

# Grant Cloud

## 903194929

### CS-4641 HW3

#### Imports

```
In [1]: 1 import numpy as np
        2 import matplotlib.pyplot, os.path
        3 from mpl_toolkits.mplot3d import Axes3D
```

Provided code, with my modified functions

In [2]:

```

1  ## Data Loading utility functions
2  def get_test_train(fname,seed,datatype):
3      '''
4      Returns a test/train split of the data in fname shuffled with
5      the given seed
6
7
8      Args:
9          fname:      A str/file object that points to the CSV file to load, p
10                     np.genfromtxt()
11          seed:       The seed passed to np.random.seed(). Typically an int or
12                     datatype:  The datatype to pass to genfromtxt(), usually int, float
13
14
15      Returns:
16          train_X:     A NxD np array of training data (row-vectors), 80% of all
17          train_Y:     A Nx1 np array of class labels for the training data
18          test_X:      A MxD np array of testing data, same format as train_X,
19          test_Y:      A Mx1 np array of class labels for the testing data
20      '''
21      data = np.genfromtxt(fname,delimiter=',',dtype=datatype)
22      np.random.seed(seed)
23      shuffled_idx = np.random.permutation(data.shape[0])
24      cutoff = int(data.shape[0]*0.8)
25      train_data = data[shuffled_idx[:cutoff]]
26      test_data = data[shuffled_idx[cutoff:]]
27      train_X = train_data[:, :-1].astype(float)
28      train_Y = train_data[:, -1].reshape(-1,1)
29      test_X = test_data[:, :-1].astype(float)
30      test_Y = test_data[:, -1].reshape(-1,1)
31      return train_X, train_Y, test_X, test_Y
32
33
34  def load_HTRU2(path='data'):
35      return get_test_train(os.path.join(path, 'HTRU_2.csv'), seed=1567708903, da
36
37  def load_iris(path='data'):
38      return get_test_train(os.path.join(path, 'iris.data'), seed=1567708904, da
39
40  ## The "digits" dataset has a pre-split set of data, so we won't do our own
41  def load_digits(path='data'):
42      train_data = np.genfromtxt(os.path.join(path, 'optdigits.tra'), delimiter=
43      test_data = np.genfromtxt(os.path.join(path, 'optdigits.tes'), delimiter=
44      return train_data[:, :-1], train_data[:, -1].reshape(-1,1), test_data[:, :-
45
46  ## You can use this dataset to debug your implementation
47  def load_test2(path='data'):
48      return get_test_train(os.path.join(path, 'data2.dat'), seed=1568572211, da
49
50  class PCA():
51      '''
52      A popular feature transformation/reduction/visualization method
53
54
55      Uses the singular value decomposition to find the orthogonal directions
56      maximum variance.

```

```

57     '''
58     def __init__(self):
59         '''
60         Initializes some data members to hold the three components of the
61         SVD.
62         '''
63         self.u = None
64         self.s = None
65         self.v = None
66         self.shift = None
67         self.data = None
68
69     def find_components(self,data):
70         '''
71         Finds the SVD factorization and stores the result.
72
73         Args:
74             data: A NxD array of data points in row-vector format.
75         '''
76         self.data = data
77         self.shift = self.data - np.mean(data, axis=0) # normalized data
78         self.u, self.s, self.v = np.linalg.svd(self.shift, compute_uv=True,
79         self.v = self.v.T
80
81     def transform(self,n_components,data=None):
82         '''
83         Uses the values computed and stored after calling find_components()
84         to transform the data into n_components dimensions.
85
86         Args:
87             n_components: The number of dimensions to transform the data into
88             data: the data to apply the transform to. Defaults to the data
89                   provided on the last call to find_components()
90
91         Returns:
92             transformed_data: a Nx(n_components) array of transformed points
93                               in row-vector format.
94         '''
95         if data is None:
96             data = self.data
97         else:
98             self.find_components(data)
99         return np.dot(data, self.v[:, :n_components])
100
101     def inv_transform(self,n_components,transformed_data):
102         '''
103         Inverts the results of transform() (if given the same arguments).
104
105         Args:
106             n_components: Number of components to use. Should match
107                           the dimension of transformed_data.
108             transformed_data: The data to apply the inverse transform to,

```

```

114         should be in row-vector format
115
116
117     Returns:
118         inv_tform_data:    a NxD array of points in row-vector format
119     '''
120     #NOTE: Don't forget to "un-center" the data
121     tX = np.dot(transformed_data, self.v[:, :n_components].T)
122     return tX + np.mean(self.data, axis=0)
123
124     def reconstruct(self, n_components, data=None):
125         '''
126         Casts the data down to n_components dimensions, and then reverses the transform,
127         returning the low-rank approximation of the given data. Defaults to the transform
128         provided on the last call to find_components().
129         '''
130         return self.inv_transform(n_components, self.transform(n_components, data))
131
132     def reconstruction_error(self, n_components, data=None):
133         '''
134         Useful for determining how much information is preserved in n_components.
135         '''
136         if data is None:
137             data = self.data
138         return np.linalg.norm(data - self.reconstruct(n_components, data), ord='fro')
139
140     def plot_2D_proj(self, data=None, labels=None):
141         '''
142         Creates a 2D visualization of the data, returning the created figure. See
143         the main() function for example usage.
144         '''
145         fig = matplotlib.pyplot.figure()
146         proj_2d_data = self.transform(2, data)
147         fig.gca().scatter(proj_2d_data[:, 0], proj_2d_data[:, 1], c=labels)
148         fig.gca().set_title('PCA 2D transformation')
149         return fig
150     def plot_3D_proj(self, data=None, labels=None):
151         '''
152         Creates a 3D visualization of the data, returning the created figure. See
153         the main() function for example usage.
154         '''
155         fig = matplotlib.pyplot.figure()
156         ax = fig.add_subplot(111, projection='3d')
157         proj_3d_data = self.transform(3, data)
158         ax.scatter(proj_3d_data[:, 0], proj_3d_data[:, 1], proj_3d_data[:, 2], c=labels)
159         fig.gca().set_title('PCA 3D transformation')
160         return fig
161
162
163
164     # You may find this method useful for providing a canonical ordering of clusters
165     def consistent_ordering(array):
166         rowvec_dists = np.linalg.norm(array, axis=1)
167         dist_order = np.argsort(rowvec_dists, kind='stable')
168         return array[dist_order, :]
169     class KMeans():
170         '''

```

```

171     A simple iterative clustering method for real-valued feature spaces.
172
173
174     Finds cluster centers by iteratively assigning points to clusters and re-
175     cluster center locations. Provided code expects self.clusters to contain
176     cluster centers in row-vector format.
177     '''
178     def __init__(self):
179         '''
180         Initializes data members.
181         '''
182         self.clusters = None
183         self.print_every = 1000
184
185     def cluster_distances(self,data):
186         '''
187         Computes the distance from each row of data to the cluster centers.
188         otherwise set self.clusters first.
189
190
191         Args:
192             data:    The data to compute distances for, in row-vector format.
193                     same number of columns as each cluster
194
195
196         Returns:
197             dists:   A Nx(len(clusters)) array, one row for each row in data,
198                     cluster center, containing the distance for each point to
199             ...
200             return np.hstack([np.linalg.norm(data-c,axis=1).reshape(-1,1) for c
201
202     def cluster_label(self,data):
203         '''
204         Returns the label of the closest cluster to each row in data. Note that
205         arbitrary, and do *not* correspond directly with class labels.
206
207
208         Args:
209             data:    Data to compute cluster labels for, in row-vector format
210
211
212         Returns:
213             c_labels: A N-by-1 array, one row for each row in data, containing
214                     corresponding to the cluster closest to each point.
215             ...
216             return np.argmin(self.cluster_distances(data),axis=1)
217
218     def cluster(self,data,k):
219         '''
220         Implements the k-Means iterative algorithm. Cluster centers are initialized
221         then on each iteration each data point is assigned to the closest cluster
222         cluster centers are re-computed by averaging the points assigned to
223         self.clusters should contain a k-by-D array of the cluster centers
224
225
226         Args:
227             data:    Data to be clustered in row-vector format. A N-by-D array

```

```

228         k:         The number of clusters to find.
229         ...
230     #This line should pick the initial clusters at random from the provi
231     #data.
232     self.clusters = data[np.random.choice(data.shape[0],k,replace=False)]
233     #An example of how to use the consistent_ordering() function, which
234     #may want to use to help determine if cluster centers have changed j
235     #iteration to the next
236     self.clusters = consistent_ordering(self.clusters)
237     #We know that k-Means will always converge, but depending on the ini
238     #conditions and dataset, it may take a long time. For debugging purp
239     #you might want to set a maximum number of iterations. When impleme
240     #correctly, none of the provided datasets take many iterations to co
241     #for most initial configurations.
242     not_done = True
243     itr = 0
244     while not_done:
245         itr += 1
246         new_clusters = np.zeros(self.clusters.shape)
247         #assign points to nearest cluster
248         labels = self.cluster_label(data)
249         #re-compute cluster centers by averaging the points assigned to
250         for i in range(new_clusters.shape[0]):
251             new_clusters[i] = data[np.argwhere(labels==i)].mean(axis=0)
252         #determine if clusters have changed from the previous iteration
253         new_clusters = consistent_ordering(new_clusters)
254         if np.array_equal(self.clusters, new_clusters):
255             not_done = False
256         #For debugging, print out every so often.
257         if itr == 0:
258             print('')
259         if itr % self.print_every == 0:
260             print("Iteration {}, change {}".format(itr,np.linalg.norm(ne
261             self.clusters = new_clusters
262         print("Converged after {} iterations".format(itr))
263
264     def normalized_mutual_information(self, data, labels):
265         '''
266         Since cluster assignments are not the same as class labels, we can't
267         compare them to measure clustering performance. However, we can meas
268         between two labelings, to see if they contain the same statistical i
269         implements the "Normalized Mutual Information Score" as described he
270
271         https://scikit-learn.org/stable/modules/clustering.html#mutual-infor
272
273         Note that this version uses arithmetic mean, when comparing output v
274         sklearn.metrics.mutual_info_score()
275         '''
276         cluster_labels = self.cluster_label(data)
277         P_cl = np.zeros(len(self.clusters))
278         P_gt = np.zeros(len(np.unique(labels)))
279         P_clgt = np.zeros((len(P_cl),len(P_gt)))
280         cl_masks = dict()
281         gt_masks = dict()
282         MI = 0.0
283         H_cl = 0.0
284         H_gt = 0.0

```

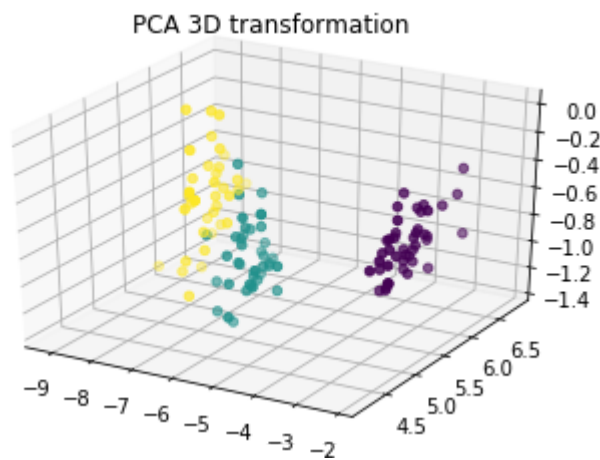
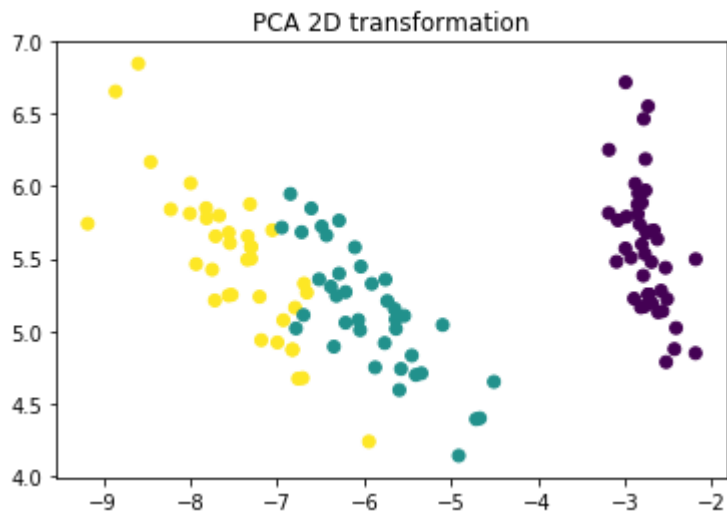
```

285     for c_cl in range(len(P_cl)):
286         cl_masks[c_cl] = (cluster_labels==c_cl).reshape(-1,1)
287         P_cl[c_cl] = (cl_masks[c_cl]).astype(int).sum()/len(data)
288         H_cl -= P_cl[c_cl]*np.log(P_cl[c_cl])
289     for c_gt in range(len(P_gt)):
290         gt_masks[c_gt] = labels == c_gt
291         P_gt[c_gt] = (gt_masks[c_gt]).astype(int).sum()/len(data)
292         H_gt -= P_gt[c_gt]*np.log(P_gt[c_gt])
293     for c_cl in range(len(P_cl)):
294         for c_gt in range(len(P_gt)):
295             P_clgt[c_cl,c_gt] = (np.logical_and(cl_masks[c_cl], gt_masks[c_gt])).sum()/len(data)
296             if P_clgt[c_cl,c_gt] == 0.0:
297                 MI += 0
298             else:
299                 MI += P_clgt[c_cl,c_gt]*np.log(P_clgt[c_cl,c_gt]/(P_cl[c_cl]*P_gt[c_gt]))
300     return MI/(np.mean([H_cl,H_gt]))
301
302 def plot_2D_clusterd(self, data, labels=None):
303     """
304     Creates a 2D visualization of the data, returning the created figure
305     the main() function for example usage.
306     """
307     if self.clusters is None or len(self.clusters)!=2:
308         self.cluster(data,2)
309     fig = matplotlib.pyplot.figure()
310     clusterd_2d_data = self.cluster_distances(data)
311     fig.gca().scatter(clusterd_2d_data[:,0],clusterd_2d_data[:,1],c=labels)
312     fig.gca().set_title('k-Means 2D cluster distance')
313     return fig
314 def plot_3D_clusterd(self, data, labels=None):
315     """
316     Creates a 3D visualization of the data, returning the created figure
317     the main() function for example usage.
318     """
319     if self.clusters is None or len(self.clusters)!=3:
320         self.cluster(data,3)
321     fig = matplotlib.pyplot.figure()
322     ax = fig.add_subplot(111,projection='3d')
323     clusterd_3d_data = self.cluster_distances(data)
324     ax.scatter(clusterd_3d_data[:,0],clusterd_3d_data[:,1],clusterd_3d_data[:,2],c=labels)
325     fig.gca().set_title('k-Means 3D cluster distance')
326     return fig

```

**1a) and 1b)**

```
In [3]: 1 pca = PCA()  
2 data = load_iris()  
3 pca.find_components(data[0])  
4 _, labels = np.unique(data[1], return_inverse=True)  
5 pca.plot_2D_proj(data[0], labels)  
6 pca.plot_3D_proj(data[0], labels)  
7 print('') # prevents Jupyter from repeat printing a figure
```





```
In [4]: 1 reconError = pca.reconstruction_error(2,data[0])
        2 print('2d reconstruction error: ',reconError)
        3 reconError = pca.reconstruction_error(3,data[0])
        4 print('3d reconstruction error: ',reconError)
```

2d reconstruction error: 82.49015254320533

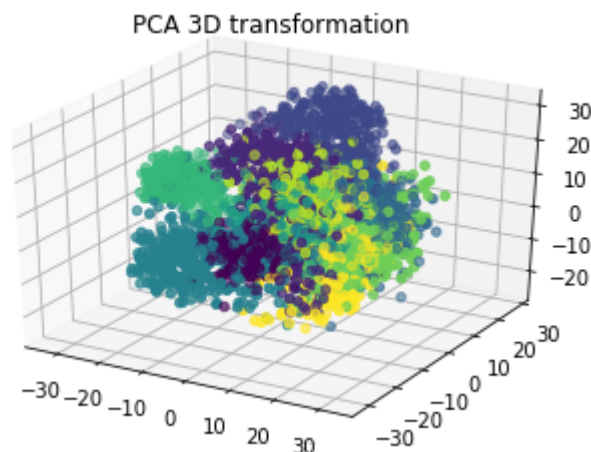
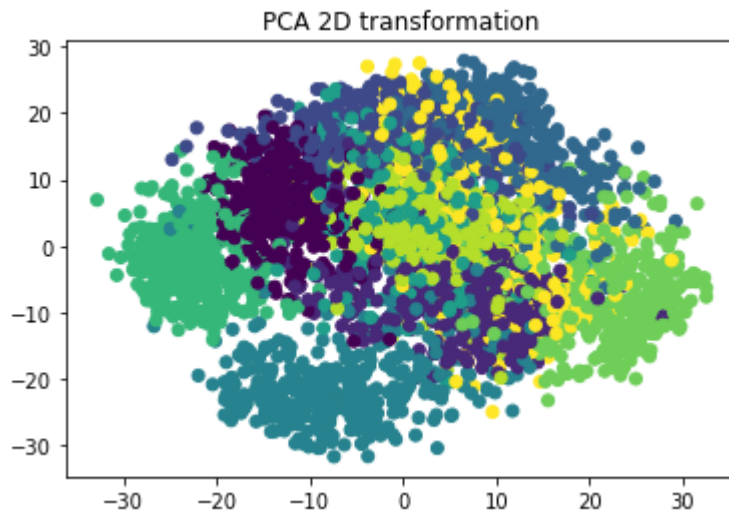
3d reconstruction error: 82.82381890639289

With the Iris dataset, we expect classification to perform better on the 2D transformed features than on the original dataset. From the 2D transformation plot we can see the purple class is easily distinguished from the other 2 classes, and the remaining 2 classes have very little overlap. This is because PCA projects onto new axis' built from the eigenvalues of the highest variance features in the original dataset (this process is called Truncated SVD), this leads to reduced noise in the transformed data. In this case a good learner would definitely improve performance since we lost little (if any) signal but removed a lot of noise.

For the 3D case, we expect classification to perform about the same or maybe slightly better on the PCA transformed data as with the original features. This is because the Iris dataset only has 4 features, meaning the 3 transformed features contain almost all the same information as the original dataset. The amount of improvement depends on how much noise was contained in the 4th feature not chosen by PCA, as well as how much of that noise was removed when the transformation occurred. From the plot the purple points are still easily classified (same as with the original features) but there looks to be more overlap in the yellow and blue classes than in the case of the 2D transformation. Overall I would expect a classifier to perform better on the 2D transformed data than the 3D transformed data (since the 3D transformed data has additional noise).

For both cases we have about the same reconstruction error, showing the added feature may hold little to no information (and if any it is probably noise).

```
In [5]: 1 data = load_digits()
2 pca.find_components(data[0])
3 _, labels = np.unique(data[1], return_inverse=True)
4 pca.plot_2D_proj(data[0], labels)
5 pca.plot_3D_proj(data[0], labels)
6 print('')
```



```
In [6]: 1 reconError = pca.reconstruction_error(2,data[0])
2 print('2d reconstruction error: ',reconError)
3 reconError = pca.reconstruction_error(3,data[0])
4 print('3d reconstruction error: ',reconError)
```

```
2d reconstruction error: 1820.7713792883224
3d reconstruction error: 1673.459624274109
```

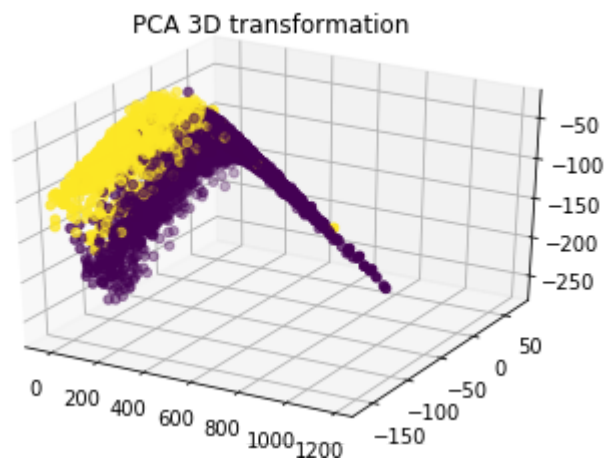
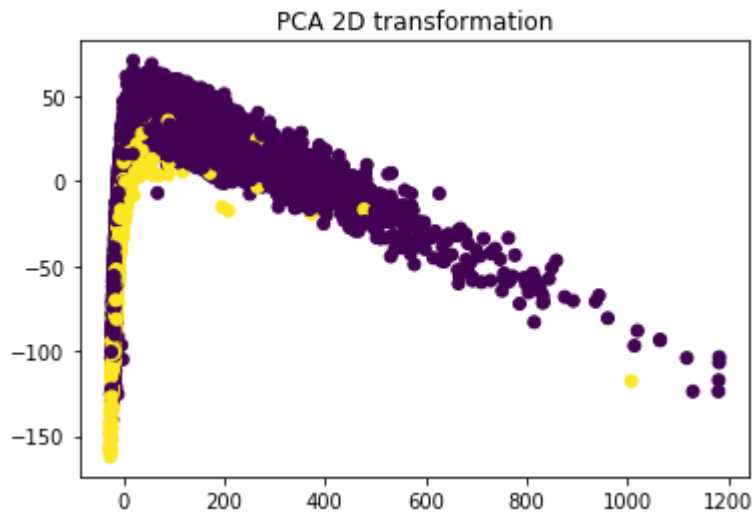
The 2D PCA transformation on the digits dataset would likely decrease classification performance. This is because the digits dataset has so many different class labels. When PCA transforms the data into the lower dimension, we lose information on a lot of features. With so many different class labels, there is bound to be different classes that have similar values for the 2 features PCA chose and the distinguishing factor between these classes might be contained in a feature that

was "squished" during the transformation. This loss of signal would decrease classification performance compared to the original feature set that still contains the important information to distinguish between these classes.

For the 3D case, we expect a classifier to perform better on the 3D transformed data than the original dataset. This is because the PCA transformation removes a lot of noise but maintains the information needed for a classifier to perform well (it looks like the additional feature selected in the 3D case contained a lot of signal). It may be the case that the best PCA transformation for the digits dataset would be one in 4D or even higher, since there are so many features and class labels, and there may be important signal in a feature that we lose when we scale down into 3D. We would also expect a classifier to perform better on the 3D PCA transformed data than the 2D PCA transformed data.

Performance on the digits dataset may also depend on what type of classifier is used. In the previous homework we saw how the high number of classes caused the logistic regression "one-vs-all" approach struggled with a high number of class labels. This would lead to lower performance on this dataset than kNN on the 3D PCA transformed data, which appears to be, in general, well clustered by class label in the plot.

```
In [7]: 1 data = load_HTRU2()
2 pca.find_components(data[0])
3 _, labels = np.unique(data[1], return_inverse=True)
4 pca.plot_2D_proj(data[0], labels)
5 pca.plot_3D_proj(data[0], labels)
6 print('')
```



```
In [8]: 1 reconError = pca.reconstruction_error(2,data[0])
2 print('2d reconstruction error: ',reconError)
3 reconError = pca.reconstruction_error(3,data[0])
4 print('3d reconstruction error: ',reconError)
```

2d reconstruction error: 13094.399418771973

3d reconstruction error: 18805.737261371978

The 2D PCA transformation on the HTRU2 dataset would decrease classifier performance. From the 2D plot we see a lot of overlap between the two classes. This occurs because PCA chooses the two highest variance features (with the underlying assumption that they provide the most information about the data), but this can be flawed when there is a high variance, noisy feature. This feature might be highly variant but actually provides little to no information about the underlying classes of the data points. This is also the most likely reason that we see both classes

occurring at both ends of the plot and a lot of overlap between classes. A classifier would struggle to classify the PCA transformed data using noisy features that contain little to no actual information about the underlying classes.

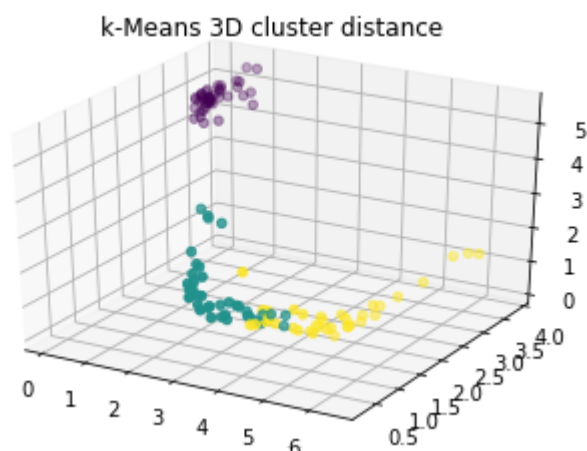
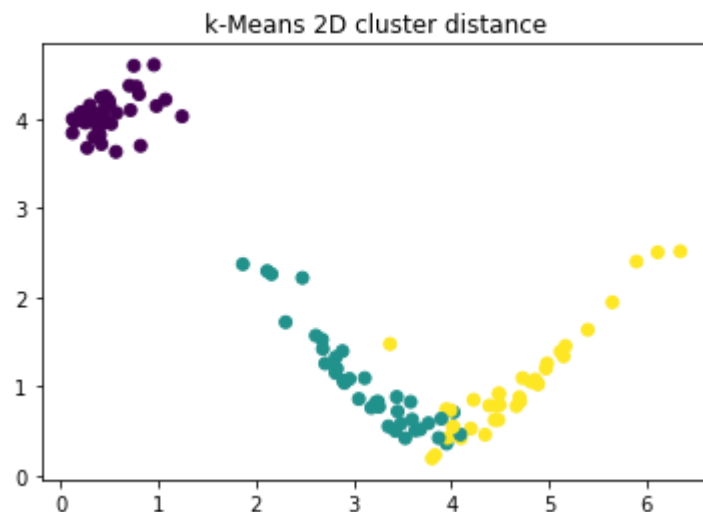
In the 3D PCA transformed data we see a much more distinguishable separation of classes. The additional feature selected by PCA clearly contains important information about the underlying class of each data point. Furthermore the transformation squished other less variant features that may contain a lot of noise. It's interesting to note that we have a higher reconstruction error and higher classifier performance with the 3D transformation. This shows that the additional "information" lost when we project onto the additional axis is mostly noise and isn't all that useful for a classifier in practice.

## **2a) and 2b)**

```
In [9]: 1 km = KMeans()
2 data = load_iris()
3 _, labels = np.unique(data[1], return_inverse=True)
4 km.plot_2D_clusterd(data[0], labels)
5 km.plot_3D_clusterd(data[0], labels)
6 print('')
```

Converged after 6 iterations

Converged after 5 iterations



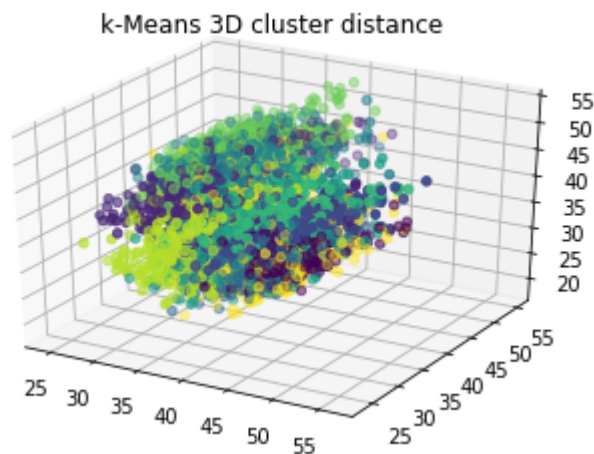
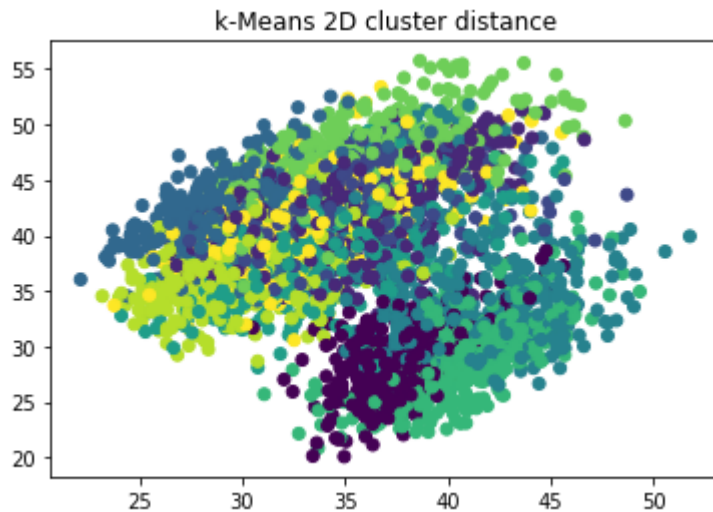
kMeans would be a better transformation than PCA in the 2D case because the clusters look to be perpendicular, whereas in the PCA transform they have similar shape and size. It's easy to distinguish them when the class labels are present but when those labels are not present it would be much harder to distinguish the different classes in the 2D PCA transformation. In general I would expect a classifier to perform better when there are less points with different class labels next to each other in the 2D visualization and minimal overlap between classes.

In the 3D case I would expect a classifier to perform about equally on the PCA transformed data as the kMeans transformed data. Both transformation have a pretty clear separation of class labels with some small overlap between two of the classes which would potentially lead to some loss in a classifier.

```
In [11]: 1 data = load_digits()
2 _, labels = np.unique(data[1], return_inverse=True)
3 km.plot_2D_clusterd(data[0], labels)
4 km.plot_3D_clusterd(data[0], labels)
5 print('')
```

Converged after 10 iterations

Converged after 25 iterations

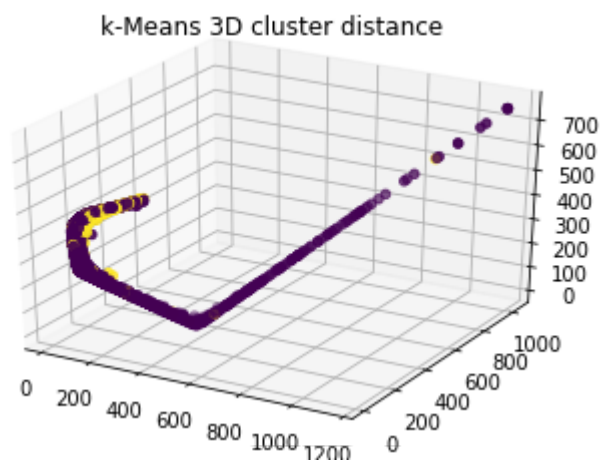
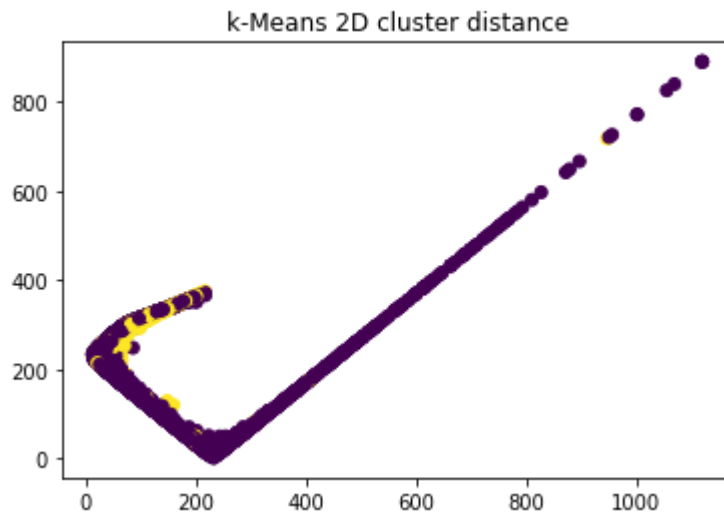


For both the 2D and 3D case in the digits dataset, the PCA transformation should be preferred to the kMeans transformation in terms of expected classifier performance. In both 2D and 3D kMeans transformations we see a large amount of overlap between different class labels and the clusters of classes are loose and wide. While the 2D PCA transformation probably wouldn't help classifier performance because of so much overlap between classes, the clustering of class labels was still tighter, and the 3D PCA transformation had clear clusters by different class labels with much less overlap. Overall a classifier would really struggle with classifying the kMeans transformed data.

```
In [12]: 1 data = load_HTRU2()
2 _, labels = np.unique(data[1], return_inverse=True)
3 print('2D : ', end='')
4 km.plot_2D_clusterd(data[0], labels)
5 print('3D : ', end='')
6 km.plot_3D_clusterd(data[0], labels)
7 print('')
```

2D : Converged after 23 iterations

3D : Converged after 44 iterations



In both the 2D and 3D cases PCA is absolutely the preferred transformation in terms of best classifier performance. The kMeans transformed data makes it virtually impossible to distinguish between class labels. There is so much overlap between different classes, a classifier would struggle to distinguish different classes. The 2D PCA transformation, while still not super helpful for



classifier performance, still has slightly easier to distinguish classes than the 2D kMeans transformed data. The 3D PCA transformation has a clear separability between classes, that separability is non-existent in the 3D kMeans transformed data plot.

## 2c)

```
In [13]: 1 km = KMeans()
2 data = load_iris()
3 _, labels = np.unique(data[1], return_inverse=True)
4 numClasses = len(np.unique(labels))
5 print('# of Classes: ', numClasses)
6 km.cluster(data[0], numClasses)
7
8 # changing dtype to float for mutual information scoring
9 i = 0
10 for classLabel in list(np.unique(data[1])):
11     data[1][data[1]==classLabel] = i
12     i += 1
13
14 print('Mutual Information Score: ', km.normalized_mutual_information(data[0],
```

```
# of Classes: 3
Converged after 14 iterations
Mutual Information Score: 0.5909010603830261
```

kMeans did an ok job of clustering by class label on the iris dataset. This is because there is overlap between different class labels in the Iris dataset. When kMeans assigns points to clusters the assigned cluster is the one with the smallest distance from the point. This causes similar points with different class labels to be a part of the same cluster, decreasing the mutual information score. In the case of the Iris dataset, the 0.591 mutual information score means there **is** a statistical relationship between class label and cluster assignments (both labelings contain similar information), but it is not a perfect indicator.

```
In [16]: 1 km = KMeans()
2 data = load_digits()
3 _, labels = np.unique(data[1], return_inverse=True)
4 numClasses = len(np.unique(labels))
5 print('# of Classes: ', numClasses)
6 km.cluster(data[0], numClasses)
7 print('Mutual Information Score: ', km.normalized_mutual_information(data[0],
```

```
# of Classes: 10
Converged after 48 iterations
Mutual Information Score: 0.7543227044445462
```

kMeans did the best job clustering according to class labels on the digits dataset, with the highest mutual information score of 0.745. This tells us that the clusters and labelings both contain similar information (there is a strong statistical relationship). We also notice that it took the most iterations to converge of all the datasets, this is because kMeans converges when cluster assignments no longer change from one iteration to the next. In the case of digits, there are 10 class labels. The high number of clusters means that more points are likely to be reassigned from one iteration to

the next (since there are more potential assignments). Digits has such a high mutual information score because data points that have the same class label have similar distances to the different cluster centers (mean of all the points assigned to a cluster).

```
In [15]: 1 km = KMeans()  
2 data = load_HTRU2()  
3 _, labels = np.unique(data[1], return_inverse=True)  
4 numClasses = len(np.unique(labels))  
5 print('# of Classes: ', numClasses)  
6 km.cluster(data[0], numClasses)  
7 print('Mutual Information Score: ', km.normalized_mutual_information(data[0],
```

# of Classes: 2

Converged after 23 iterations

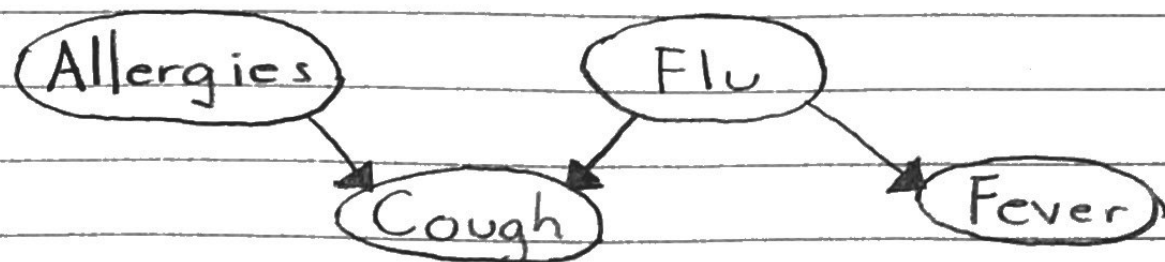
Mutual Information Score: 0.024972522126462778

kMeans did the worst job clustering according to class labels on the HTRU2 dataset, with a mutual information score of 0.025. This tells us that there is almost no statistical relationship between clusters and class labels. This also means the clusters are not representative of class labels. It's worth noting that kMeans can get stuck in local optima, it's possible that the converged cluster centers are not global optima and we would need to implement random restarts to help address this issue. The poor mutual information score most likely is caused by noisy features in the dataset, these features would have a large impact on cluster assignments (because of high distances) and no correlation to class labels.

```
In [ ]:
```

```
1
```

3a)



$$\begin{aligned}
 b) \quad & P(Cgh, Fev, Alg, Flu) \quad P(x_1, \dots, x_D) = \prod_{i=1}^D P(x_i | x_{1:i-1}) \\
 & = P(Cgh | Fev, Alg, Flu) \cdot P(Fev | Alg, Flu) \cdot P(Alg | Flu) \\
 & = \underset{Cgh \perp Fev}{P(Cgh | Alg, Flu)} \cdot \underset{Fev \perp Alg}{P(Fev | Flu)} \cdot \underset{Alg \perp Flu}{P(Alg)} \cdot P(Flu)
 \end{aligned}$$

$$c) \text{ Prior: } P(z_0 = j) = \pi_j$$

$$\text{Transition: } P(z_t = k | z_{t-1} = j, u_t = i) = A_{(ij)k}$$

$$\text{Observation: } P(x_t = l | z_t = j, u_t = i) = B_{(ij)l}$$

d) Proof by Contradiction

Assume  $P(x_1, x_2 | y) = P(x_1 | y)P(x_2 | y)$

Let  $x_1 = 1$ ,  $x_2 = 1$  and  $y = 1$

We have  $P(x_1 | y) = P(x_1) = 0.5$

$P(x_2 | y) = P(x_2) = 0.5$

$P(x_1 | y) \cdot P(x_2 | y) = 0.5^2 = 0.25$

Based on the nature of XOR

we can see

$x_1$	$x_2$	$y$
0	0	0
1	0	1
0	1	1
1	1	0

$\rightarrow P(x_1=1, x_2=1 | y=1) = 0$

So when  $x_1 = 1$ ,  $x_2 = 1$  and  $y = 1$

we have

$P(x_1, x_2 | y) = P(x_1 | y)P(x_2 | y)$

$0 \neq 0.25$

This is a contradiction so we know

$P(x_1, x_2 | y) \neq P(x_1 | y)P(x_2 | y)$