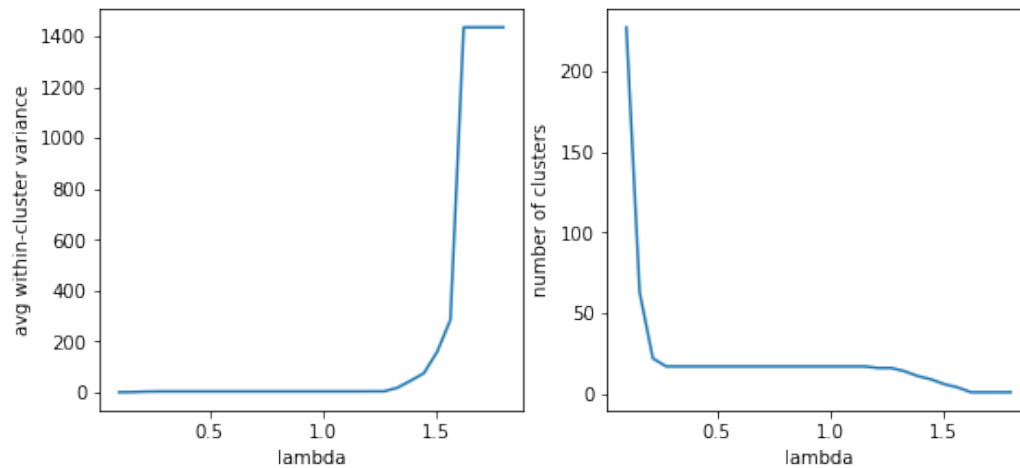


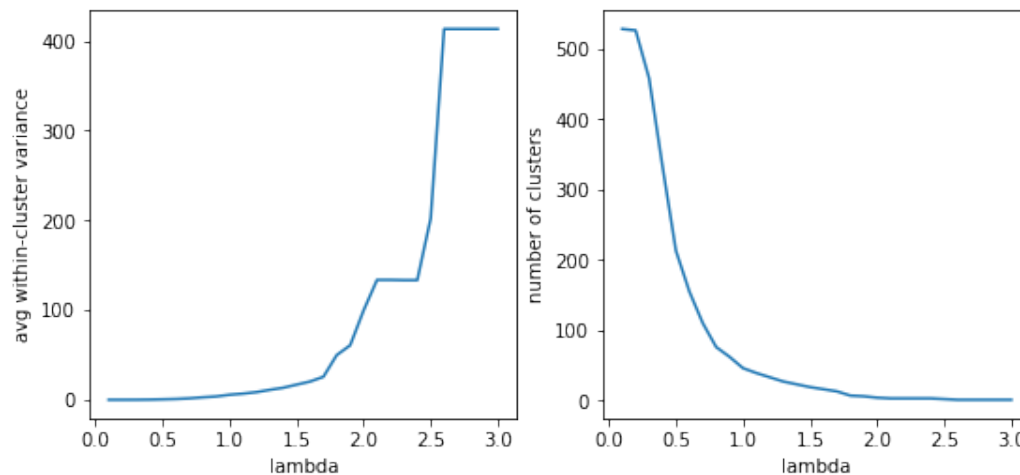
Yanbing Wang
11/18/19
K-means and clustering algorithm

Programming

Numerical data



Water treatment data



Discussion and results

The shape of the curves for both datasets look similar: there is a range for λ to yield optimal results of the cluster. The lower bound is the bending point of the λ vs. # clusters curve; the upper bound is the bending point of the λ vs. variance curve. From this observation, we see that the optimal λ range for the numerical dataset is 0.3-1.2, and for the water treatment data is 1.0-1.5.

Written

- 1) The total variance is a non-increasing function of k because:
k-means converges because the within-cluster variance monotonically decreases or stays the same in each iteration. First, the variance decreases in the reassignment step since each vector is assigned to the closet centroid, so the distance it contributes to the variance decreases. Second, it decreases in the re-computation step because the new centroid is the vector for which the variance reaches its minimum.

Analytically:

$$RSS_k(\mu_k) = \sum_{k=1}^K \sum_{x_i \in S_k} \|x_i - \mu_k\|_2^2$$
$$\frac{\partial RSS_k(\mu_k)}{\partial \mu_k} = \sum_{x_i \in S_k} 2(\mu_k - x_i)$$

$\mu_k = \frac{1}{n_k} \sum_{x_i \in S_k} x_i$, where n_k is the number of points in the cluster k .

Thus, we minimize RSS_k when the old centroid is replaced with the new centroid. RSS , the sum of the RSS_k , must then also decrease during recomputation.

- 2) See attachment

b). $x \rightarrow \psi(x)$

$$\sum_{k=1}^K \sum_{x_i \in S_k} \|\psi(x_i) - a_k\|_2^2$$

$$= \sum_k \sum_{x_i \in S_k} (\psi(x_i) - a_k)^T (\psi(x_i) - a_k)$$

$$= \sum_k \sum_{x_i \in S_k} \psi(x_i)^T \psi(x_i) + \underline{a_k^T a_k} - \underline{2 \psi(x_i)^T a_k}$$

$$a_k^T a_k = \frac{1}{|S_k|^2} (\sum_{x_i \in S_k} \psi(x_i))^T (\sum_{x_i \in S_k} \psi(x_i))$$

$$= \frac{1}{|S_k|^2} \sum_{i,j \in S_k} K(x_i, x_j)$$

$$\psi(x_i)^T a_k = \psi(x_i)^T \frac{1}{|S_k|} \sum_{j \in S_k} \psi(x_j) = \frac{1}{|S_k|} \sum_{j \in S_k} K(x_i, x_j)$$

\therefore the original objective fcn. becomes:

$$\min_{S=\{S_1, \dots, S_K\}} \sum_{k=1}^K \sum_{x_i \in S_k} \left[K(x_i, x_i) + \underbrace{\frac{1}{|S_k|^2} \sum_{i,j \in S_k} K(x_i, x_j) + \frac{2}{|S_k|} \sum_{j \in S_k} K(x_i, x_j)}_{\substack{\text{fcn. of } K(x_i, x_j), \text{ measure of} \\ \text{within-cluster variance}}} \right]$$

\downarrow
 do not change