(a)

We choose k = 37, since the letters C, I, J, K, L, M, O, P, S, U, V, W, X, Y, Z have similar uppercase and lowercase, their images will be quite similar and hence should be put in the same cluster.

(b)

First, we preprocess the datasets using prep.py:

```python
1.  #!/usr/bin/env python
2.
3.  import struct
4.  import pickle
5.
6.  # SIZE: 124800 for training; 20800 for testing
7.  # DIM: 28*28=784 pixels
8.
9.  train_images = "data/emnist-letters-train-images-idx3-ubyte"
10. train_labels = "data/emnist-letters-train-labels-idx1-ubyte"
11. test_images = "data/emnist-letters-test-images-idx3-ubyte"
12. test_labels = "data/emnist-letters-test-labels-idx1-ubyte"
13.
14. def load_img(filepath):
15.     data = []
16.     print 'Loading: %s' % (filepath)
17.     with open(filepath, 'rb') as f:
18.         magic, size, rows, cols = struct.unpack('>IIII', f.read(16))
19.         for img in range(size):
20.             data.append(struct.unpack('>784B', f.read(784)))
21.             if (img+1) % 10000 == 0:
22.                 print 'Loaded: %s' % (str(img+1))
23.     return data
24.
25. def load_lbl(filepath):
26.     print 'Loading: %s' % (filepath)
27.     data = []
28.     with open(filepath, 'rb') as f:
29.         magic, size = struct.unpack('>II', f.read(8))
30.         for img in range(size):
31.             data.append(struct.unpack('>B', f.read(1))[0])
32.             if (img+1) % 10000 == 0:
33.                 print 'Loaded: %s' % (str(img+1))
34.     return data
35.
36. def write_to_file(filepath, data):
37.     print 'Writing %s to file...' % (filepath)
38.     with open(filepath, 'w') as f:
39.         for i, img in enumerate(data):
40.             f.write(str(i))
41.             if (i+1) % 10000 == 0:
42.                 print 'Writing img: %s' % (str(i+1))
43.             for pixel in img:
44.                 f.write(' ' + str(pixel))
45.             f.write('\n')
46.
47. def write_to_pkl(filepath, data):
48.     print 'Writing %s to pkl...' % (filepath)
49.     with open(filepath, 'wb') as f:
50.         pickle.dump(data, f)
51.
52. write_to_file("train_images", load_img(train_images))
53. write_to_pkl("train_labels.pkl", load_lbl(train_labels))
54. write_to_file("test_images", load_img(test_images))
```

```
55. write_to_pkl("test_labels.pkl", load_lbl(test_labels))
```

We use init_cens.py to generate k = 37 random centroids:

```
1.  #!/usr/bin/env python
2.
3.  import numpy as np
4.  import pickle
5.
6.  def init_centroids(k = 37, dim = 28*28):
7.      cens = {}
8.      for i in range(k):
9.          cens[i] = np.random.uniform(low=0, high=256, size=(dim,))
10.     return cens
11.
12. init_cens = init_centroids()
13. with open('init_cens.pkl', 'wb') as f:
14.     pickle.dump(init_cens, f)
```

For each iteration of k-means algorithm, we run the following MapReduce job:

mapper.py

```
1.  #!/usr/bin/env python
2.  import sys
3.  import numpy as np
4.  import math
5.  import pickle
6.
7.  def euclid_dist(v1, v2):
8.      return math.sqrt(np.dot(v1-v2, v1-v2))
9.
10. # Get the list of old centroids
11. with open('old_cens.pkl', 'rb') as f:
12.     cens = pickle.load(f)
13.
14. partial_sum = {}
15.
16. for line in sys.stdin:
17.     dataline = line.strip().split()
18.     # First element is only for numbering the images
19.     img = np.asarray(dataline[1:]).astype(int)
20.
21.     # Get the closest centroid
22.     min_dist = 1e10 # Magic number for inf
23.     for i, cen in cens.items():
24.         dist = euclid_dist(img, cen)
25.         if dist < min_dist:
26.             min_dist = dist
27.             assigned_cen = i
28.
29.     # Add the current sample to the partial sum
30.     if assigned_cen not in partial_sum:
31.         partial_sum[assigned_cen] = {'count': 1, 'img': img}
32.     else:
33.         partial_sum[assigned_cen]['count'] += 1
34.         partial_sum[assigned_cen]['img'] += img
35.
36. # Print the partial sums from this mapper
37. for cen in partial_sum:
38.     img = ' '.join([str(pixel) for pixel in partial_sum[cen]['img']])
39.     print '%s\t%s %s' % (str(cen), str(partial_sum[cen]['count']), img)
```

reducer.py

```python
1.  #!/usr/bin/env python
2.
3.  import sys
4.  import numpy as np
5.
6.  curr_cen = None
7.
8.  for line in sys.stdin:
9.      dataline = line.strip().split()
10.     img = np.asarray(dataline).astype(int)
11.     cen = img[0]
12.     count = img[1]
13.
14.     if curr_cen == None:
15.         curr_count = count
16.         curr_sum = img[2:]
17.         curr_cen = cen
18.     elif curr_cen == cen:
19.         # Summing points in the same cluster
20.         curr_sum += img[2:]
21.         curr_count += count
22.     else:
23.         # Print out the new centroids: cen + pixel
24.         avg = ' '.join([str(round(float(pixel)/curr_count, 3)) for pixel in curr_su
    m])
25.         print '%s\t%s' % (str(curr_cen), avg)
26.
27.         curr_count = count
28.         curr_sum = img[2:]
29.         curr_cen = cen
30.
31. # Print out the new centroids
32. avg = ' '.join([str(round(pixel/curr_count, 3)) for pixel in curr_sum])
33. print '%s\t%s' % (str(curr_cen), avg)
```

At the end of each round, we check whether the total distances between the old and the new centroids are below a certain threshold using check.py:

```python
1.  #!/usr/bin/env python
2.
3.  import numpy as np
4.  import math
5.  import pickle
6.
7.  with open('old_cens.pkl', 'rb') as f:
8.      old_cens = pickle.load(f)
9.
10. # Convert new_cens to numpy array, and dump it to pkl file
11. new_cens = {}
12. with open("new_cens", 'r') as f:
13.     for line in f.readlines():
14.         dataline = line.strip().split()
15.         cen = int(dataline[0])
16.         new_cens[cen] = np.asarray(dataline[1:]).astype(float)
17. # If centroids not updated
18. for cen in old_cens:
19.     if cen not in new_cens:
20.         new_cens[cen] = old_cens[cen]
21.
22. # Write new centroids to pickle file
23. with open('new_cens.pkl', 'wb') as f:
24.     pickle.dump(new_cens, f)
```

```
25.
26. # Calculate total distance between corresponding old and new centroids
27. diff = 0
28. for cen in new_cens:
29.     v1 = new_cens[cen]
30.     v2 = old_cens[cen]
31.     diff += math.sqrt(np.dot(v1-v2, v1-v2))
32.
33. # Logging the diff
34. with open("diff.log", 'a+') as f:
35.     f.write(str(diff) + '\n')
36.
37. # Set stopping threshold
38. THRESHOLD = 30
39.
40. # If diff falls below given threshold, stop
41. if diff <= THRESHOLD:
42.     print 'STOP'
```

We wrap the shell commands into mr.sh, and limit the number of iterations to at most 150 (to avoid running for too long):

```
1.  #!/bin/bash
2.  rm new_cens.pkl old_cens.pkl init_cens.pkl diff.log new_cens
3.  python init_cens.py
4.  cp init_cens.pkl old_cens.pkl
5.
6.  MAX_ITER=150
7.
8.  for ITER in {1..150}
9.  do
10.     hdfs dfs -rm -r out
11.     hadoop jar /usr/hdp/2.4.2.0-258/hadoop-mapreduce/hadoop-streaming.jar -
    D mapred.job.name="Round_$ITER/$MAX_ITER" -D mapred.map.tasks=21 -
    D mapred.reduce.tasks=7 -file mapper.py -mapper mapper.py -file reducer.py -
    reducer reducer.py -input train_images -output out -file old_cens.pkl
12.     rm new_cens
13.     hdfs dfs -cat out/* >> new_cens
14.     if [[ $(python check.py) == "STOP" ]]
15.     then
16.         printf "CONVERGED: BELOW THRESHOLD\n"
17.         break
18.     else
19.         rm old_cens.pkl
20.         mv new_cens.pkl old_cens.pkl
21.     fi
22. done
23. mv new_cens final_cens
24. rm new_cens.pkl old_cens.pkl
25. hdfs dfs -rm -r out
26. hadoop jar /usr/hdp/2.4.2.0-258/hadoop-mapreduce/hadoop-streaming.jar -
    D mapred.job.name="Assign" -D mapred.map.tasks=10 -D mapred.reduce.tasks=2 -
    file assign_mapper.py -mapper assign_mapper.py -file assign_reducer.py -
    reducer assign_reducer.py -input train_images -output out -file final_cens
27. hdfs dfs -cat out/* >> tmp
28. sort -k1 -n tmp >> assign
29. rm tmp
30. python accuracy.py
```

Finally, we use the centroids from the algorithm to assign the data points accordingly and output the statistics using a simple MapReduce job:

assign_mapper.py

```python
1.  #!/usr/bin/env python
2.
3.  import numpy as np
4.  import math
5.  import sys
6.
7.  def euclid_dist(v1, v2):
8.      return math.sqrt(np.dot(v1-v2, v1-v2))
9.
10. # Get the list of final centroids
11. cens = {}
12. with open('final_cens', 'r') as f:
13.     for line in f.readlines():
14.         dataline = line.strip().split()
15.         cens[int(dataline[0])] = np.asarray(dataline[1:]).astype(float)
16.
17. for line in sys.stdin:
18.     dataline = line.strip().split()
19.     number = dataline[0]
20.     dataline = dataline[1:]
21.     img = np.asarray(dataline).astype(int)
22.
23.     # Get the closest centroid
24.     min_dist = 1e10 # Magic number for inf
25.     for i, cen in cens.items():
26.         dist = euclid_dist(cen, img)
27.         if dist < min_dist:
28.             min_dist = dist
29.             assigned_cen = i
30.
31.     print '%s\t%s' % (number, str(assigned_cen))
```

assign_reducer.py

```python
1.  #!/usr/bin/env python
2.
3.  import sys
4.
5.  for line in sys.stdin:
6.      dataline = line.strip().split()
7.      for c in dataline:
8.          print '%s' % (c),
9.      print ''
```

(c)

We use accuracy.py to calculate accuracy for each output of the corresponding random seeds:

```python
1.  #!/usr/bin/env python
2.
3.  import pickle
4.
5.  # Load labels
6.  with open("train_labels.pkl", 'rb') as f:
7.      labels = pickle.load(f)
8.
9.  def file_to_dict(path):
10.     dict = {}
11.     with open(path, 'r') as f:
12.         for line in f.readlines():
13.             dataline = line.strip().split()
14.             dict[int(dataline[0])] = int(dataline[1])
15.     return dict
```

```
16.
17. NUM_CLUSTER = 37
18. NUM_LABEL = 26
19.
20. # Load cluster assignment
21. path = 'assign'
22.
23. assignment = file_to_dict(path)
24. # Record clusters and the assigned img with their true label
25. # Note: EMNIST labels start from 1, not 0
26. count_cluster = {cluster: {label: [] for label in range(1, NUM_LABEL+1)}\
27.                   for cluster in range(NUM_CLUSTER)}
28.
29. # Loop through all assigned images
30. for i, c in assignment.items():
31.     print 'Count: %s' % (str(i))
32.     # Increase the corr. label count of current cluster by 1
33.     # Note: EMNIST labels start from 1, not 0
34.     count_cluster[c][labels[i]].append(i)
35.
36. # Get the cluster label and calculate its accuracy
37. accuracy = {cluster: {'label': -1, 'total': 0, 'correct': 0, 'acc': -1}\
38.                     for cluster in range(NUM_CLUSTER)}
39. for c, cluster in count_cluster.items():
40.     print 'Calc acc: %s' % (str(c))
41.     max_label = -1
42.     best_label = -1
43.     for l, label in cluster.items():
44.         accuracy[c]['total'] += len(label)
45.         if len(label) > max_label:
46.             max_label = len(label)
47.             best_label = l
48.     # Note: EMNIST labels start from 1, not 0
49.     accuracy[c]['label'] = best_label
50.     accuracy[c]['correct'] = max_label
51.     if accuracy[c]['total'] == accuracy[c]['correct'] and accuracy[c]['total'] == 0
    :
52.         accuracy[c]['acc'] = -1
53.     else:
54.         accuracy[c]['acc'] = round(float(accuracy[c]['correct']) / accuracy[c]['tot
    al'] * 100, 2)
55.
56. with open('accuracy', 'w') as f:
57.     total = 0
58.     correct = 0
59.     for c, cluster in accuracy.items():
60.         total += cluster['total']
61.         correct += cluster['correct']
62.         f.write(str(c) + ' ' + str(cluster['total']) + ' ' + str(cluster['label'])
    +
63.                 ' ' + str(cluster['correct']) + ' ' + str(cluster['acc']) + '\n')
64.     f.write("TOTAL: " + str(total) + '\n' +  "CORRECT: " + str(correct) + '\n' + "O
    VERALL ACCURACY: " + str(round(float(correct)/total*100, 2)) + '\n')
```

We report the corresponding statistics for each random seed in the following tables:

Table. 1. The Accuracy of Clustering Performance with Random Seed 1

| Cluster Number | # train images belongs to the cluster | Label of the cluster | # correctly clustered images | Classification Accuracy (%) |
|---|---|---|---|---|
| 0 | 0 | NA | 0 | NA |
| 1 | 4938 | 10 | 1743 | 35.3 |
| 2 | 5488 | 24 | 2298 | 41.87 |
| 3 | 3671 | 19 | 2878 | 78.4 |
| 4 | 0 | NA | 0 | NA |
| 5 | 4350 | 18 | 2324 | 53.43 |
| 6 | 6103 | 6 | 1647 | 26.99 |
| 7 | 4215 | 8 | 1297 | 30.77 |
| 8 | 0 | NA | 0 | NA |
| 9 | 0 | NA | 0 | NA |
| 10 | 3152 | 26 | 2372 | 75.25 |
| 11 | 3799 | 13 | 3353 | 88.26 |
| 12 | 3059 | 14 | 1701 | 55.61 |
| 13 | 4181 | 16 | 1545 | 36.95 |
| 14 | 0 | NA | 0 | NA |
| 15 | 3562 | 15 | 1695 | 47.59 |
| 16 | 4204 | 4 | 1352 | 32.16 |
| 17 | 3955 | 21 | 2023 | 51.15 |
| 18 | 4946 | 3 | 1196 | 24.18 |
| 19 | 5073 | 23 | 1022 | 20.15 |
| 20 | 3485 | 11 | 2001 | 57.42 |
| 21 | 4838 | 26 | 1140 | 23.56 |
| 22 | 4083 | 15 | 2089 | 51.16 |
| 23 | 4208 | 5 | 2144 | 50.95 |
| 24 | 5239 | 10 | 1252 | 23.9 |
| 25 | 0 | NA | 0 | NA |
| 26 | 0 | NA | 0 | NA |
| 27 | 3701 | 17 | 1321 | 35.69 |
| 28 | 0 | NA | 0 | NA |
| 29 | 4230 | 16 | 811 | 19.17 |
| 30 | 6692 | 9 | 2194 | 32.79 |
| 31 | 4212 | 8 | 1325 | 31.46 |
| 32 | 3688 | 23 | 2869 | 77.79 |
| 33 | 3718 | 22 | 2267 | 60.97 |
| 34 | 0 | NA | 0 | NA |
| 35 | 4789 | 14 | 1264 | 26.39 |
| 36 | 7221 | 9 | 1747 | 24.19 |
| Total Set | 124800 | NA | 50870 | 40.76 |

Table. 2. The Accuracy of Clustering Performance with Random Seed 2

| Cluster Number | # train images belongs to the cluster | Label of the cluster | # correctly clustered images | Classification Accuracy (%) |
|---|---|---|---|---|
| 0 | 0 | NA | 0 | NA |
| 1 | 5964 | 9 | 2110 | 35.38 |
| 2 | 3077 | 19 | 2409 | 78.29 |
| 3 | 3725 | 25 | 1042 | 27.97 |
| 4 | 4467 | 23 | 3527 | 78.96 |
| 5 | 3642 | 1 | 831 | 22.82 |
| 6 | 0 | NA | 0 | NA |
| 7 | 3736 | 4 | 1766 | 47.27 |
| 8 | 2804 | 22 | 1820 | 64.91 |
| 9 | 5054 | 6 | 1477 | 29.22 |
| 10 | 4227 | 5 | 1881 | 44.5 |
| 11 | 2468 | 21 | 847 | 34.32 |
| 12 | 4808 | 10 | 1871 | 38.91 |
| 13 | 3002 | 6 | 1073 | 35.74 |
| 14 | 2580 | 26 | 1500 | 58.14 |
| 15 | 2795 | 17 | 996 | 35.64 |
| 16 | 3989 | 18 | 2242 | 56.2 |
| 17 | 3848 | 3 | 1007 | 26.17 |
| 18 | 2826 | 14 | 1731 | 61.25 |
| 19 | 2380 | 26 | 1832 | 76.97 |
| 20 | 4181 | 15 | 2431 | 58.14 |
| 21 | 3304 | 11 | 1943 | 58.81 |
| 22 | 3890 | 16 | 1806 | 46.43 |
| 23 | 6536 | 9 | 1713 | 26.21 |
| 24 | 3320 | 19 | 1160 | 34.94 |
| 25 | 5193 | 24 | 2275 | 43.81 |
| 26 | 3308 | 8 | 915 | 27.66 |
| 27 | 3453 | 8 | 1087 | 31.48 |
| 28 | 4865 | 14 | 1223 | 25.14 |
| 29 | 3867 | 15 | 930 | 24.05 |
| 30 | 2401 | 22 | 1493 | 62.18 |
| 31 | 4682 | 13 | 3768 | 80.48 |
| 32 | 3301 | 21 | 1898 | 57.5 |
| 33 | 4094 | 16 | 789 | 19.27 |
| 34 | 3013 | 8 | 931 | 30.9 |
| 35 | 0 | NA | 0 | NA |
| 36 | 0 | NA | 0 | NA |
| Total Set | 124800 | NA | 54324 | 43.53 |

Table. 3. The Accuracy of Clustering Performance with Random Seed 3

| Cluster Number | # train images belongs to the cluster | Label of the cluster | # correctly clustered images | Classification Accuracy (%) |
|---|---|---|---|---|
| 0 | 0 | NA | 0 | NA |
| 1 | 3551 | 22 | 1326 | 37.34 |
| 2 | 3774 | 1 | 844 | 22.36 |
| 3 | 3577 | 21 | 1820 | 50.88 |
| 4 | 3680 | 3 | 954 | 25.92 |
| 5 | 4617 | 23 | 3551 | 76.91 |
| 6 | 4383 | 10 | 1448 | 33.04 |
| 7 | 3220 | 19 | 2467 | 76.61 |
| 8 | 6394 | 9 | 1763 | 27.57 |
| 9 | 2197 | 26 | 1692 | 77.01 |
| 10 | 4767 | 13 | 3776 | 79.21 |
| 11 | 2573 | 21 | 846 | 32.88 |
| 12 | 4028 | 8 | 1213 | 30.11 |
| 13 | 4914 | 14 | 1221 | 24.85 |
| 14 | 3422 | 11 | 2008 | 58.68 |
| 15 | 3708 | 16 | 1501 | 40.48 |
| 16 | 4053 | 8 | 1405 | 34.67 |
| 17 | 2178 | 26 | 1619 | 74.33 |
| 18 | 4160 | 18 | 2300 | 55.29 |
| 19 | 3460 | 25 | 1240 | 35.84 |
| 20 | 0 | NA | 0 | NA |
| 21 | 0 | NA | 0 | NA |
| 22 | 2557 | 22 | 1784 | 69.77 |
| 23 | 0 | NA | 0 | NA |
| 24 | 2895 | 10 | 1266 | 43.73 |
| 25 | 0 | NA | 0 | NA |
| 26 | 3325 | 15 | 1966 | 59.13 |
| 27 | 5235 | 24 | 2131 | 40.71 |
| 28 | 5868 | 16 | 1634 | 27.85 |
| 29 | 3043 | 14 | 1726 | 56.72 |
| 30 | 3200 | 4 | 1187 | 37.09 |
| 31 | 4128 | 16 | 870 | 21.08 |
| 32 | 6413 | 9 | 2107 | 32.86 |
| 33 | 3984 | 15 | 1779 | 44.65 |
| 34 | 3218 | 17 | 1187 | 36.89 |
| 35 | 3992 | 3 | 785 | 19.66 |
| 36 | 4286 | 5 | 2254 | 52.59 |
| Total Set | 124800 | NA | 53670 | 43 |

Best random seed: Random seed 2 and 3 have fewer clusters with no member compared to random seed 1, hence, together with better overall accuracy, they are somewhat better than random seed 1. Random seed 2 has slightly better overall accuracy (43.53%) compared to random seed 3 (43%), hence random seed 2 is the best seed (but only being better than seed 3 by a very small margin).

(d)

We will use the code of part (b) and (c), in particular the assign_mapper.py & assign_reducer.py (MapReduce job for assignment) and accuracy.py, to determine the accuracy of the model with random seed 2 on test set. Run the following shell script with corresponding input, which are the test dataset and the final centroids with random seed 2:

```
1.  hdfs dfs -rm -r out
2.  hadoop jar /usr/hdp/2.4.2.0-258/hadoop-mapreduce/hadoop-streaming.jar -
    D mapred.job.name="Assign" -D mapred.map.tasks=10 -D mapred.reduce.tasks=2 -
    file assign_mapper.py -mapper assign_mapper.py -file assign_reducer.py -
    reducer assign_reducer.py -input test_images -output out -file final_cens
3.  hdfs dfs -cat out/* >> tmp
4.  sort -k1 -n tmp >> assign
5.  rm tmp
6.  python accuracy.py
```

We report the accuracy in the following table:

Table. 4. The Accuracy of Clustering Performance on the test data

| Cluster Number | # train images belongs to the cluster | Label of the cluster | # correctly clustered images | Classification Accuracy (%) |
|---|---|---|---|---|
| 0 | 0 | NA | 0 | NA |
| 1 | 1040 | 9 | 350 | 33.65 |
| 2 | 523 | 19 | 410 | 78.39 |
| 3 | 613 | 25 | 176 | 28.71 |
| 4 | 768 | 23 | 613 | 79.82 |
| 5 | 607 | 1 | 154 | 25.37 |
| 6 | 0 | NA | 0 | NA |
| 7 | 605 | 4 | 302 | 49.92 |
| 8 | 497 | 22 | 311 | 62.58 |
| 9 | 846 | 6 | 242 | 28.61 |
| 10 | 703 | 5 | 326 | 46.37 |
| 11 | 423 | 21 | 134 | 31.68 |
| 12 | 774 | 10 | 306 | 39.53 |
| 13 | 490 | 6 | 165 | 33.67 |
| 14 | 426 | 26 | 250 | 58.69 |
| 15 | 462 | 17 | 180 | 38.96 |
| 16 | 617 | 18 | 348 | 56.4 |
| 17 | 608 | 3 | 157 | 25.82 |
| 18 | 479 | 14 | 289 | 60.33 |

| 19 | 414 | 26 | 325 | 78.5 |
|---|---|---|---|---|
| 20 | 702 | 15 | 420 | 59.83 |
| 21 | 539 | 11 | 312 | 57.88 |
| 22 | 641 | 16 | 313 | 48.83 |
| 23 | 1074 | 9 | 287 | 26.72 |
| 24 | 606 | 19 | 190 | 31.35 |
| 25 | 902 | 24 | 389 | 43.13 |
| 26 | 568 | 8 | 160 | 28.17 |
| 27 | 568 | 8 | 186 | 32.75 |
| 28 | 832 | 14 | 219 | 26.32 |
| 29 | 590 | 3 | 145 | 24.58 |
| 30 | 396 | 22 | 248 | 62.63 |
| 31 | 752 | 13 | 623 | 82.85 |
| 32 | 560 | 21 | 334 | 59.64 |
| 33 | 683 | 1 | 130 | 19.03 |
| 34 | 492 | 8 | 131 | 26.63 |
| 35 | 0 | NA | 0 | NA |
| 36 | 0 | NA | 0 | NA |
| Total Set | 20800 | NA | 9125 | 43.87 |