

Chapter 10.7: Exercises

1.) (a.) Prove (10.12):

$$a.) \text{ Prove (10.12):} \\ 2 \sum_{i \in C_K} \sum_{j=1}^p (x_{ij} - \bar{x}_{ij})^2 = 2 \sum_{i \in C_K} \sum_{j=1}^p \left(\frac{|C_K| x_{ij}}{|C_K|} - \frac{\sum_{i \in C_K} x_{ij}}{|C_K|} \right)^2$$

$$= 2 \sum_{i \in C_k} \sum_{j=1}^p \left(\frac{(|C_k|-1)x_{ij} - \sum_{\ell \neq k, \ell \in C_k} x_{\ell j}}{|C_k|} \right)^2$$

$$= \alpha \sum_{i \in C_K} \sum_{j=1}^p \left(\frac{\sum_{(i' \in C_K) \neq i} (x_{ij} - x_{i'j})}{|C_K|} \right)$$

$$= \left(\frac{1}{2} \right) \left(\prod_{i=1}^k x_i^{c_{ki}} \right) \left[(x_{1j} - x_{2j})^2 + (x_{1j} - x_{3j})^2 + \dots + (x_{1j} - x_{nj})^2 + (x_{2j} - x_{3j})^2 + \dots + (x_{2j} - x_{nj})^2 + \dots + (x_{3j} - x_{nj})^2 + \dots + (x_{1j} - x_{1kij})^2 + (x_{2j} - x_{1kij})^2 + \dots + (x_{nj} - x_{1kij})^2 \right]$$

A hand-drawn graph on lined paper showing a periodic wave function. The vertical axis has two horizontal grid lines. The wave starts at a positive value, crosses the top grid line, reaches a maximum, crosses the bottom grid line, reaches a minimum, and returns to its starting value. The period of the wave is approximately 10 units of time. The amplitude of the wave is approximately 4 units.

$$= 2 \left(\frac{1}{|C_K|} \right)^2 \sum_{j=1}^{|C_K|} \left[(x_{1j} - \bar{x}_{1j})^2 + (x_{2j} - \bar{x}_{2j})^2 + \dots + (x_{ij} - \bar{x}_{ij})^2 + (x_{1j} - \bar{x}_{1Kj})^2 + \dots + (x_{ij} - \bar{x}_{1Kj})^2 \right]$$

$$(x_{ij} - x_{1ckl})^2 + \dots + (x_{ij} - x_{(ckl-1)j})^2 + \dots + (x_{1j} - x_{1ckl})^2 + \dots + (x_{1j} - x_{(ckl-1)j})^2$$

↓ BACK ↓

$$= (2) \left(\frac{1}{|C_k|} \right)^2 \sum_{j=1}^{|C_k|} (x_{ij}^2 - 2x_{ij}x_{2j} + x_{2j}^2) + (x_{ij}^2 - x_{ij} \cdot x_{3j} - \cancel{x_{2j}x_{ij} - x_{2j}x_{3j}}) + \dots + \} \}_{i=1} \\ (x_{|C_k|j}^2 - 2x_{|C_k|j} \cdot x_{ij} + x_{ij}^2) + (x_{|C_k|j}^2 - x_{|C_k|j} \cdot x_{2j} - x_{ij} \cdot x_{|C_k|j} + x_{2j}x_{ij}) + \dots \} \}_{i=1}^{|C_k|}$$

$$= 2 \left(\frac{1}{|C_k|} \right) \left(\frac{1}{|C_k|} \right) \sum_{j=1}^{|C_k|} (|C_k|(|C_k|-1)) (x_{ij}^2 + x_{2j}^2 + \dots + x_{|C_k|j}^2 - x_{ij}x_{2j} - x_{ij}x_{3j} - \dots - x_{ij}x_{|C_k|j} - \dots - x_{|C_k|j}x_{|C_k|j})$$

$$= \left(\frac{1}{|C_k|} \right) \sum_{j=1}^{|C_k|} (2) [(x_{ij}^2 - 2x_{ij}x_{2j} + x_{2j}^2) + (x_{ij}^2 - 2x_{ij}x_{3j} + x_{3j}^2) + \dots + (x_{ij}^2 - 2x_{ij}x_{|C_k|j} + x_{|C_k|j}^2) + \dots]$$

is used $(|C_k|-1)$ times, etc.

$$= \left(\frac{1}{|C_k|} \right) \sum_{j=1}^{|C_k|} (2) [(x_{ij} - x_{2j})^2 + \dots + (x_{ij} - x_{|C_k|j})^2 + \dots \} \}_{i=1}^{i-1} \quad (i-1) \text{ terms} \\ (x_{2j} - x_{3j})^2 + \dots + (x_{2j} - x_{|C_k|j})^2 + \dots \} \}_{i=2}^{i-2} \quad (i-2) \text{ terms} \\ + (x_{(|C_k|-1)j} - x_{|C_k|j})^2 \quad 1 \text{ term}$$

$$= \left(\frac{1}{|C_k|} \right) \sum_{j=1}^{|C_k|} \sum_{i, i' \in C_k} (x_{ij} - x_{i'j})^2$$

$$= \left(\frac{1}{|C_k|} \right) \sum_{i, i' \in C_k} \sum_{j=1}^{|C_k|} (x_{ij} - x_{i'j})^2$$

Chapter 10.7, Exercise 1 (continued)

(b) For any iteration of Algorithm 10.1, we begin the iteration with a current assignment.

Then:

- * Per Algorithm 10.1 Step 2, Part (a), we compute the centroid for each cluster.
- * Given, the above centroids, we then compute γ = the sum of the within-cluster variations, over all clusters.
- * Then, per Algorithm 10.1 Step 2, Part (b), we assign each observation to the cluster which is closest to it.
- * Thus, after this latest assignment, each observation will either be closer to (or the same distance from) its new cluster centroid. Thus, this either reduces (or keeps the same, if algo has converged) the quantity:

$$\sum_{k=1}^K 2 \sum_{i \in C_k} \sum_{j=1}^p (x_{ij} - \bar{x}_{kj})^2$$

- * Then, by identity (10.12) which I proved in (a), this either reduces or keeps constant the quantity:

$$\sum_{k=1}^K \frac{1}{|C_k|} \sum_{i, i' \in C_k} \sum_{j=1}^p (x_{ij} - x_{i'j})^2$$

which is precisely the objective (10.11) we wish to minimize.

Thus I have shown that K-Means clustering algorithm decreased (or keeps constant) objective (10.11) at each iteration. 