# Clustering

I. מבוא ואלגוריתמים בסיסיים

מוטיבציה: הרבה פעמים רוצים ללמוד ולהבין דברים מעל מבלי שיהיה לנו מורה

Unsupervised ⟸

נרצה לחקור את הדביקל שלנו ולמצוא בו משהו – נרצה מבנה

ננסה את הדביקל שלב בל גוף שנראה מסביר ונמצא על כך הבנה.

דוגמא שנרצה לענות: כמה "טיפוסים" שונים של דביקל יש לנו?

מה הקריטריון (הנתונים) דביקל?

נגדיר את הבעיה:

יהי    X = אוסף נקודות    ← נקודות

d = פונקציית מרחק R:X×X→R    ונדרוש שהיא תהיה מטריקה.

C = קלאסטרים 1-K קבוצות כך שמתקיים:

∀i,j  $C_i \cap C_j = \emptyset$   התאמך רק

$\bigcup_{i=1}^{K} C_i = X$   הנקודות נמצאות

פונקציית מטרה לפשטנית:

לרוב נרצה למינימום את סכם המרחקים בין הנקודות בתוך הקלאסטרים

$$\text{argmin} \sum_{K=1}^{K} \sum_{x_i,x_j \in C_k} d(x_i,x_j)$$
    $C$

הנחה: ארגינו ידעים את K מראש.

אבל יש פה פונקציה ישירה, פשוט נעשה ונסרב שהתחלף מה...

שאלה של John Kleinberg ב-NIPS 2002 שואל פורמלית →

של קלאסטרים ומראה ש-3 התנאים שנרצה לבקש אל תמיד מתקיימות ביחד:

– Scale Invariance → שאם נכפיל את "מחיר" הנקודה על המרחק לשנות את הפיצול

– Richness → שאם נסדר את הנקודות נמצא p שהיפוטזה הפלט תהיה זה הפל

– Consistency → אם נשנה את הנקודות בתוך ה קלאסטר (a) שלהם

המרחק שלהם בין קלאסטרים, הקלסטרו לל ישונה

<u>I</u> שאלון ותאוריה: כב"ם עינטרים: תרגול

אנו רוצים לחלק את ... Bottom-Up → מתחילים נקולטם של נקודה ונחבר
Top-Down → מתחילים מקלוסטר אחד ונחלקם וכו'.

→ <u>Hierarchical Agglomerative Clustering</u> *

Bottom-Up          איך (מתי) ארמד מין הקלוסטרים?

# Single Linkage : לפי (המינימ):

$$d(A,B) = \min_{\substack{a \in A \\ b \in B}} d(a,b)$$

# Average Linkage : לפי (הממוצע):

$$d(A,B) = \frac{\sum_{a \in A} \sum_{b \in B} d(a,b)}{|A| \cdot |B|}$$

בזמן הריצה על ... יהיה בערך $O(n^3)$

→ <u>Lloyd's Algorithm</u> (k-means) :

המטרה שלנו:          Top-Down

$$C = \underset{C}{\arg\min} \sum_{k=1}^{K} \sum_{x_i, x_j \in C_k} \|x_i - x_j\|_2 = \underset{C}{\arg\min} \sum_{k=1}^{K} \sum_{x_i \in C_k} \|x_i - \mu_k\|_2$$

                                                     ↑ הצנטרויד.

נקולטרים זה תהליך דמוי EM אבל עם מתודת hard-assignment
בשלב ה...

האלגוריתם: -נגריל של k, את (הצנטרויד $\mu_k$
נשאף ע"י גיבור k מרכז ... בתחילה.
- של k :

* נעבור לכלוטר $C_k$ את ה-... ל-$\mu_k$ הכי ...
קרוב אליו. ונמצא: $C_k = \{x: k = \arg\min \|x - \mu_k\|\}$

* נחשב מחדש את $\mu_k$

<u>שיפורים</u>: ① ... הרבה פעמים וכל פעם נבחר ... התוצאה ...
② (נבחר את החל... (subsample)
③ (נבחר ... נבחר טוב ... עם k-means ++

<u>k-means</u>: נקולטר ... של k וכל פעם את k וכך ... הכי ...
... בתוצאה ...

$$d_{ij} = \|x_i - x_j\|_2$$

data point

distance  $f$

X  Dist  W  Laplacian

K-means ← normalization ← V

eigen vectors

אל תלכו לקטוע את הקווים, לפי שזה אבר מיכה אל קווים (Hebrew handwritten notes near Dist arrow)

אפשר סתירה סלילי, מ סתירה משתנים, מ סתירה לשבצ מקבל
לכל (זפל את פ)
קווים, כמו ב $D - Q$
ה בתורה ג, אול ספל קים לבי. (Hebrew handwritten notes near W arrow)

# איך עוברים מ- Dist ל- W מתור סריב הקווים?

<u>שיטה 1:</u> Dist & <u>thresholding</u>: הקווים אשר לסל מ-N-ר כפלי 1 וכסל 0.



וכמו ס.

איך נקפל את ספל ס ?

(סרל אל הכיגותכת סל הנקודים אכילוגיות.

'לפס ס, פ סרל יש מאבדי הסרו אל קפל סאל.

<u>שיטה 2:</u> וכמל מקפל כללי, פ ל-

$$\omega_{ij} = \exp\left(\frac{-\|x_i - x_j\|^2}{2\sigma^2}\right)$$



אל ס (סרל כולל צונ אר שורך קווים.

(סל) אל ס $\Rightarrow$ ... סל ישר סעל

אל רצי מ לכלכס בסילו וכגל אר רי קפר קסר בצל וכיבל אל W.

\# איך נמצא רכיבים קשירים בגרף?

נגדיר מטריצה נוספת:



$$d_{ii} = \sum_{j} w_{ij} = \begin{array}{c}\text{סכום המשקלים}\\ \text{היוצאים מ-}i\end{array}$$

\# אפשר להכניס אל תוך הקלסטרינג חשיבה של מיזוג בין רכיבי הגרף (אולי לא קשירים)

ציור של קבוצה —



הגדרות:

A                                        B

נגדיר:

$$|A| = \sum_{i \in A} \mathbb{1}$$

$$vol(A) = \sum_{i \in A} d_{ii} = \sum_{i \in A} \sum_{j} w_{ij}$$

נשים לב שה- weak clustering אפשר להגדיר כ:

$$\underset{A,B}{argmin} \quad CUT(A,B)$$

$$\underset{\substack{i \in A \\ j \in B}}{\sum} w_{ij}$$

• יש כאן אלגוריתם שפותר את זה → פורד-פולקרסון $O(VE^2)$
עם אלגוריתמים...

מה הבעיה? outliers!

אם יש נקודה קיצונית רחוקה מהשאר, פעמים רבות נקבל אותה
כקלסטר ואת (השאר) הקלסטר...

איך נפתור את זה? נורמל רכיבי הגרף לגודל של הקלסטרים

$$\underset{A,B}{argmin} \quad CUT(A,B) \cdot \left[\frac{1}{|A|} + \frac{1}{|B|}\right]$$

— קרוב ל-1 כאשר עוצמה של שוויון של הקבוצות
וקרוב ל-0 כאשר עוצמה של אי-שוויון כי הקבוצות.

נשים לב שה-NP קשה למצוא אופטימום לקרוב ל-$\frac{1}{2}$

אפשרות אחרת להגדיר: נחלק ל-$A,B$ ביחס ל-volume:

$$\underset{A,B}{\arg\min} \ \ CUT(A,B) \cdot \left[ \frac{1}{vol(A)} + \frac{1}{vol(B)} \right]$$

וגם כאן הבעיה כמו קודם זו NP-קשה לפתור.

## <u>Graph Laplacian</u> #

נעזר על ידי מטריצה (סימטרית): $\qquad L = D - W$

$\underset{\text{laplacian}}{\nearrow} \qquad \underset{\text{degree}}{\uparrow} \qquad \underset{\text{adjacency}}{\nwarrow}$

לכל גרף יש תכונות ל-$L$ כמו: $0$ הוא ערך עצמי $L$-ו $\vec{1}$ הוא וקטור עצמי

$$\begin{bmatrix} L \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

(כאשר סכום השורה שווה $0=$ ⟸ $\lambda=0$.

נוכל, $L$ גם ערכים עצמיים.

האם ניתן לומר $PSD$ ?

⟸ $Positive \ Semi \ Definite = PSD$ ⟵ המטריצה

⟸ כל הערכים העצמיים חיוביים או שווה ל-$0$

לכל וקטור $f$, $\qquad f^T L f \ge 0$.

על אם נזיז את הטרנספורמציה $L$-ה וקטור $f$-ל

נשנה את $f$ בפחות מ-$90°$.

נוכיח שהיא $PSD$ :

$$f^T L f = f^T D f - f^T W f = \sum_i d_i f_i^2 - \sum_{i,j} f_i f_j \omega_{ij} =$$

$$= \sum_i \sum_j \omega_{ij} f_i^2 - \sum_{i,j} f_i f_j \omega_{ij} =$$

$$= \frac{1}{2} \left( \sum_i \sum_j \omega_{ij} f_i^2 - 2 \sum_i \sum_j f_i f_j \omega_{ij} + \sum_j \sum_i \omega_{ij} f_j^2 \right) =$$

$$= \frac{1}{2} \sum_{i,j} \underbrace{\omega_{ij}}_{\ge 0} \cdot \underbrace{(f_i - f_j)^2}_{\ge 0} \ge 0$$

⟸ ומכאן נובע $\qquad 0 = \lambda_1 \le \lambda_2 \le \ \dots \ \le \lambda$

$\qquad\qquad\qquad\qquad \underset{\mathbb{1}}{\uparrow} \quad \underset{Fiedler \ Vector}{\uparrow}$

נושא: <u>Spectral Clustering II</u>

\# מ"י את הערכים העצמיים מחזיר השלב? #

אם הגרף הוא אכן מורכב מ רכיבים קשירים, שהם נקודות אות אחד,
על אזור הקשור הנתון בכל אחד בתוכונות העצמיים, ערכים קצרה
אחד ואפסים בכל הארכי, וזה מה יהיה עם הערך העצמי 0.
תמיד קיים ווקטור עצמי של הלפלסיאן קיים וקטור שמתאים לערך 0 =>

מספר הערכים העצמיים העצמיים ל"ע 0 הוא ווא מספר הקלסטרים.

↑
אם כל הים הוא ל"ע 0 אחד מחובר
אבל קשיר ל-0 זה הם כי ל-ש הקשורים
כל הקשורים בנקודה.

\# אבל מה זה הדרורים של לינון של לפלסיאן? <u>Normalized Laplacian</u>

1. Shi & Malik , 2000 :          $L_n = D^{-1} L$

זה נרמול על ד שוורה שומר את זה הסימטריות.

(נקודות לחלונקן = נקודה = Random Walk )

2. Ng , Jordan & Weiss , 2002 :

$$L_{sym} = D^{-\frac{1}{2}} \cdot L \cdot D^{-\frac{1}{2}} = I - D^{-\frac{1}{2}} W D^{-\frac{1}{2}}$$

מחלוקת ל עוצמה בשורה  ↑ מחלוקת ל שורה בשורה הערכה
ההגדרה הזו נאור סימטרי, ולכן PSD

| קטגות של ה-notes |
| זו |

\# עכשיו, האלגוריתם בשלבו:

נתון: $X$ ,  $x_i \in \mathbb{R}^d$

חשב 1. מטריצת $S_{n \times n}$ המרחקים : $\forall i j \quad S_{ij} = \|x_i - x_j\|$

חשב 2. מטריצת $W_{n \times n}$ הדמיון שערך : $W_{ij} = \exp\left(\frac{-S_{ij}^2}{2\sigma^2}\right)$

כאשר $\sigma \overset{\Delta}{=}$ ith percentile of $S$

חשב 3. מטריצת $D_{n \times n}$ ודרגות בדרך : $D_{ii} = \sum_j W_{ij}$
$D_{ij} = 0 , i \neq j$

חשב 4. מטריצת $L_{sym}$ $L_{sym} = I - D^{-\frac{1}{2}} W D^{-\frac{1}{2}}$ :

5. חשב את k הערכים העצמיים הקטנים הקטנים:
$U = \begin{bmatrix} | & | & & | \\ u_1 & u_2 & \dots & u_k \\ | & | & & | \end{bmatrix}$

כאשר k-את לפי ה-eigen-gap.

6. חשב $T$ תהיה נרמול של הדורים של $U$ : $T_i = \frac{u_i}{\|u_i\|}$ → כל הקוטרים של הקוטרים של צורה הזה
→ ולחרב הקוטרים אולחרבשאה כיון יהי לוקים.

7. הרץ k-means על להקל השורות של $T$.

t-distribution Stochastic Neighbour Embedding t-SNE III

\# נתון לנו פסט נקודות $X$. נקודות $X_i \in \mathbb{R}^p$

ונרצה למצוא embedding של $X$ ל־$Y$: $Y_i \in \mathbb{R}^d$  $d << p$

כך שהמרחקים בין זוגי נקודות ישמרו, כלומר  $d'_{ij} \approx d_{ij}$

$$d_{ij} = \|x_i - x_j\|$$

כך זו כמו שרצה על שנה את ד ננסה ע"ד החישובים של אלה את זה (ריאלי)

נרצה שנקודות "שכנות" של אחד מהם ... על נקודות שמות.

\# נסתכל על זה כמו כין ענין נסמן:

$$P_{j|i} = \frac{\exp(-d_{ij}^2/2\sigma_i)}{\sum_{k \neq i} \exp(-d_{ik}^2/2\sigma_i)}$$

ההסתברות אשר הנק' $j$... $i$

וכל' נרצה שתמרחק על ההתפלגות... יהיה קטן, ולמדוד את זה

נוכל ע"י מרחק $KL$.

$$\underset{Y}{argmin}\ C = \sum_i KL(P_i \| Q_i)$$

$X$ ־ב $i$-ה נקודת של ההתפלגות   ...   $Y$ ־ב $i$-ה נקודת של ההתפלגות

$$Q_{j|i} = \frac{\exp(-d_{ij}'^2/2\sigma_i')}{\sum_{k \neq i} \exp(-d_{ik}'^2/2\sigma_i')}$$

גדרות לעזרה:  perplexity: $2^{H(P_i)}$ ← זה פסק מיצג את כמות השכנים $H(P_i) = \sum_j P_{ij} \log_2(P_{ij})$

$KL$:  $KL(P_i\|Q_i) = \sum_j P_{j|i} \log \frac{P_{j|i}}{Q_{j|i}}$

$$\Rightarrow \underset{Y}{argmin} \sum_i KL(P_i\|Q_i) = \sum_i \sum_j P_{j|i} \log \frac{P_{j|i}}{Q_{j|i}}$$

← נשים לב כי פה יש ... אסימטריה: אם שתי נקודות אשר רחוקות $Q_{j|i}$

של $P_{j|i}$ יהיה מאוד נמוך (מרחק קטן וקרוב), זה יהיה

מסבר קטן ... => מחיר ... => זה לא ... מסימה.

\# נמצא $Y$ ... אשר ימזער את המרחק את ההתפלגות נגזור את $C$ ונעשה הורדה:

$$\frac{\partial C}{\partial y_i} = 2 \cdot \sum_j (P_{j|i} - Q_{j|i} + P_{i|j} - Q_{i|j})(y_i - y_j)$$

כלל:  $Y^{(t)} = Y^{(t-1)} + \eta \frac{\partial C}{\partial Y} + \alpha(t)[Y^{(t-1)} - Y^{(t-2)}]$

SNE

מה הקנייה ב-SNE שהלגן מאוד קשים ?

‎1. מתמך- outliers

‎2. במידה של (הנית מן מקדרים את נקודות, כן (יד קשר
‎את ‎t הקרוב את לו שהיה על underestimation כן נקודות-
‎קמעת, או שיהיה ‎over-estimation כן נקודות רחוק-

את הנקרה השנייה מתיר ה-t-SNE מבעב t-distribution (המפלנת t)
המפלנת ‎t היו מן המפלאות נורמלי יכל עב כלב כbd.
‎את ‎נוראb הקנן של (המפלנת t :

$$q_{j|i} = \frac{(1 + \|y_i - y_j\|^2)^{-1}}{\sum_k (1 + \|y_i - y_k\|^2)^{-1}}$$

בעל, מרחקים להמח של המטרו נורמיי עם נקודת i, הם,
‎הפרש של המפלוויי נשיטו נקודת i ל-j.

$$c = \underset{Y}{\text{argmin}} \ KL(P \| Q)$$

‎? על קבל מנ i עם נקודת נורמיי עולמית (המפלא ←

$$\frac{\partial c}{\partial y_i} = 2 \sum_j (P_{j|i} - Q_{j|i}) \cdot (Y_i - Y_j) \cdot [term]$$

‎ובכה ‎רב נילא מתח נשן נקודת ל-outliers.