

sheet04_prog

November 11, 2018

1 Weighted K-Means

In this exercise we will simulate finding good locations for production plants of a company in order to minimize its logistical costs. In particular, we would like to place production plants near customers so as to reduce shipping costs and delivery time.

We assume that the probability of someone being a customer is independent of its geographical location and that the overall cost of delivering products to customers is proportional to the squared Euclidean distance to the closest production plant. Under these assumptions, the K-Means algorithm is an appropriate method to find a good set of locations. Indeed, K-Means finds a spatial clustering of potential customers and the centroid of each cluster can be chosen to be the location of the plant.

Because there are potentially millions of customers, and that it is not scalable to model each customer as a data point in the K-Means procedure, we consider instead as many points as there are geographical locations, and assign to each geographical location a weight w_i corresponding to the number of inhabitants at that location. The resulting problem becomes a weighted version of K-Means where we seek to minimize the objective:

$$J(c_1, \dots, c_K) = \frac{\sum_i w_i \min_k \|x_i - c_k\|^2}{\sum_i w_i},$$

where c_k is the k th centroid, and w_i is the weight of each geographical coordinate x_i . In order to minimize this cost function, we iteratively perform the following EM computations:

- **Expectation step:** Compute the set of points associated to each centroid:

$$\forall 1 \leq k \leq K : \quad \mathcal{C}(k) \leftarrow \left\{ i : k = \arg \min_k \|x_i - c_k\|^2 \right\}$$

- **Minimization step:** Recompute the centroid as a the (weighted) mean of the associated data points:

$$\forall 1 \leq k \leq K : \quad c_k \leftarrow \frac{\sum_{i \in \mathcal{C}(k)} w_i \cdot x_i}{\sum_{i \in \mathcal{C}(k)} w_i}$$

until the objective $J(c_1, \dots, c_K)$ has converged.

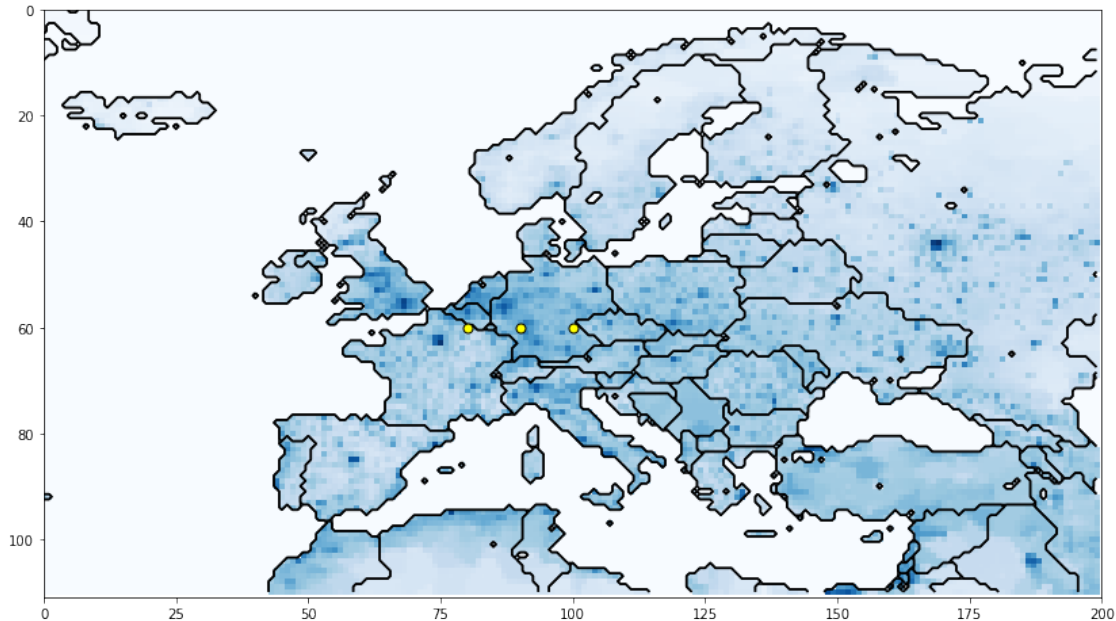
1.1 Getting started

In this exercise we will use data from <http://sedac.ciesin.columbia.edu/>, that we store in the files `data.mat` as part of the zip archive. The data contains for each geographical coordinates (latitude and longitude), the number of inhabitants and the corresponding country. Several variables and methods are provided in the file `utils.py`:

- `utils.population` A 2D array with the number of inhabitants at each latitude/longitude.
- `utils.countries` A 2D array with the country indicator at each latitude/longitude.
- `utils.nx` The number of latitudes considered.
- `utils.ny` The number of longitudes considered.
- `utils.plot(latitudes,longitudes)` Plot a list of centroids given as geographical coordinates in overlay to the population density map.

The code below plots three factories (white squares) with geographical coordinates (60,80), (60,90),(60,100) given as input.

```
In [1]: import utils
import matplotlib.pyplot as plt
%matplotlib inline
utils.plot([60,60,60],[80,90,100])
```



1.2 Initializing Weighted K-Means (15 P)

Because K-means has a non-convex objective, choosing a good initial set of centroids is important. Centroids are drawn from the following discrete probability distribution:

$$P(x,y) = \frac{1}{Z} \cdot \text{population}(x,y)$$

where Z is a normalization constant. Furthermore, to avoid identical centroids, we add a small Gaussian noise to the location of centroids, with standard deviation 0.01.

Tasks:

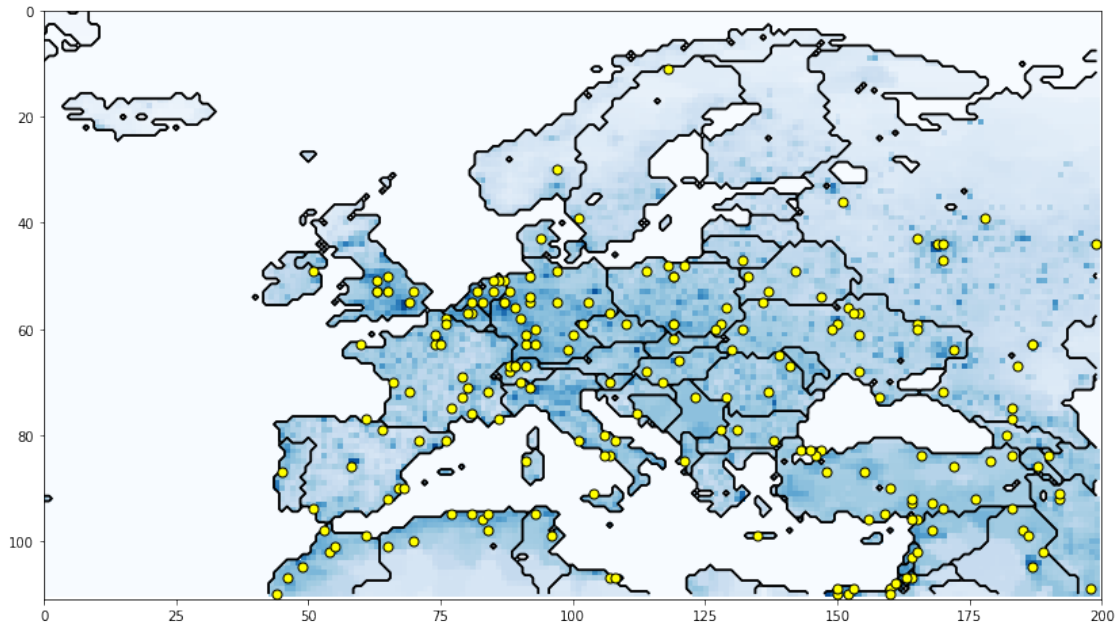
- **Implement the initialization procedure above.**
- **Run the initialization procedure for $K=200$ clusters.**
- **Visualize the centroids obtained with your initialization procedure using `utils.plot`.**

```
In [4]: # YOUR CODE HERE
        %matplotlib inline
        import numpy as np
        xys = np.array(np.meshgrid(np.arange(utils.nx), np.arange(utils.ny))).T.reshape(-1,2)

        def P(x, y):
            noise = np.random.normal(scale = 0.01, size=utils.population.shape)
            population = utils.population * (1 + noise)
            return (population/population.sum())[x,y]

        def get_centroids(K = 200):
            # randomly chose centroids, with the associated pdf
            # generated from P() above, in mind
            centroid_ids = np.random.choice(xys.shape[0], K, p=P(*xys.T))
            return xys[centroid_ids]

        centroids = get_centroids()
        utils.plot(*centroids.T)
        #utils.plot(*np.argwhere(utils.population).T)
```



1.3 Implementing Weighted K-Means (30 P)

Tasks:

- Implement the weighted K-Means algorithm as described in the introduction.
- Run the algorithm with $K=200$ centroids until convergence (stop if the objective does not improve by more than 0.01). Convergence should occur after less than 50 iterations. If it takes longer, something must be wrong.
- Print the value of the objective function at each iteration.
- Visualize the centroids at the end of the training procedure using the methods `utils.plot`.

```
In [5]: def dist(xy, c):
        # compute distance from one (x, y)-vector to all centroids
        # @return shape (amount centroids,) vector
        return np.linalg.norm(xy - c, axis=1)

        def dists(xys, c):
            # compute distances for all (x, y)-vectors to all centroids
            # @return shape (amount xypairs, amount centroids)
            return np.apply_along_axis(dist, axis=1, arr=xys, c=c)

        def assign(centroids):
            # find the index k for all the centroids which are closest
            # to all vectors
```

```

    # @return shape (111, 200)
    return dists(xys, centroids).argmin(axis=1).reshape(utils.nx, utils.ny)

def normalize_cluster(k, assoc):
    # finds all vectors associated with a centroids
    # @return total weighted sum of all vectors
    cluster_members = assoc == k
    xys = np.argwhere(cluster_members)
    ws = utils.population[cluster_members].reshape(-1,1)
    return (ws * xys).sum(axis=0) / (ws.sum() or 1)

def update_clusters(assoc, K = 200):
    # updates each cluster with a new weight
    # @return updated clusters
    ks = np.arange(0, K).reshape(K, 1)
    return np.apply_along_axis(normalize_cluster, axis=1, arr=ks, assoc=assoc)

def J(c):
    # calculate the weighted sum for all (x,y)-vectors
    is_ = utils.population > .0
    xs = np.argwhere(is_)
    ks = dists(xs, c).argmin(axis=1)
    ws = utils.population[is_]
    return (ws * ks).sum(axis=0) / utils.population.sum()

J_old = 0
iterations = 0
while iterations < 50:
    centroid_associations = assign(centroids)
    centroids = update_clusters(centroid_associations)
    J_new = J(centroids)
    print("J = {:.2f}".format(J_new))
    if abs(J_old - J_new) <= .01:
        break
    J_old = J_new
    iterations += 1
    #plt.imshow(assigned)

utils.plot(*centroids.T)

```

```

J = 101.77
J = 101.40
J = 101.13
J = 101.28
J = 101.14
J = 101.13
J = 101.41
J = 101.46

```

J = 101.22
J = 101.03
J = 100.95
J = 100.95

