
192             result = (1.0 * numerator) / ((clust_members - 1) *
193                                         denominator)
194             if result < 0:
195                 return 0;
196             if result > 1:
197                 return 1
198             return (1.0 * numerator) / ((clust_members - 1) *
199                                         denominator)
200
201 def _cal_global_attr_freq(X, npoints, nattrs, categorical):
202     :
203         # global_attr_freq is a list of lists with
204         # dictionaries that contain the
205         # frequencies of attributes.

```

```

197     global_attr_freq = [defaultdict(float) for _ in
 categorical]
198     # global_attr_freq = {}
199     # for i in categorical:
200     #     global_attr_freq[i] = defaultdict(float)
201
202     for ipoint, curpoint in enumerate(X):
203         attrIndex = 0
204         for iattr, curattr in enumerate(curpoint):
205             if iattr in categorical:
206                 global_attr_freq[attrIndex][curattr] += 1.
207                 attrIndex += 1
208             attrIndex = 0
209             for iattr in range(nattrs):
210                 if iattr in categorical:
211                     for key in global_attr_freq[attrIndex].keys():
212                         global_attr_freq[attrIndex][key] /= npoints
213                         attrIndex += 1
214     return global_attr_freq
215
216 ...
217 This function is used to calculate the centroid center at
each attribute.
218
219 * ldb is lambda
220 * cl_attr_freq_attr is cl_attr_freq[clust][iattr], is the
number of attribute at the index iattr in the cluster
clust.
221 * clust_members is the number of data objects in the
cluster
222 * global_attr_count is the number of attribute at the
index iattr in the whole dataset X.
223 ...
224
225 def cal_centroid_value(lbd, cl_attr_freq_attr,
cluster_members, attr_count):
226     '''Calculate centroid value at iattr'''
227     assert cluster_members >= 1, "Cluster has no member,
why?"
228
229     keys = cl_attr_freq_attr.keys()
230     vjd = defaultdict(float)
231     for odl in keys:
232         vjd[odl] = lbd / attr_count + (1 - lbd) * (1.0 *
cl_attr_freq_attr[odl] / cluster_members) #Noted by Tai
Dinh equation 12, page 121
233     return vjd
234
235 def cal_mean_value(X, indexAttr):

```

```

236     # print(X.iloc[:,indexAttr])
237     meanValue = mean(np.asarray(X.iloc[:,indexAttr], dtype
= float))
238     return meanValue
239
240 ...
241 This function is the loop for the k-CMM algorithm
242 For each vector curpoint with the index ipoint in X, the
purpose is to find the nearest centroid with this vector.
243 ...
244
245 def _k_CMM_iter(X, categorical, centroids, cl_attr_freq,
membership, global_attr_freq, lbd, use_global_attr_count):
246     '''Single iteration of the k_CMM clustering algorithm
...
247     moves = 0
248     for ipoint, curpoint in enumerate(X):
249         clust, distance = vectors_matching_dissim(
centroids, categorical, curpoint, global_attr_freq)
250         if membership[clust, ipoint]:
251             # Sample is already in its right place
252             continue
253
254             # Move point and update old/new cluster
frequencies and centroids
255             ...
256             moves is the number of moving vectors between
cluster
257             old_clust is the old index of vector curpoint
258             ...
259             moves += 1
260             old_clust = np.argwhere(membership[:, ipoint])[0][
0]
261
262             ...
263             Move vector with index ipoint from old_clust to
clust, meanwhile recalculate the probability of attributes
in the corresponding clusters.
264             ...
265             cl_attr_freq, membership =
move_point_between_clusters(
266                 curpoint, ipoint, clust, old_clust,
cl_attr_freq, membership,categorical)
267
268             # In case of an empty cluster, reinitialize with a
random point
269             # from the largest cluster.
270             ...
271             After moving vectors from old_clust to new_clust,
if the old_clust is empty,

```

```

272     then get an arbitrary vector from the largest
273     cluster to this cluster to avoid empty clusters.
274     ...
275     if sum(membership[old_clust, :]) == 0:
276         from_clust = membership.sum(axis = 1).argmax()
277         choices = \
278             [ii for ii, ch in enumerate(membership[
279                 from_clust, :]) if ch]
280         rindex = np.random.choice(choices)
281
282         cl_attr_freq, membership =
283         move_point_between_clusters(
284             X[rindex], rindex, old_clust, from_clust,
285             cl_attr_freq, membership, categorical)
286
287         # Re-calculate lambda of changed centroid
288         for curc in (clust, old_clust):
289             lbd[curc] = cal_lambda(cl_attr_freq[curc], sum(
290                 membership[curc, :]))
291
292         # Update new and old centroids by choosing mode of
293         # attribute.
294         attrIndex = 0
295         for iattr in range(len(curpoint)):
296             if iattr in categorical:
297                 for curc in (clust, old_clust):
298                     cluster_members = sum(membership[curc,
299                         :])
300                     if use_global_attr_count:
301                         centroids[curc][iattr] =
302                             cal_centroid_value(lbd[curc], cl_attr_freq[curc][attrIndex],
303                                 cluster_members, len(global_attr_freq[attrIndex]))
304                     else:
305                         attr_count = len(cl_attr_freq[curc][attrIndex].keys())
306                         centroids[curc][iattr] =
307                             cal_centroid_value(lbd[curc], cl_attr_freq[curc][attrIndex],
308                                 cluster_members, attr_count)
309                         attrIndex += 1
310             else:
311                 comSetDF = pd.DataFrame(X)
312                 centroids[curc][iattr] = cal_mean_value(
313                     comSetDF, iattr)
314
315             return centroids, moves, lbd
316
317
318     ...
319
320 This function is used to calculate the sum of distances
321 between vectors inside X and centroids of clusters after

```

```

307 each step.
308 Labels is the label of vector in X, labels[x] = c means
vector with index is x is belonged to the cluster that has
its index is c.
309 Cost is the sum of dissimilarity.
310 ...
311
312 def _labels_cost(X, categorical, centroids,
global_attr_freq):
313     ...
314     Calculate labels and cost function given a matrix of
points and
315         a list of centroids for the k-CMM algorithm.
316         ...
317
318     npoints, nattrs = X.shape
319     cost = 0.
320     labels = np.empty(npoints, dtype = 'int64')
321     for ipoint, curpoint in enumerate(X):
322         ...
323             For every vector ipoint (its value is curpoint) in
X, find out nearest cluster with it by using the function
vectors_matching_dissim.
324             Then calculate the distance by using the
dissimilarity for a mixed object and cluster)
325             ...
326             clust, diss = vectors_matching_dissim(centroids,
categorical, curpoint, global_attr_freq)
327             assert clust != -1, "Why there is no cluster for
me?"
328             labels[ipoint] = clust
329             cost += diss
330
331     return labels, cost
332
333 # This function is used to split numeric and categorical
attributes in the complete dataset to two subsets.
334 def splitNumCat(x, categorical):
335     xNum = x.loc[:, [ii for ii in range(x.shape[1])
if ii not in categorical]]
336     xCat = x.loc[:, categorical]
337     return xNum, xCat
338
339
340 # This function is used to impute numeric attributes
341 def numericImputation(inSetDFNum, comSetDFNum, IDList):
342     inObjectNum = inSetDFNum.loc[0]
343     # listOfObjects = []
344     # for i in range(0, len(IDList)):
345     #     listOfObjects.append(comSetDFNum.iloc[[IDList[i
]]])

```

```

346     listOfObjects = comSetDFNum.iloc[IDList]
347     imputeValues = []
348     # Calculate the mean of each numeric attributes in the
349     # list of objects
350     for j in range(0, comSetDFNum.shape[1]):
351         imputeValues.append(round(mean(pd.to_numeric(
352             listOfObjects.iloc[:, j])), 2))
353     indexList = inObjectNum.index.T.values
354     l=0
355     for k in range(0, len(indexList)):
356         if inObjectNum[indexList[k]] == '?':
357             inObjectNum[indexList[k]] = imputeValues[l]
358             l += 1
359     return inObjectNum
360
361 # This function is used to concatenate the incomplete and
362 # complete objects into a mixed object
363 def concatenationTwoObjects(inObject, comObject):
364     mixedObject = inObject.append(comObject, ignore_index=
365     False)
366     # mixedObjectDF = mixedObject.reset_index(drop=False)
367     indexList = np.sort(mixedObject.index.T.values)
368     mixedObject = mixedObject.loc[indexList]
369     return mixedObject
370
371 # This function is used to calculate Euclidean distance
372 # for numeric attributes.
373 def euclidean_dissim(a, b, **_):
374     """Euclidean distance dissimilarity function"""
375     # if np.isnan(a).any() or np.isnan(b).any():
376     #     raise ValueError("Missing values detected in
377     # numerical columns.")
378     return np.sum((a - b) ** 2, axis=1)
379
380 # This function is to calculate the means of all numeric
381 # columns in the complete numeric dataset.
382 def meansColumn(comSetDFNum, column):
383     return round(mean(pd.to_numeric(comSetDFNum.iloc[:, column])), 2)
384
385 ...
386
387 The k-CMM algorithm for clustering a mixed dataset with
388 missing values X into k clusters.
389 * init is the initial method for the algorithm (in this
390 case, we used randomness)
391 * n_init is the number of running algorithm with different
392 initialization
393 * max_iter is the number of executing the algorithm
394 * verbose == 1 print information, 0 is otherwise
395 ...

```

```

385
386 def perform_kCMM(X, categorical, n_clusters, init, verbose
387   , use_global_attr_count):
388     '''k-CMM algorithm'''
389     inSetDF, comSetDF = readData(X)
390     # inSet = np.asanyarray(inSetDF)
391     comSet = np.asanyarray(comSetDF)
392
393     inSetDFNum, inSetDFCat = splitNumCat(inSetDF,
394       categorical)
395     comSetDFNum, comSetDFCat = splitNumCat(comSetDF,
396       categorical)
397
398     npoints, nattrs = comSet.shape
399     assert n_clusters < npoints, "More clusters than data
400     points?"
401
402     all_centroids = []
403     all_labels = []
404     all_costs = []
405
406     if init == 'random':
407         seeds = np.random.choice(range(npoints),
408           n_clusters)
409         centroids = comSet[seeds]
410     else:
411         raise NotImplementedError
412
413     itr = 0
414     converged = False
415     cost = np.Inf
416
417     # For each categorical object in incomplete dataset
418     S_2^c do
419       # for idx in range(0, len(inSetDF)):
420       while len(inSetDFCat) != 0:
421         inObjectCat = inSetDFCat.loc[0]
422         #Create a set that contains all missing attributes
423         # of an object, initialize an empty set
424         inSetAttr = []
425         comSetAttr = []
426         # Create a set that contains decision tree for
427         # missing attributes
428         setDT = {}
429         # Create a set that contains complete dataset
430         # after change index of columns
431         setComSetDFCatAfterChangeIndex = {}
432         #Put missing attributes of object into the setAttr
433         tmpCount = 0
434         inObjectCatAfterChangeIndex = {}

```

```

426         indexList = inObjectCat.index.T.values
427         for k in range(0, len(indexList)):
428             if inObjectCat[indexList[k]] == '?':
429                 inSetAttr.append(indexList[k])
430             else:
431                 inObjectCatAfterChangeIndex[str(tmpCount)]
432                     = inObjectCat[indexList[k]]
433                     tmpCount +=1
434                     # Store index of complete attribute to
435                     build the decision tree
436                     comSetAttr.append(indexList[k])
437                     # print("Finish finding missing attributes for
438                     object "+ str(index))
439                     # Check if the object contains no missing
440                     categorical values. In case of missing values in numeric
441                     attributes
442                     if (len(inSetAttr) == 0):
443                         inObjectNum = inSetDFNum.loc[0]
444                         indexList = inObjectNum.index.T.values
445                         for k in range(0, len(indexList)):
446                             if inObjectNum[indexList[k]] == '?':
447                                 inObjectNum[indexList[k]] =
448                                     meansColumn(comSetDFNum, k)
449                         else:
450                             for index in inSetAttr:
451                                 # This step is to move the missing
452                                 attribute as the class attribute in the comSetDFCat
453                                 colsAfterChangeIndex = comSetAttr.copy()
454                                 colsAfterChangeIndex.append(index)
455                                 tmpList = comSetDFCat[colsAfterChangeIndex
456 ]
457                                 comSetDFCatAfterChangeIndex = comSetDFCat[
458                                     colsAfterChangeIndex]
459                                 # print(comSet)
460                                 tmpList = tmpList.reset_index(drop=True).
461                                     values.tolist()
462                                 # Build a decision tree for a missing
463                                 attribute of an object in incomplete set
464                                 decisionTree = c45.growDecisionTreeFrom(
465                                     tmpList)
466                                 # c45.plot(decisionTree)
467                                 # Add decision tree to the set of DT
468                                 setDT[index] = decisionTree
469                                 setComSetDFCatAfterChangeIndex[index] =
470                                     comSetDFCatAfterChangeIndex
471                                 # print("Finish building DTs for missing
472                                 attributes")
473
474                                 #Create a table that contains all correlated
475                                 objects in a DT

```

```

461         tableOfObjects = []
462         for key, value in setDT.items():
463             # c45.plot(value)
464             # getAllLeavesInDT(value)
465             comSetDFCatAfterChangeIndex =
466                 setComSetDFCatAfterChangeIndex[key]
467                 listBestPath = findSuitableLeave(value,
468                     inObjectCatAfterChangeIndex)
469                     for path in listBestPath:
470                         listOfObject =
471                             assignCompleteObjectsIntoDT(path,
472                             comSetDFCatAfterChangeIndex.T.reset_index(drop=True).T)
473                             tableOfObjects.append(listOfObject)
474                             # Merge correlated complete objects with
475                             # incomplete object into one collection
476                             IDList = list(itertools.chain.from_iterable(
477                             tableOfObjects))
478                             IDList = list(dict.fromkeys(IDList))
479                             IDList.sort()
480                             listOfObjects = []
481                             for i in IDList:
482                                 listOfObjects.append(comSetDFCat.iloc[[i]])
483                             )
484
485                             # This step is to impute missing values for
486                             # categorical attributes
487                             imputeValues = IS_MCS_Measure(inObjectCat,
488                               comSetAttr, inSetAttr, listOfObjects)
489
490                             indexList = inObjectCat.index.T.values
491                             i = 0;
492                             for j in range (0, len(indexList)):
493                                 if inObjectCat[indexList[j]] == '?':
494                                     inObjectCat[indexList[j]] =
495                                         imputeValues[i]
496                                         i +=1
497
498                             # This step is to fill in missing values for
499                             # numeric attributes
500                             inObjectNum = numericImputation(inSetDFNum,
501                               comSetDFNum, IDList)
502
503                             mixedObject = concatenationTwoObjects(inObjectNum,
504                               inObjectCat)
505                             #This step is to add imputed object into complete
506                             #dataset and remove it to complete datasets
507                             comSetDF = pd.concat([comSetDF, mixedObject.
508                               to_frame().T], ignore_index=True)
509                             comSet = np.asarray(comSetDF)
510                             inSetDFNum = inSetDFNum.drop(inSetDFNum.index[0])

```

```

496     inSetDFNum = inSetDFNum.reset_index(drop=True)
497     inSetDFCat = inSetDFCat.drop(inSetDFCat.index[0])
498     inSetDFCat = inSetDFCat.reset_index(drop=True)
499     inSetDF = inSetDF.drop(inSetDF.index[0])
500     inSetDF = inSetDF.reset_index(drop=True)
501     # Then, this step is to perform clustering process
502     npoints, nattrs = comSet.shape
503     global_attr_freq = _cal_global_attr_freq(comSet,
504     npoints, nattrs, categorical)
504     cl_attr_freq, membership = _init_clusters(comSet,
505     centroids, n_clusters, nattrs, npoints, verbose,
506     categorical)
507     centroids = [[defaultdict(float) for _ in range(
507     nattrs)] for _ in range(n_clusters)]
508     # Perform initial centroid update
509     lbd = np.zeros(n_clusters, dtype='float')
510     for ik in range(n_clusters):
511         cluster_members = sum(membership[ik, :])
512         attrIndex = 0
513         for iattr in range(nattrs):
514             if iattr in categorical:
515                 centroids[ik][iattr] =
516                     cal_centroid_value(lbd[ik], cl_attr_freq[ik][attrIndex],
516                     cluster_members,
517                     len(global_attr_freq[attrIndex]))
518                     attrIndex += 1
519             else:
520                 centroids[ik][iattr] = cal_mean_value(
520                 comSetDF, iattr)
521                 # print(comSet)
522                 # print("\n")
523                 # ''
524                 # Bước lặp chính của thuật toán
525                 # 1. Tính các vector trung tâm, lambda
526                 # 2. Nếu dissimilarity mới (cost) nhỏ hơn thì cập
527                 nhật và tiếp tục thực
528                 # hiện thuật toán lần nữa (từ bước 1).
529                 # Nếu lớn hơn thì kết thúc thuật toán.
530                 # ''
531                 # while itr <= max_iter and not converged:
532                 # while not converged:
533                 #     itr += 1
534                 #     if verbose:
535                 #         print("...k-center loop")
535                 centroids, moves, lbd = _k_CMM_iter(comSet,
535                 categorical, centroids, cl_attr_freq, membership,

```

```

535 global_attr_freq, lbd,
536     use_global_attr_count)
537         labels, ncost = _labels_cost(comSet, categorical,
538             centroids, global_attr_freq)
538     cost = ncost
539     # Store result of current run
540     all_centroids.append(centroids)
541     all_labels.append(labels)
542     all_costs.append(cost)
543     # while itr <= max_iter and not converged:
544     converged = (moves == 0) or (ncost >= cost)
545     if not converged:
546         while not converged:
547             all_costs = []
548
549             if verbose:
550                 print("... k-center loop")
551                 centroids, moves, lbd = _k_CMM_iter(comSet,
552                     centroids, cl_attr_freq, membership, global_attr_freq, lbd
552
552             use_global_attr_count)
553             if verbose:
554                 print("...Update labels, costs")
555                 labels, ncost = _labels_cost(comSet,
556                     categorical, centroids, global_attr_freq)
556                 converged = (moves == 0) or (ncost >= cost)
557                 cost = ncost
558                 # if verbose:
559                 #     print("Run {}, iteration: {} / {}, moves
559 : {}, cost: {}"
560                 #         .format(init_no + 1, itr, max_iter
560 , moves, cost))
561                 # Store result of current run
562                 all_centroids.append(centroids)
563                 all_labels.append(labels)
564                 all_costs.append(cost)
565             ...
566             Return the labels that contains exactly the number of
566             data points in the original dataset X
567             ...
568             # print(all_labels[-1])
569             return all_centroids[-1], all_labels[-1], all_costs[-1]
570
571
572 class KCMM(object):
573
574     '''k-CMM clustering algorithm for mixed numeric and

```

```

574 categorical data with missing values
575
576     Parameters
577     -----
578         K : int, optional, default: 8
579             The number of clusters to form as well as the
580             number of
581                 centroids to generate.
582
583         categorical: ndarray
584             Indicate the list of categorical
585             attributes in a dataset
586
587         max_iter : int, default: 300
588             Maximum number of iterations of the k-modes
589             algorithm for a
590                 single run.
591
592         n_init : int, default: 10
593             Number of time the k-modes algorithm will be run
594             with different
595                 centroid seeds. The final results will be the best
596             output of
597                 n_init consecutive runs in terms of cost.
598
599         init : 'random'
600             'random': choose k observations (rows) at random
601             from data for
602                 the initial centroids.
603
604         verbose : boolean, optional
605             Verbosity mode.
606
607     Attributes
608     -----
609         cluster_centroids_ : array, [K, n_features]
610             Categories of cluster centroids
611
612         labels_ :
613             Labels of each point
614
615         cost_ : float
616             Clustering cost, defined as the sum distance of
617             all points to
618                 their respective cluster centroids.
619
620         ...
621
622
623
624     def __init__(self, categorical, n_clusters, init='
625         random', n_init=3, max_iter = 10,
626             verbose=1, use_global_attr_count=1):

```

```

616     if verbose:
617         print("Number of clusters: {0}" . format(
618             n_clusters))
618         print("Init type: {0}" . format(init))
619         print("Max iterations: {0}" . format(max_iter)
619     )
620         print("Use global attributes count: {0}" .
620             format(use_global_attr_count > 0))
621
622     if hasattr(init, '__array__'):
623         n_clusters = init.shape[0]
624         init = np.asarray(init, dtype = np.float64)
625
626     self.categorical = categorical
627     self.n_clusters = n_clusters
628     self.init = init
629     self.n_init = n_init
630     self.verbose = verbose
631     self.max_iter = max_iter
632     self.use_global_attr_count = use_global_attr_count
633
634 def fit(self, X, **kwargs):
635     '''Compute k-CMM clustering.
636
637     Parameters
638     -----
639     X : array-like, shape=[n_samples, n_features]
640     ...
641
642     self.cluster_centroids_, self.labels_, self.cost_
642 = \
643         perform_kCMM(X, self.categorical, self.
643         n_clusters, self.init, self.verbose, self.
643         use_global_attr_count)
644     return self
645
646 def fit_predict(self, X, **kwargs):
647     '''Compute cluster centroids and predict cluster
647     index for each sample.
648
649     Convenience method; equivalent to calling fit(X)
649     followed by
650         predict(X).
651     ...
652     return self.fit(X, **kwargs).labels_
653
654 #This function is used to read a data with missing values
654 and then split it into two parts: completed dataset and
654 missing dataset
655 def readData(x):

```

```

656     inSet = []
657     comSet = []
658     for i in range(0, len(x)):
659         object = list(x[i])
660         for j in range(0, len(object)):
661             flag = False
662             if(object[j] == '?'):
663                 inSet.append(object)
664                 flag = True
665                 break;
666             if(flag == False):
667                 comSet.append(object)
668     inSetDF = pd.DataFrame(inSet)
669     comSetDF = pd.DataFrame(comSet)
670     return inSetDF, comSetDF
671
672 def getLeafNodes(decisionTree):
673     result = []
674     for key, value in decisionTree.results.items():
675         result.append(key)
676     return result
677
678 # This function is used to get all paths in a decision
679 # tree where they have the same categorical values in many
680 # column
681 def getAllPathsInDT(decisionTree):
682     result = []
683     firstTime = True
684     def recursiveCall(decisionTree, list, firstTime):
685         if decisionTree.results != None: # leaf node
686             leafNodeList = getLeafNodes(decisionTree)
687             for i in range(len(leafNodeList)):
688                 tmpList = list.copy()
689                 # tmpList.append(str(decisionTree.results
690                 # .split()[0].split("{}")[1])
691                 tmpList.append(leafNodeList[i])
692                 result.append(tmpList)
693             return
694         else:
695             # index column of each decision node is added
696             # to distinguish the same categories
697             # appearing in attributes.
698             decision = decisionTree.value
699             decisionTrue = str(decisionTree.col) + '-' +
str(decision)
700             decisionFalse = str(decisionTree.col) + '-!' +
str(decision)
701             # if(decisionTree.trueBranch.value == None):
702             if(firstTime == True):
703                 tmpTrue = [decisionTrue]

```

```

700                 tmpFalse = [decisionFalse]
701                 firstTime = False
702             else:
703                 tmpTrue = list.copy()
704                 tmpTrue.append(decisionTrue)
705                 tmpFalse = list.copy()
706                 tmpFalse.append(decisionFalse)
707                 recursiveCall(decisionTree.trueBranch, tmpTrue
708 , firstTime)
709                 recursiveCall(decisionTree.falseBranch,
710 tmpFalse, firstTime)
710             recursiveCall(decisionTree, None, firstTime)
711             # print(result)
712             return result
713
714     def splitNodeInDT(node):
715         if node.find('—') == -1:
716             return -1, node
717         else:
718             result = node.split('—')
719             return result[0], result[1]
720
721     def findSuitableLeave(decisionTree, inObject):
722         # print(object)
723         bestPath = []
724         listOfPaths = getAllPathsInDT(decisionTree)
725         for i in range(0, len(listOfPaths)):
726             # path = pd.DataFrame(listOfPaths[i])
727             path = listOfPaths[i]
728             # print(type(path))
729             # print(path)
730             for j in range(0, len(path)):
731                 nodeIndex,nodeValue = splitNodeInDT(path[j])
732                 if nodeIndex == -1:
733                     bestPath.append(path)
734                     break
735                 if(nodeValue[0] != '!!'):
736                     if(inObject[nodeIndex] !=nodeValue):
737                         break
738                    nodeValue =nodeValue[1:]
739                     if(inObject[nodeIndex]==nodeValue):
740                         break
741             return bestPath
742
743 #This function is to assign objects in the complete
dataset into leaves in DT
744     def assignCompleteObjectsIntoDT(path, comSet):
745         listOfObject = []
746         for i in range(0, len(comSet)):
```

```

747     comObject = comSet.loc[i]
748     for j in range(0, len(path)):
749         nodeIndex, nodeValue = splitNodeInDT(path[j])
750         if nodeIndex == -1:
751             listOfObject.append(i)
752             break
753         if (nodeValue[0] != '!'):
754             if (comObject[int(nodeIndex)] != nodeValue
755             ):
756                 break
757             else:
758                 nodeValue = nodeValue[1:]
759                 if (comObject[int(nodeIndex)] == nodeValue
760                 ):
761                     break
762     return listOfObject
763
764 # This function is used to calculate the IS measure that
765 # is the correlation between set of attributes with non-
766 # missing values C and set
767 # of attributes with missing values M within a record. The
768 # IS measure measures the degree of associations between
769 # two sets of attribute values.
770 # IS = Support(C,M)/sqrt(Support(C)*Support(Q)) where
771 # Support(C,M) = |C,M|/Q, |C,M| is the number of records
772 # that contain both C and M, Q is the size of dataset
773
774 def getFrequency(comSetAttr, inSetAttr, listOfObjects):
775     valueSet = {} # Use to store: key: complete values,
776     values: values at corresponding missing attributes in
777     complete object
778     comObjectSet = {} # Use to store: key: complete values
779     , values: indices of complete objects
780     # that will be used for calculating the MSC measure
781     freqC = {}
782     freqM = {}
783     freqCM = {}
784     for i in range(0, len(listOfObjects)):
785         comObject = listOfObjects[i]
786         # Note that the variable missingValues is the
787         # values at the missing attribute in the complete object
788         missingValues = comObject[inSetAttr]
789         missingValuesDF = missingValues.copy()
790         missingIndexList = missingValues.columns.values
791         for index in missingIndexList:
792             missingValuesDF.loc[:,index] = str(index)+'-'+index
793             missingValuesDF.loc[:,index]
794         missingValuesTuple = [tuple(x) for x in
795 
```

```

782 missingValuesDF.values][0]
783     #
784     completeValues = comObject[comSetAttr]
785     completeValuesDF = completeValues.copy()
786     completeIndexList = completeValues.columns.values
787     for index in completeIndexList:
788         completeValuesDF.loc[:,index] = str(index)+'-' + completeValues.loc[:,index]
789     completeValuesTuple = [tuple(x) for x in completeValuesDF.values][0]
790     valueSet[completeValuesTuple] = missingValuesTuple
791     if completeValuesTuple in comObjectSet:
792         objects = comObjectSet[completeValuesTuple]
793         if i not in objects:
794             objects.append(i)
795         comObjectSet[completeValuesTuple] = objects
796     else:
797         objects = [i]
798         comObjectSet[completeValuesTuple] = objects
799     CMTuple = completeValuesTuple + missingValuesTuple
800     if missingValuesTuple in freqM:
801         count = freqM[missingValuesTuple]
802         count += 1
803         freqM[missingValuesTuple] = count
804     else:
805         freqM[missingValuesTuple] = 1
806     if completeValuesTuple in freqC:
807         count = freqC[completeValuesTuple]
808         count += 1
809         freqC[completeValuesTuple] = count
810     else:
811         freqC[completeValuesTuple] = 1
812     if CMTuple in freqCM:
813         count = freqCM[CMTuple]
814         count += 1
815         freqCM[CMTuple] = count
816     else:
817         freqCM[CMTuple] = 1
818
819     return comObjectSet, freqM, freqC, freqCM, valueSet
820
821 # This function is used to calculate the frequency of
822 # categorical values in listOfObjects
823 def categoricalValuesFrequency(inObject, comSetAttr,
824     listOfObjects):
825     inObjectCat = pd.DataFrame(inObject).T
826     listOfObjects.append(inObjectCat)
827     result = {}
828     # print(listOfObjects[0])
829     # indexList = inObject.index.T.values

```

```

828     for i in range(0, len(listOfObjects)):
829         comObject = listOfObjects[i]
830         for j in range(0, len(comSetAttr)):
831             categoryValue = str(comSetAttr[j]) + '-' +
comObject[comSetAttr[j]]
832             categoryTuple = tuple(categoryValue)
833             if categoryTuple in result:
834                 count = result[categoryTuple]
835                 count = count +1
836                 result[categoryTuple] = count
837             else:
838                 result[categoryTuple] = 1
839     return result, len(listOfObjects)
840
841 def IS_MCS_Measure(inObject, comSetAttr, inSetAttr,
listOfObjects):
842     resultIS = {}
843     q = len(listOfObjects)
844     comObjectSet, freqM, freqC, freqCM, valueSet =
getFrequency(comSetAttr, inSetAttr, listOfObjects)
845     for key, value in freqC.items():
846         supC = value/q
847         supM = freqM[valueSet[key]]/q
848         _key = list(key)
849         CM = _key.copy()
850         for j in range(0, len(valueSet[key])):
851             CM.append(valueSet[key][j])
852         CMTuple = tuple(CM)
853         supCM = freqCM[CMTuple]/q
854         IS = supCM/math.sqrt(supC*supM)
855         resultIS[key]= IS
856     # The following code is used to calculate the MCS
measure that quantifies the similarity between an object
with missing values
857     # and an object with no missing values
858     frequency, cardinality = categoricalValuesFrequency(
inObject, comSetAttr, listOfObjects)
859     del listOfObjects[-1]
860     resultMCS = {}
861     for key, value in comObjectSet.items():
862         msc = 0
863         for j in range(0, len(comSetAttr)):
864             # Note that the variable missValue below is
the category incomplete attribute of incomplete object
865             # print(inObject[comSetAttr[j]])
866             missValue = str(comSetAttr[j]) + '-' +
inObject[comSetAttr[j]]
867             missValueTuple = tuple([missValue])
868             compValueTuple = tuple([key[j]])
869             if missValueTuple == compValueTuple:

```

```

870             msc += 1
871         else:
872             freqMissValue = frequency[missValueTuple]
873             / cardinality
874             freqComValue = frequency[compValueTuple] /
875             cardinality
876             sumFreq = (frequency[missValueTuple] +
877             frequency[compValueTuple]) / cardinality
878             msc += (2 * math.log(sumFreq)) / (math.log
879             (freqComValue) + math.log(freqMissValue))
880             resultMCS[key] = msc
881             sumIS_MCS = []
882             # Calculate the affinity degree for the IS and MCS
883             measures
884             for key, value in resultMCS.items():
885                 sumIS_MCS.append((resultIS[key]+value)/2.0)
886             sumValue = sum(sumIS_MCS)
887             randomSampling = []
888             for i in range(0,len(sumIS_MCS)):
889                 randomSampling.append(sumIS_MCS[i]/sumValue)
890             maxIndexValue = randomSampling.index(max(
891             randomSampling))
892             tmp = list(resultIS)[maxIndexValue]
893             imputeValues = valueSet[tmp]
894             result = []
895             for value in imputeValues:
896                 result.append(value.split('-')[1])
897             return result
898
899 def do_kr(x, y, nclusters, verbose, use_global_attr_count,
900 n_init):
901     start_time = time()
902     tracemalloc.start()
903     categorical = [0, 3, 4, 5, 6, 8, 9, 11, 12]
904     kr = KCMM(categorical, n_clusters = nclusters, init='
905     random',
906             n_init = n_init, verbose = verbose,
907             use_global_attr_count = use_global_attr_count)
908             kr.fit_predict(x)
909             # print(kr.labels_)
910
911             ari = evaluation.rand(kr.labels_, y)
912             nmi = evaluation.nmi(kr.labels_, y)
913             purity = evaluation.purity(kr.labels_, y)
914             end_time = time()
915             elapsedTime = timedelta(seconds=end_time - start_time)
916             .total_seconds()
917             memoryUsage = tracemalloc.get_tracemalloc_memory() /
918             1024 / 1024
919             if verbose == 1:

```

```

909         print("Purity = {:.3f}".format(purity))
910         print("NMI = {:.3f}".format(nmi))
911         print("Adjusted rand index = {:.3f}" . format(ari)
912     ))
913         print("Elapsed Time = {:.3f} secs".format(
914     elapsedTime))
915         print("Memory usage = {:.3f} MB".format(
916     memoryUsage))
917         print("Finish all task!")
918     tracemalloc.stop()
919     return [round(purity, 3), round(nmi, 3), round(ari, 3)
920 , round(elapsedTime, 3), round(memoryUsage, 3)]
921
922 def run(argv):
923     max_iter = 10
924     ifile = "data/mixed_credit.csv"
925     ofile = "output/credit.csv"
926     use_global_attr_count = 0
927     use_first_column_as_label = False
928     verbose = 1
929     delim = ","
930     n_init = 3
931
932     # Get samples & labels
933     if not use_first_column_as_label:
934         x = np.genfromtxt(ifile, dtype = str, delimiter =
935     delim)[:, :-1]
936         y = np.genfromtxt(ifile, dtype = str, delimiter =
937     delim, usecols = -1)
938     else:
939         x = np.genfromtxt(ifile, dtype = str, delimiter =
940     delim)[:, 1:]
941         y = np.genfromtxt(ifile, dtype = str, delimiter =
942     delim, usecols = 0)
943
944     from collections import Counter
945     nclusters = len(list(Counter(y)))
946     result = []
947     for i in range(max_iter):
948         if verbose:
949             print("\n=====Run {0}/{1} times
950 =====\n" . format(i + 1, max_iter))
951         result.append(do_kr(x, y, nclusters, verbose =
952     verbose, use_global_attr_count = use_global_attr_count,
953     n_init = n_init))
954     resultDF = pd.DataFrame(result)
955     tmpResult = []
956     for i in range(0,5):
957         tmpResult.append(cal_mean_value(resultDF,i))
958     finalResult = [['Purity', 'NMI', 'Adjusted rand index']

```

```
947 , "Elapsed Time", "Memory Usage"]]  
948     finalResult.append(tmpResult)  
949     import csv  
950     with open(ofile, 'w') as fp:  
951         writer = csv.writer(fp, delimiter = ',')  
952         writer.writerows(finalResult)  
953  
954  
955 if __name__ == "__main__":  
956     run(sys.argv[1:])
```