

מבוא לבינה מלאכותית – תרגול 7

נושאים:

- (סיום אלגוריתמי האשכול - Fuzzy C-means , GMM)
 - **זיהוי קהילות (Community Detection) – אשכול היררכי, Girvan-Newman, Louvain**
-

זיהוי קהילות (Community Detection):

למעשה זה המשך די ישיר לנושא של אשכול. הפעם, נתון לנו גרף כלשהו, אפשר עם משקלים על הקשתות שנסמלים את מידת הדמיון שבין קודקודים שונים אך אפשר גם ללא משקלים כלשהו, והמשימה שלנו היא לחלק את קודקודי הגרף לקהילות. האלגוריתמים השונים מחליטים על סמך מה בדיקן לבצע את החלוקה לקהילות (למה קודקוד אחד יהיה באותה קהילה עם קודקוד אחר, ולא יהיה בקהילה עם קודקוד שלישי), אך הרעיון המרכזי של האלגוריתמים הוא שני קודקודים שקרובים ומקשרים חזק אחד לשני אמורים להיות דומים יותר זה לזה מאשר שני קודקודים שהם רחוקים זה מזה או שהקשר שביניהם נעשה על ידי קשתות במשקלים נמוכים, ולכןעדיף לשירות את הזוג הראשון לאוותה קהילה על פני שיוור של הזוג השני.

Community Detection באמצעות אשכול היררכי:

כמו בשיעור בעבר, בהינתן ממד דמיון בין כל זוג קודקודים, אפשר לבצע חלוקה לקהילות באמצעות אשכול היררכי.

ספציפית לתרגול זהה, אתן דוגמה לממד הדמיון הבא:

בhinaten GRAPH (לא ממושך) עם מטריצת שכנות A , כמות המסלולים שבין קודקוד i לקודקוד j שאורכם k היא $\sum_{k=1}^{\infty} A^k$. לכן, כמות המסלולים הכלולות שבין i ו- j בכל אורך היא $\sum_{i,j} \sum_{k=1}^{\infty} A^k$. הכמות זו מתבדרת, ולכן במקומות זה ניקח α קטן מספיק ונחשב $\sum_{i,j} \alpha^k (A^k)$, כך שהפעם זו יתכנס. זה בעצם טור הנדסי של מטריצות, שמתכנס למטריצה הבאה $(\alpha A - I)^{-1} = D$. המטריצה זו תהיה מטריצת הדמיון שבין כל זוג קודקודים, ובאמצעותה ניתן לבצע אשכול היררכי ולהלך את הגרף.

אלגוריתם Girvan-Newman

האלגוריתם זהה מתמקד בהסירה איטרטיבית של קשתות שמהוות "גשרים" שבין קהילות. שימוש לבשוא, כמו זה שלמעלה, מקבל גרפף שאינו בהכרח ממושקל. מחשבים לכל קשת מدد שקובע כמה היא "מרכזית", ובכך איטרטיבית הקשתות ה"מרכזיות" ביותר יורדות מהגרף עד שלא נשארות כללה. את התהיליך אפשר לתאר בתור בניה מלמעלה למטה של dendrogram (שראינו בשיעור שעבר בנושא של אשכול היררכי, אבל שם הבנייה הייתה *bottom up*), ואת הבחירה הסופית של החלוקה לקהילות אפשר ללקחת בתור שלב מסוים בעז, כמו באשכול היררכי.

בצורה פורמלית יותר:

אלגוריתם:

כל עוד יש קשתות בגרף:

1. לכל קשת בגרף, מחשבים – edge centrality

$$c(e) = \sum_{u,v} \frac{\sigma(u, v|e)}{\sigma(u, v)}$$

כאשר e קשת, u, v זוג קודקודים, (u, v) σ מספר המסלולים המינימליים (ביחס לאורך, ללא משקלים) שבין u ל- v בגרף, $(u|v, e)$ σ מספר המסלולים המינימליים שבין u ל- v בגרף, שעוביים דרך e .

2. מורידים את הקשת בעלת ה- edge centrality הגבוהה ביותר. רכיבי הקשרות שנותרו הם ה- communities יחסית לשלב זהה, ומהם בונים את ה-dendrogram.
3. חוזרים ל- 1 עבור הגרף הנותר.

חסרון בולט של שני האלגוריתמים שראינו עד כה הוא שאין בהם פונקציית מטרת מוגדרת. האלגוריתם השני טוב יותר מהראשון, כי הוא לא צפוי להיתקל בבעיה שהאשכול היררכי bottom up נוטה להיתקל בה (בנייה אשכול גדול שבסך פעם מצטרפת אליו נקודה), אך מצד שני הוא יקר יחסית – הסיבוכיות שלו היא בסדר גודל של מספר הקשתות (מספר הורדות הקשתות) כפול מספר הקודקודים בריבוע (עלות חישוב ה- edge centrality לכל קשת).

אלגוריתם Louvain:

סתם לידע כללי, אלגוריתם זה הומצא ע"י Blondel ושותפים שלו, והוא נקרא על שם האוניברסיטה בה למדנו (בבלגיה).

זהו אלגוריתם חמדן שמחפש מינימום לפונקציה מוגדרת, שנקראה modularity ומוגדרת כ:

$$L = -\frac{1}{2m} \sum_{i,j} \left[A_{ij} - \frac{k_i k_j}{2m} \right] C(i,j)$$

כאשר:

- m סכום כל הקשתות בגרף (כולל משקל אם יש).
- k_i הדרגה (המושקלת) של קודקוד i .
- $C(i,j)$ היא למשה הדلتא של Kronecker לגבי האם הקודקודים j,i , שייכים לאותה קהילה (1 אם כן, 0 אחרת).

אינטרואיציה: נגיד שיש לנו שני קודקודים שאינם באותו רכיב קשירות. אם נבחר לשים אותם באותה קהילה $(C(i,j) = 1)$, נוסיף סתם $\frac{1}{4m^2} k_i k_j$ לערך הפונקציה, וזה מרחיק אותנו ממחיפוש המינימום. אם j,i מחוברים זה לזה בקשת עם משקל A_{ij} גדול מאוד, אז כדאי לנו ששניהם יהיו באותה קהילה, כי אז יקטין את ערך הפונקציה, וכך נעדיף שבמקרה זה $C(i,j) = 1$.

הערה – לאלגוריתם קיימת גרסה עם קבוע שירתי γ שכפול את $\frac{k_i k_j}{2m}$. התפקיד שלו: אם γ קטן, אז נעדיף לקבץ הרבה הרכבה קודקודים יחד ולקבל מעט קהילות גדולות יחסית, כי $\left(A_{ij} - \frac{k_i k_j}{2m} \right) \gamma$ יהיה חיובי וזה יקטין את ערך הפונקציה. אם γ גדול, בדיק ההפר – נקבל מספר קהילות גדול ומעט קודקודים בכל קהילה.

איך מוצאים את חלוקה שמהווה מינימום לפונקציית ה- modularity ?

אלגוריתם:

1. כל קודקוד מאוחתל בתור שיר לקהילה משל עצמו.
2. בכל שלב, עברור על כל הקהילות הקיימות.
לקהילה i , נבדוק עבור כל קהילה אחרת j מה ההפרש בפונקציית ה- modularity כאשר "נציא" את הקהילה i ו"נכenis" את הקודקודים שלה לקהילה j . עבור הקהילה שמתקבל בה הרוחח הגדל ביותר עבורנו (הפונקציה תרד בהכי הרבה), אם אכן קיים רוחח זהה, נאחד את הקהילות ונמשיך בתהליך עבור הקהילה הבאה. אם אין איחוד רוחח, לא מחדדים.
3. כאשר אין עוד צעדי איחוד שמקטין את ה- modularity, מסייםים.

אלגוריתם זה מהיר יותר מ predecessors, עלות החישוב שלו עבור גרף עם n קודקודים היא ($n^2 \log^2 n$).
חו"ץ מזה, כאן באמת יש פונקציית מטרה שימושים בה, וזה עוזר להבין מה קורה באלגוריתם
ולשנות בו.

שימוש לבשהאלגוריתם הזה תלוי בסדר הריצה על הקהילות, שכן הוא יכול להתכנס לערכי מינימום
שוניים עבור הרצות שונות.