

EECS 445 HW3

JunCheng An

$$1. (a) \frac{dL}{d\hat{y}_c} = \frac{d}{d\hat{y}_c} - \sum_{i=1}^3 y_i \log \hat{y}_i$$

$$= -y_c \cdot \frac{1}{\hat{y}_c}$$

$$(b) \frac{d\hat{y}_c}{dZ_k} = \frac{d}{dZ_k} \frac{\exp(Z_c)}{\sum_{i=1}^3 \exp(Z_i)}$$

$$\cdot c=k : = \frac{\exp(Z_c) \cdot \sum_{i=1}^3 \exp(Z_i) - \exp(Z_c)}{(\sum_{i=1}^3 \exp(Z_i))^2}$$

$$= \frac{\exp(Z_c)}{(\sum_{i=1}^3 \exp(Z_i))} - \frac{\exp(Z_c)^2}{(\sum_{i=1}^3 \exp(Z_i))^2}$$

$$= \hat{y}_c - \hat{y}_c^2 \quad \text{or} \quad \hat{y}_k - \hat{y}_k^2$$

$$\cdot c \neq k : = \frac{-\exp(Z_c) \cdot \exp(Z_k)}{(\sum_{i=1}^3 \exp(Z_i))^2}$$

$$= \left( \frac{\exp(Z_c)}{\sum_{i=1}^3 \exp(Z_i)} \right) \cdot \left( \frac{\exp(Z_k)}{\sum_{i=1}^3 \exp(Z_i)} \right) \cdot (-1)$$

$$= -\hat{y}_c \cdot \hat{y}_k$$

(c) • using indicator function to write cb1:

$$c=k: \frac{d\hat{y}_c}{dZ_k} = \hat{y}_c - \hat{y}_c \cdot \hat{y}_k$$

$$\Rightarrow \frac{d\hat{y}_c}{dZ_k} = [c=k] \cdot \hat{y}_c - \hat{y}_c \cdot \hat{y}_k$$

$$\begin{aligned} \cdot \frac{dL}{dZ_k} &= \sum_{c=1}^3 \frac{d\hat{y}_c}{dZ_k} \cdot \frac{dL}{d\hat{y}_c} \\ &= \sum_{c=1}^3 ([c=k] \hat{y}_c - \hat{y}_c \hat{y}_k) \cdot y_c \cdot \frac{-1}{\hat{y}_c} \end{aligned}$$

$$c_{t,c} = \sum_{k=1}^3 (\llbracket t == k \rrbracket - \hat{q}_k) \cdot (-1)^{\llbracket t == c \rrbracket}$$

$$= \sum_{k=1}^3 \llbracket t == c \rrbracket \hat{q}_k - \llbracket k == c \rrbracket \llbracket t == c \rrbracket$$

$$\Rightarrow \llbracket k == c \rrbracket \cdot \llbracket t == c \rrbracket = 1 \text{ while } k = t = c$$

= 0 otherwise

$$\therefore \sum_{c=1}^3 \llbracket k == c \rrbracket \llbracket t == c \rrbracket = \llbracket k == t \rrbracket$$

$$\Rightarrow \sum_{c=1}^3 \llbracket t == c \rrbracket = 1$$

$$\therefore = \hat{q}_k - \llbracket k == t \rrbracket$$

$$\therefore \frac{dL}{dk} = \hat{q}_k - \llbracket k == t \rrbracket$$

2. (a)

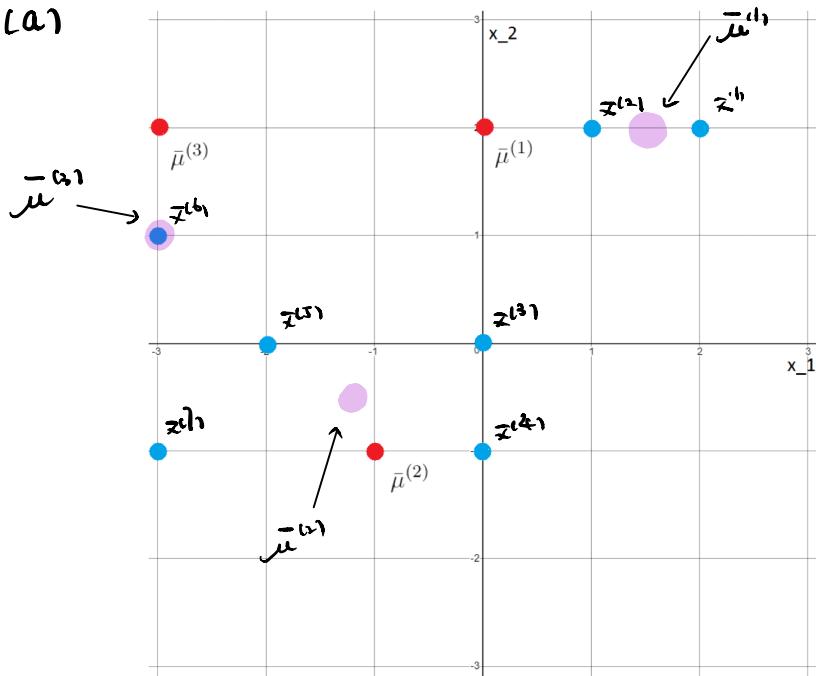


Figure 2: Blue — points, Red — means

- $\bar{x}^{(1)}$ 's nearest mean is  $\bar{\mu}^{(1)}$   $\Rightarrow \bar{\mu}_{\text{sum}} = [2.2]$
- $\bar{x}^{(2)}$ 's nearest mean is  $\bar{\mu}^{(1)}$   $\Rightarrow \bar{\mu}_{\text{sum}} = [3.4]$
- $\bar{x}^{(3)}$ 's nearest mean is  $\bar{\mu}^{(2)}$   $\Rightarrow \bar{\mu}_{\text{sum}} = [0.0]$
- $\bar{x}^{(4)}$ 's nearest mean is  $\bar{\mu}^{(2)}$   $\Rightarrow \bar{\mu}_{\text{sum}} = [0.0]$
- $\bar{x}^{(5)}$ 's nearest mean is  $\bar{\mu}^{(3)}$   $\Rightarrow \bar{\mu}_{\text{sum}} = [-2.0]$
- $\bar{x}^{(6)}$ 's nearest mean is  $\bar{\mu}^{(3)}$   $\Rightarrow \bar{\mu}_{\text{sum}} = [-3.0]$
- $\bar{x}^{(7)}$ 's nearest mean is  $\bar{\mu}^{(1)}$   $\Rightarrow \bar{\mu}_{\text{sum}} = [-5.0]$
- $\therefore \bar{\mu}^{(1)} = [3.4] \div 2 = [1.5, 2]$
- $\bar{\mu}^{(2)} = [-5.0] \div 4 = [-1.25, -0.5]$
- $\bar{\mu}^{(3)} = [-3.0]$ 
  - position is marked as

(b) only 1 iteration and converge, no more update !

$$\begin{aligned}
 \text{(c)} \quad & \frac{d}{d\bar{\mu}} \|\bar{\mu}\|^2 + \sum_{i \in C} \|\bar{x}^i - \bar{\mu}\|^2 \\
 &= 2\bar{\mu} - \sum_{i \in C} 2 \cdot (\bar{x}^i - \bar{\mu}) = 0 \\
 &\bar{\mu} + \sum_{i \in C} \bar{\mu} = \sum_{i \in C} \bar{x}^i \\
 &(1+|C|) \cdot \bar{\mu} = \sum_{i \in C} \bar{x}^i \\
 \therefore \bar{\mu}^* &= \frac{\sum_{i \in C} \bar{x}^i}{1+|C|}
 \end{aligned}$$

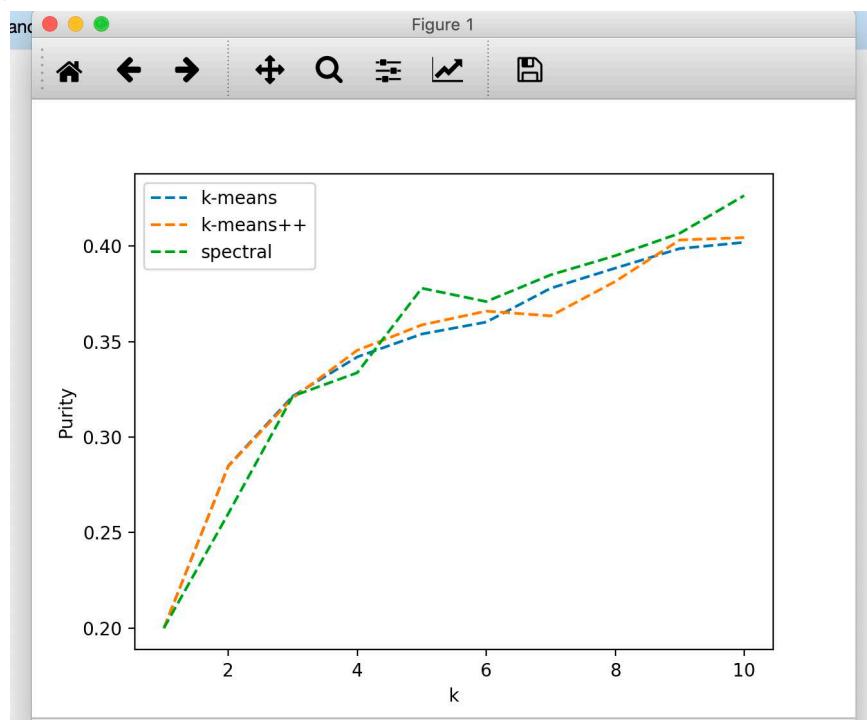
- (d)
- As the formula calculated in (c)  
and see normal formula  $\bar{\mu}^* = \frac{\sum_{i \in C} \bar{x}^i}{|C|}$
  - We can find that we are just adding a point

$$\bar{x}^* = [0, 0]$$

3. (a) Because k-mean clustering is not guaranteed to converge to a global minimum, a poor initialization would result a poor local minimum.

(b) ~ (c) Code Only

(h) graph as shown below



- Theoretically, we are looking for a "elbow" on the graph.
    - According to the image, I would choose
      - $k=7$  for k-means
      - $k=6$  for k-means++
      - $k=5$  for spectral
- (pick the point that starts to increase slower)

	K	Ave	MAX	MEAN
K-means	7	0.3720	0.4075	0.3450
K-means++	6	0.3687	0.3975	0.3350
spectral	5	0.3765	0.3825	0.3650

- By the graph, we can find that generally spectral performs better than k-mean methods. And among these k-mean methods, different initialization almost make no difference, their performances are close to each other.
- By the chart given, we will also choose spectral because it is more stable (low distance between min and max) and it has highest average value.
- We have 5 labels in our dataset so random assignment would result in purity = 0.2  
 $\Rightarrow$  Obviously, our algorithms are doing better

```
[0.375, 0.3775, 0.375, 0.3825, 0.3825, 0.365, 0.3825, 0.375, 0.3775, 0.3725]
k_final: Ave: 0.3720, MAX: 0.4075, MIN: 0.3450
p_final: Ave: 0.3687, MAX: 0.3975, MIN: 0.3350
s_final: Ave: 0.3765, MAX: 0.3825, MIN: 0.3650
```

4. (a) similarity matrix :

$$\begin{bmatrix} 1 & 0.75 & 0.5 & 0 & 0.25 \\ 0.75 & 1 & 0 & 0.5 & 0 \\ 0.5 & 0 & 1 & 0.8 & 0.4 \\ 0 & 0.5 & 0.8 & 1 & 0 \\ 0.25 & 0 & 0.4 & 0 & 1 \end{bmatrix} = w$$

Degree matrix:

$$\begin{bmatrix} 2.5 & 0 & 0 & 0 & 0 \\ 0 & 2.25 & 0 & 0 & 0 \\ 0 & 0 & 2.7 & 0 & 0 \\ 0 & 0 & 0 & 2.3 & 0 \\ 0 & 0 & 0 & 0 & 1.65 \end{bmatrix} = D$$

Laplacian  $L = D - w$

$$= \begin{bmatrix} 1.5 & -0.75 & -0.5 & 0 & -0.25 \\ -0.75 & 1.25 & 0 & -0.5 & 0 \\ -0.5 & 0 & 1.7 & -0.8 & -0.4 \\ 0 & -0.5 & -0.8 & 1.3 & 0 \\ -0.25 & 0 & -0.4 & 0 & 0.65 \end{bmatrix}$$

(b) 1<sup>st</sup> eigenvalue :  $2.91337056 \times 10^{-1}$

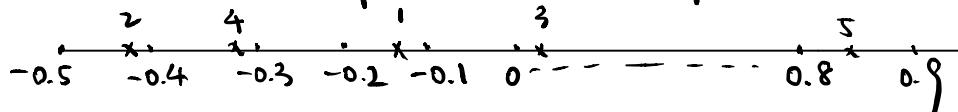
eigenvector :  $[-0.4472136, -0.4472136, -0.4472136, -0.4472136, -0.4472136]$

• 2<sup>nd</sup> eigenvalue :  $6.745586 \times 10^{-1}$

eigenvector :  $[-0.121416, -0.42937, 0.02441, -0.31203, 0.838413]$

(c) first dimension coordinate, which is the value of first eigenvector, are all same.

$\Rightarrow$  Thus, we only consider 2<sup>nd</sup> eigenvector values



• Therefore cluster 1 : 1.2.3.4

cluster 2 : 5

(d) 1<sup>st</sup> eigenvalue :  $2.91337056 \times 10^{-6}$

eigenvector :  $[ -0.4472136, -0.4472136, -0.4472136, -0.4472136, -0.4472136 ]$

• 2<sup>nd</sup> eigenvalue :  $6.745586 \times 10^{-1}$

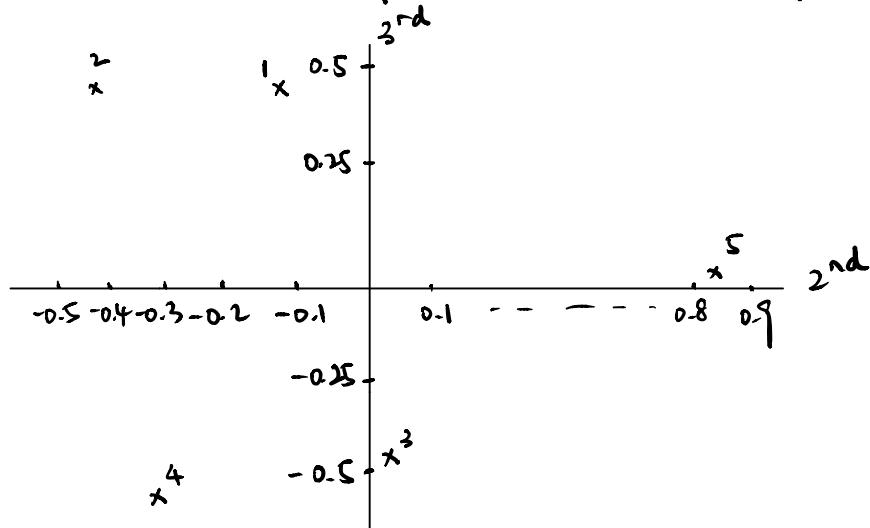
eigenvector :  $[ -0.121416, -0.42937, 0.02441, -0.31203, 0.838413 ]$

• 3<sup>rd</sup> eigenvector :  $1.149806$

eigenvector :  $[ 0.4758, 0.46409, -0.46659, -0.62071, 0.0874 ]$

i.e., first dimension coordinate, which is the value of first eigenvector, are all same.

$\Rightarrow$  Thus, we only consider 2<sup>nd</sup> & 3<sup>rd</sup> eigenvector values



• Therefore : cluster 1 : 1, 2

cluster 2 : 3, 4

cluster 3 : 5

# CODE APPENDIX

## clustering\_classes.py

```
"""
EECS 445 - Introduction to Machine Learning
Winter 2019 - Homework 3
Clustering Classes
"""

import numpy as np
from scipy import stats

class Point(object):
    """
    Represents a data point
    """

    def __init__(self, features, label=None):
        """
        Initialize label and attributes
        """
        self.features = features
        self.label = label

    def dimensionality(self):
        """Returns dimension of the point"""
        return len(self.features)

    def get_features(self):
        """Returns features"""
        return self.features

    def distance(self, other):
        """
        other: point, to which we are measuring distance to
        Return Euclidean distance of this point with other
        """
        # TODO: Implement this function
        return np.linalg.norm(self.features - other.features)

    def get_label(self):
        """Returns label"""
        return self.label

class Cluster(object):
    """
    A Cluster is defined as a set of elements
    """

    def __init__(self, points):
        """
        Elements of a cluster are saved in a list, self.points
        """
        self.points = points
```

```

def get_points(self):
    """Returns points in the cluster as a list"""
    return self.points

def get_label(self):
    """Returns label of the cluster, which is determined by the
       mode of labels"""
    labels = [point.get_label() for point in self.points]
    cluster_label, count = stats.mode(labels)
    return cluster_label[0]

def get_purity(self):
    """Returns number of points in cluster and the number of points
       with the most common label"""
    labels = [point.get_label() for point in self.points]
    cluster_label, count = stats.mode(labels)
    return len(labels), np.float64(count)

def get_centroid(self):
    """Returns centroid of the cluster"""
    # TODO: Implement this function
    if len(self.points) == 0: return Point(0)
    features = [point.get_features() for point in self.points]
    return Point(np.sum(features, axis = 0) / len(features))

def equivalent(self, other):
    """
    other: Cluster, what we are comparing this Cluster to
    Returns true if both Clusters are equivalent, or false otherwise
    """
    if len(self.get_points()) != len(other.get_points()):
        return False
    matched = []
    for p1 in self.get_points():
        for point2 in other.get_points():
            if p1.distance(point2) == 0 and point2 not in matched:
                matched.append(point2)
    return len(matched) == len(self.get_points())

class ClusterSet(object):
    """
    A ClusterSet is defined as a list of clusters
    """

    def __init__(self):
        """
        Initialize an empty set, without any clusters
        """
        self.clusters = []

    def add(self, c):
        """
        c: Cluster
        Appends a cluster c to the end of the cluster list
        only if it doesn't already exist in the ClusterSet.
        """

```

```

If it is already in self.clusters, raise a ValueError
"""
if c in self.clusters:
    raise ValueError
self.clusters.append(c)

def get_clusters(self):
    """Returns clusters in the ClusterSet"""
    return self.clusters[:]

def get_centroids(self):
    """Returns centroids of each cluster in the ClusterSet as a list"""
    # TODO: Implement this function

    centroids = [cluster.get_centroid() for cluster in self.clusters]
    return centroids

def get_score(self):
    """
    Returns accuracy of the clustering given by the clusters
    in ClusterSet object
    """

    total_correct = 0
    total = 0
    for c in self.clusters:
        n, n_correct = c.get_purity()
        total = total + n
        total_correct = total_correct + n_correct

    return total_correct / float(total)

def num_clusters(self):
    """Returns number of clusters in the ClusterSet"""
    return len(self.clusters)

def equivalent(self, other):
    """
    other: another ClusterSet object
    Returns true if both ClusterSets are equivalent, or false otherwise
    """

    if len(self.get_clusters()) != len(other.get_clusters()):
        return False
    matched = []
    for c1 in self.get_clusters():
        for c2 in other.get_clusters():
            if c1.equivalent(c2) and c2 not in matched:
                matched.append(c2)
    return len(matched) == len(self.get_clusters())

```

## clustering.py

```
"""
EECS 445 - Introduction to Machine Learning
Winter 2019 - Homework 3
Clustering
"""

import numpy as np
import random
import matplotlib.pyplot as plt
from sklearn.cluster import SpectralClustering
from sklearn.metrics.pairwise import pairwise_distances
from operator import methodcaller

from clustering_classes import Cluster, ClusterSet, Point
from data.landmarks import LandmarksDataset
from utils import denormalize_image

def random_init(points, k):
    """
    Arguments:
        points: a list of point objects
        k: Number of initial centroids/medoids
    Returns:
        List of k unique points randomly selected from points
    """
    # TODO: Implement this function
    random_result = random.sample(range(0, len(points)), k)
    result = []
    for i in range(k) : result += [points[random_result[i]]]
    return result

def k_means_pp_init(points, k):
    """
    Arguments:
        points: a list of point objects
        k: Number of initial centroids/medoids
    Returns:
        List of k unique points selected from points
    """
    # TODO: Implement this function
    result = [points[random.randint(0, len(points) - 1)]]

    while len(result) != k:
        distance = np.array([])
        for point in points:
            min_dis = float("inf")
            for centroid in result: min_dis = min(min_dis, point.distance(centroid))
            distance = np.append(distance, min_dis ** 2)
        distance = distance / np.sum(distance)
        result += [np.random.choice(points, 1, p = distance)[0]]

    return result
```

```

def k_means(points, k, init='random'):
    """
    Clusters points into k clusters using k_means clustering.
    Arguments:
        points: a list of Point objects
        k: the number of clusters
        init: The method of initialization. One of ['random', 'kpp'].
            If init='kpp', use k_means_pp_init to initialize clusters.
            If init='random', use random_init to initialize clusters.
            Default value 'random'.
    Returns:
        Instance of ClusterSet with k clusters
    """
    # TODO: Implement this function
    centroids = []
    if init == 'random': centroids = random_init(points, k)
    if init == 'kpp': centroids = k_means_pp_init(points, k)
    points_cluster = [[] for i in range(k)]

    for n in range(len(points)):
        min_dis, min_cluster = float("inf"), 0
        for i in range(k):
            if min_dis > points[n].distance(centroids[i]):
                min_cluster = i
                min_dis = points[n].distance(centroids[i])
        points_cluster[min_cluster] += [points[n]]

    clusters = [Cluster(point) for point in points_cluster]
    Cluster_new = ClusterSet()
    for cluster in clusters: Cluster_new.add(cluster)
    Cluster_set = ClusterSet()

    while (not Cluster_new.equivalent(Cluster_set)):
        Cluster_set = Cluster_new
        centroids = Cluster_set.get_centroids()
        points_cluster = [[] for i in range(k)]

        for n in range(len(points)):
            min_dis, min_cluster = float("inf"), 0
            for i in range(k):
                if min_dis > points[n].distance(centroids[i]):
                    min_cluster = i
                    min_dis = points[n].distance(centroids[i])
            points_cluster[min_cluster] += [points[n]]

        clusters = [Cluster(point) for point in points_cluster]
        Cluster_new = ClusterSet()
        for cluster in clusters: Cluster_new.add(cluster)

    return Cluster_new

def spectral_clustering(points, k):
    """
    Uses sklearn's spectral clustering implementation to cluster the input

```

```

data into k clusters
Arguments:
    points: a list of Points objects
    k: the number of clusters
Returns:
    Instance of ClusterSet with k clusters
"""

X = np.array([point.get_features() for point in points])
spectral = SpectralClustering(
    n_clusters=k, n_init=1, affinity='nearest_neighbors', n_neighbors=50)
y_pred = spectral.fit_predict(X)
clusters = ClusterSet()
for i in range(k):
    cluster_members = [p for j, p in enumerate(points) if y_pred[j] == i]
    clusters.add(Cluster(cluster_members))
return clusters

def plot_performance(k_means_scores, kpp_scores, spec_scores, k_vals):
    """
    Uses matplotlib to generate a graph of performance vs. k
    Arguments:
        k_means_scores: A list of len(k_vals) average purity scores from
            running the k-means algorithm with random initialization
        kpp_scores: A list of len(k_vals) average purity scores from running
            the k-means algorithm with k-means++ initialization
        spec_scores: A list of len(k_vals) average purity scores from running
            the spectral clustering algorithm
        k_vals: A list of integer k values used to calculate the above scores
    """
    # TODO: Implement this function
    plt.xlabel("k")
    plt.ylabel("Purity")
    plt.plot(k_vals, k_means_scores, label = "k-means", linestyle = '--')
    plt.plot(k_vals, kpp_scores, label = "k-means++", linestyle = '--')
    plt.plot(k_vals, spec_scores, label = "spectral", linestyle = '--')
    plt.legend()
    plt.show()

def get_data():
    """
    Retrieves the data to be used for the k-means clustering as a list of
    Point objects
    """
    landmarks = LandmarksDataset(num_classes=5)
    X, y = landmarks.get_batch('train', batch_size=400)
    X = X.reshape((len(X), -1))
    return [Point(image, label) for image, label in zip(X, y)]

def visualize_clusters(kmeans, kpp, spectral):
    """
    Uses matplotlib to generate plots of representative images for each
    of the clustering algorithm. In each image, every row is from the same
    cluster, and from leftmost image is the medoid. Intra-cluster distance
    increases as we go from left to right.
    Arguments:
        - kmeans, kpp, and spectral: ClusterSet instances
    """

```

```

def get_medoid_and_neighbors(points, num=4):
    D = pairwise_distances([p.features for p in points])
    distances = D.mean(axis=0)
    return np.array(points)[np.argsort(distances)][:num].tolist()

names = ['k-means', 'k-means++', 'spectral']
cluster_sets = [kmeans, kpp, spectral]
clusters_s = [sorted(cs.get_clusters(),
                     key=methodcaller('get_label')) for cs in cluster_sets]

for i, clusters in enumerate(clusters_s):
    num = 4
    k = len(clusters)
    fig, axes = plt.subplots(nrows=k, ncols=num, figsize=(8,8))
    plt.suptitle(names[i])
    for j in range(k):
        pts = get_medoid_and_neighbors(clusters[j].get_points(), num)
        for n in range(len(pts)):
            axes[j,n].imshow(denormalize_image(
                np.reshape(pts[n].features,(32,32,3))), interpolation='bicubic')
            axes[j,0].set_ylabel(clusters[j].get_label())

    for ax in axes.flatten():
        ax.set_xticks([])
        ax.set_yticks([])

    plt.savefig('4j_clusters_viz_{}.png'.format(names[i]),
                dpi=200, bbox_inches='tight')

def main():
    points = get_data()
    # TODO: Implement this function
    # for 3.h and 3.i

    """ 3.j """
    # Display representative examples of each cluster for clustering algorithms
    k_final, p_final, s_final = np.zeros(10), np.zeros(10), np.zeros(10)
    k_vals = [i for i in range(1,11)]
    # np.random.seed(42)
    for i in range(10):
        kmeans = np.array([k_means(points, i, 'random').get_score() for i in k_vals])
        kpp = np.array([k_means(points, i, 'kpp').get_score() for i in k_vals])
        spectral = np.array([spectral_clustering(points, i).get_score() for i in k_vals])
        k_final += kmeans
        p_final += kpp
        s_final += spectral
    k_final = k_final / 10
    p_final = p_final / 10
    s_final = s_final / 10
    plot_performance(k_final, p_final, s_final, k_vals)

    k_final, p_final, s_final = [], [], []
    k_vals = [i for i in range(1,11)]
    for i in range(10):
        k_final += [k_means(points, 7, 'random').get_score()]
        p_final += [k_means(points, 6, 'kpp').get_score()]
        s_final += [spectral_clustering(points, 5).get_score()]

```

```
print(s_final)
print("k_final: Ave: %.4f, MAX: %.4f, MIN: %.4f" % (np.mean(k_final), max(k_final), min(k_final)))
print("p_final: Ave: %.4f, MAX: %.4f, MIN: %.4f" % (np.mean(p_final), max(p_final), min(p_final)))
print("s_final: Ave: %.4f, MAX: %.4f, MIN: %.4f" % (np.mean(s_final), max(s_final), min(s_final)))

if __name__ == '__main__':
    main()
```