

Chap. 12

本レジュムは (12.9) 式の導出を目指すものである。教科書とは異なる表現^{*}を用いるため、まず各種前提から整理させてもらう。

※ i 番目のデータをベクトル x_i で表現します。教科書では x_{ik} のようにベクトル成分を露骨に書いていますが、抜いづらくなるため、あえてベクトル表示を使います。

2つのクラスター A, B を考える。各クラスターに所属する i 番目のデータは A に所属する場合; $x_i (i \in A)$, データ数は n_A コ

B に所属する場合; $x_i (i \in B)$, データ数は n_B コ

のように表現できる。今2つのクラスター A, B を併合して、新しいクラスターを作る、この新しいクラスター $A \cup B$ に所属するデータは

$$x_i (i \in A \cup B)$$

と書くことができる。

各クラスター内における平方和の def. は

$$A; \sum_{i \in A} \left| x_i - \sum_{j \in A} \frac{x_j}{n_A} \right|^2$$

$$B; \sum_{i \in B} \left| x_i - \sum_{j \in B} \frac{x_j}{n_B} \right|^2$$

$$A \cup B; \sum_{i \in A \cup B} \left| x_i - \sum_{j \in A \cup B} \frac{x_j}{n_A + n_B} \right|^2$$

となる。いずれの平方和も各クラスター内のデータの平均からの偏差平方和として約束されている。

ここで

$$\Delta = \sum_{i \in A \cup B} \left| x_i - \sum_{j \in A \cup B} \frac{x_j}{n_A + n_B} \right|^2 - \sum_{i \in A} \left| x_i - \sum_{j \in A} \frac{x_j}{n_A} \right|^2 - \sum_{i \in B} \left| x_i - \sum_{j \in B} \frac{x_j}{n_B} \right|^2$$

のように Δ を定義する。 Δ は merging cost と呼ばれ、この Δ を minimize するようなクラスタリング方法を「ワード法 (Ward Method)」と呼ぶ。

案は Δ は concise に表現することが可能で

$$\Delta = \frac{n_A n_B}{n_A + n_B} (m_A - m_B)^2, \text{ with } m_A = \sum_{i \in A} \frac{x_i}{n_A} \text{ and } m_B = \sum_{i \in B} \frac{x_i}{n_B}$$

となる。以下では Δ の def. から上式を導くことにする。

$$\begin{aligned} \Delta &= \sum_{i \in A \cup B} \left\{ |x_i|^2 - 2x_i \cdot \frac{\sum_{j \in A \cup B} x_j}{n_A + n_B} + \left| \frac{\sum_{j \in A \cup B} x_j}{n_A + n_B} \right|^2 \right\} \\ &= \sum_{i \in A} \left\{ |x_i|^2 - 2x_i \cdot \frac{\sum_{j \in A} x_j}{n_A} + \left| \frac{\sum_{j \in A} x_j}{n_A} \right|^2 \right\} \\ &\quad - \sum_{i \in B} \left\{ |x_i|^2 - 2x_i \cdot \frac{\sum_{j \in B} x_j}{n_B} + \left| \frac{\sum_{j \in B} x_j}{n_B} \right|^2 \right\} \end{aligned}$$

ここで、以降の計算のポイントを抑える。

① $|x_i|^2$ の項はキャンセル

② $\sum_{j \in A \cup B}$ の項は $\sum_{j \in A}$ と $\sum_{j \in B}$ に分離

③ m_A, m_B が現われるように変形

④, ⑤ を実行すると

$$\begin{aligned} \frac{\sum_{j \in A \cup B} x_j}{n_A + n_B} &= \frac{\sum_{j \in A} x_j + \sum_{j \in B} x_j}{n_A + n_B} = \frac{1}{n_A + n_B} \left(n_A \sum_{j \in A} \frac{x_j}{n_A} + n_B \sum_{j \in B} \frac{x_j}{n_B} \right) \\ &= \frac{1}{n_A + n_B} (n_A m_A + n_B m_B) \quad \dots (\star) \end{aligned}$$

① と (\star) によ、

$$\begin{aligned} \Delta &= -2(n_A + n_B) \sum_{i \in A \cup B} \frac{x_i}{n_A + n_B} \cdot \sum_{j \in A \cup B} \frac{x_j}{n_A + n_B} + (n_A + n_B) \left| \frac{\sum_{j \in A \cup B} x_j}{n_A + n_B} \right|^2 \\ &\quad + 2n_A \sum_{i \in A} \frac{x_i}{n_A} \cdot \sum_{j \in A} \frac{x_j}{n_A} - n_A \left| \sum_{j \in A} \frac{x_j}{n_A} \right|^2 \\ &\quad + 2n_B \sum_{i \in B} \frac{x_i}{n_B} \cdot \sum_{j \in B} \frac{x_j}{n_B} - n_B \left| \sum_{j \in B} \frac{x_j}{n_B} \right|^2 \quad (\because \textcircled{1}) \end{aligned}$$

$$= -2(n_A + n_B) \frac{1}{n_A + n_B} (n_A m_A + n_B m_B) \cdot \frac{1}{n_A + n_B} (n_A m_A + n_B m_B)$$

$$+ (n_A + n_B) \left| \frac{1}{n_A + n_B} (n_A m_A + n_B m_B) \right|^2$$

$$+ 2n_A m_A \cdot m_A - n_A |m_A|^2 + 2n_B m_B \cdot m_B - n_B |m_B|^2 \quad (\because (\star))$$

$$= - \frac{2}{n_A + n_B} (n_A m_A + n_B m_B)^2 + \frac{1}{n_A + n_B} (n_A m_A + n_B m_B)^2$$

$$+ n_A m_A^2 + n_B m_B^2$$

$$= \frac{- (n_A m_A + n_B m_B)^2 + (n_A + n_B) n_A m_A^2 + (n_A + n_B) n_B m_B^2}{n_A + n_B}$$

$$= \dots = \frac{n_A n_B}{n_A + n_B} (m_A - m_B)^2$$

