

## תרגול 13 : אישכול

### אלגוריתם K-Means

זהו אלגוריתם איטרטיבי אשר מחלק וקטורים נתונים ל- $K$  קבוצות, ולכל קבוצה מוצא נקודת "מרכז מסה" שמייצגת את הקבוצה כולה (נקודה מרכזית או סנטרואיד, centroid). בינתיים נניח כי מספר הקבוצות  $K$  נתון מראש.

סימונים:  $\mu_i$  הוא הסנטרואיד של קבוצה  $G_i$  ( $i = 1, \dots, K$ ).

האלגוריתם:

▪ אתחול: בחירת  $K$  נקודות מרכזיות  $\{\mu_i^{(0)}\}_{i=1}^K$ .  $t = 0$ .

▪ התהליך:

1. סיווג הנקודות הקיימות באמצעות אלגוריתם 1-NN ביחס

לסנטרואידים. כלומר, נקודה  $x$  שייכת לקבוצה  $G_i^{(t)}$  אם

$$i = \arg \min_{j=1, \dots, C} \{\|x - \mu_j^{(t)}\|\}$$

נניח כי במקרה של שוויון תמיד נבחר את הקבוצה עם האינדקס הקטן.

$$2. \text{ מציאת הסנטרואידים החדשים: } \mu_i^{(t+1)} = \frac{1}{|G_i^{(t)}|} \sum_{x \in G_i^{(t)}} x$$

כאשר  $|G_i^{(t)}|$  הוא מספר האיברים בקבוצה ה- $i$ .

$$\text{אם } |G_i^{(t)}| = 0, \mu_i^{(t+1)} = \mu_i^{(t)}$$

3.  $t \leftarrow t + 1$  וחזרה לשלב 1 עד להתכנסות ( $\mu_i^{(t+1)} \approx \mu_i^{(t)}$  לכל  $i$ ).

$$\text{האלגוריתם שואף למזער את סכום השגיאות הריבועיות: } \sum_{i=1}^K \sum_{x \in G_i} \|x - \mu_i\|^2$$

לאלגוריתם זה מובטחת התכנסות למינימום מקומי. בניסיון מתקבל שתהליך זה עמיד (*robust*) לתנאי ההתחלה, אך מומלץ לבחור את תנאי ההתחלה על פי מחשבה ושימוש במידע מקדים במידת האפשר.

## קביעת מספר הקבוצות $K$ ע"פ המידע עצמו

לעתים קרובות, לא יודעים מראש את ערכו של  $K$  ומעוניינים לקבוע  $K$  "טבעי", תוך כדי ריצת האלגוריתם.

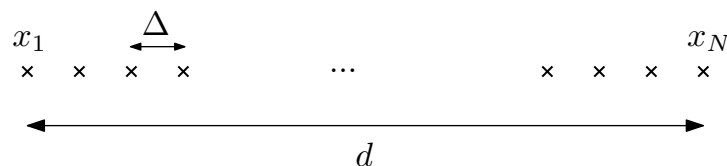
$$E(K) = \sqrt{\sum_{i=1}^K \sum_{x \in G_i} \|x - \mu_i\|^2}$$

נגדיר את שגיאת האשכול  $E(K)$  ככל שמגדילים את  $K$ , השגיאה קטנה. ניתן לראות שאם נהפוך כל נקודה לאשכול עצמאי, השגיאה תהיה אפס, אבל לא הרווחנו שום ידע. אחת השיטות לקביעת מספר קבוצות "טבעי" היא ע"י הגדלה הדרגתית של  $K$ , וחישוב את  $E(K)$  בסיום

כל שלב. עוצרים את התהליך ב- $K$  שנותן (למשל)  $1 - \frac{E(K)}{E(K-1)} < \varepsilon$ , כאשר  $\varepsilon$  הוא סף מסוים.

## תרגיל

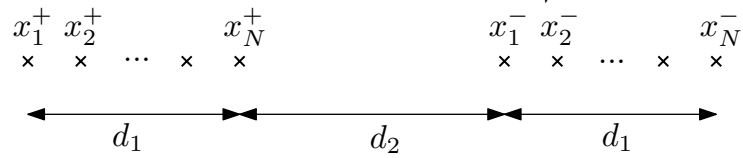
נתבונן בבעיית "האשכול" החד-מימדית הבאה:



כאשר הנקודות  $\{x_j\}_{j=1}^N$  ממוקמות באופן אחיד באינטרוול  $[0, d]$  ומספרן  $N \rightarrow \infty$  (וכמובן  $\Delta \rightarrow 0$ ).

א. הראו כי האלגוריתם K-Means עם  $K = 2$  מתכנס למינימום הגלובלי של השגיאה הריבועית מכל תנאי התחלה סביר (כלומר, המרכזים ההתחלתיים ממוקמים באינטרוול  $[0, d]$ ).

נתבונן כעת בבעיית אשכול קצת יותר אמתית:



כלומר נתונים שני אשכולות  $\{x_j^+\}_{j=1}^N$ ,  $\{x_j^-\}_{j=1}^N$ , כאשר הנקודות בכל אשכול ממוקמות באופן אחיד, ושוב נניח כי  $N \rightarrow \infty$ .

ב. היעזרו בסעיף א' על מנת להראות כי גם במקרה זה האלגוריתם K-Means עם  $K = 2$  מתכנס למינימום הגלובלי של השגיאה הריבועית מכל תנאי התחלה סביר.

ג. האם התוצאות של סעיפים א' וב' עדיין תקפות לכל תנאי התחלה של המרכזים ב- $\mathbb{R}$ ?