

K-MEANS אלגוריתם

אלגוריתם להצברה (CLUSTERING) של נקודות ב- \mathbb{R}^d ל- K צברים. מטרתו היא למצוא את מרכזי הצברים כך שסכום המרחקים הריבועיים בין כל הנקודה ומרכז הצבר אליו שוייכה יהיה הקטן ביותר.

שלבי האלגוריתם:

1. בוחרים באקראי K מרכזים μ_1, \dots, μ_k השייכים ל- \mathbb{R}^d .

2. משייכים כל נקודה למרכז הקרוב אליה ביותר:

$$c(x_i) = \operatorname{argmin}_j (||x_i - \mu_j||)$$

כאשר $c(x_i)$ משמעותו לאיזה צבר נשייך את x_i (מקבל ערכים מ-1 עד K).

3. מעדכנים כל מרכז להיות הממוצע (מרכז כובד) של כל הנקודות ששויכו אליו:

$$\mu_j = \frac{\sum_{i:c(x_i)=j} x_i}{\sum_{i:c(x_i)=j} 1}$$

המכנה: $\sum_{i:c(x_i)=j} 1$ הוא בעצם מספר הנקודות ששויכו ל- μ_j .

4. חוזרים על צעדים 2-3 עד שאין שינוי באף אחד מהמרכזים.

קשר בין K-MEANS ל-EM של תערובת גאוסיינים:

אם ב-EM נחליף את $P(h_i = m | x_i, \hat{\theta}_t)$ בפונקציה שמקבלת 1 עבור m שמביא את ההסתברות הנ"ל למקסימום ו-0 לשאר נקבל כלל שיוך קשה (HARD CLUSTERING). אם בנוסף נניח שהפריורים (c_j) שווים ושהווריאנס ידוע ושווה למטריצת היחידה $\Sigma = I$, נקבל את אלגוריתם K-MEANS. במצב כזה אנחנו בעצם ממקסמים את הפונקציה הבאה:

$$\log \left[P \left(x_i | h_i, \hat{\theta}_t \right) \right] \propto -(x_i - \hat{\mu}_t)^T \begin{matrix} \Sigma^{-1} \\ \parallel \\ I \end{matrix} (x_i - \hat{\mu}_t) = -||x_i - \hat{\mu}_t||^2$$