

# עבודה סמינריונית בנושא

# Spectral Graph Theory

בהנחיית: פרופ' אוריאל פייגה

מגישה: אילנה ברזניק

אוקטובר 2012

## תוכן עניינים :

- 3.....מבוא
- 4.....מהו גרף
- 5.....הצגת גרף בעזרת מטריצה
- 8.....אינטרפרטציה חדשה למושג הערך העצמי במטריצות
- 9.....Spectral Graph Theory
- 16.....Power Method של מטריצה בעזרת ה
- 18.....הרעיון של קליינברג
- 19.....שאלות ומקורות מוסמכים
- 19.....ניתוח הרשת המתקשרת
- 20.....בניית גרף מרוכז באינטרנט
- 22.....חישוב רכזים ודפים איכותיים
- 23.....האלגוריתם האיטראטיבי
- 26.....רשימה ביבליוגרפית

## מבוא

מטרתי בעבודה זו היא להציג את אחת השיטות מהשנים האחרונות לדרוג הדפים באינטרנט. השיטה הוצגה על ידי קליינברג במאמרו JON M. KEINBERG/Authoritative Sources in a Hyperlinked Environment/*Journal of the ACM, Vol. 46, No. 5, September 1999, pp.*

604-632. טענתו המרכזית של קליינברג היא כי ניתן ללמוד הרבה על תוכנה של הסביבה המתוקשרת בעזרת ניתוח המבנה שלה. תנאי הכרחי ללמידה זו הוא הימצאותם של הכלים המתאימים. בעבודתו המוצגת במאמר מציג קליינברג כלים אלגוריתמיים שמטרתם לנתח את מבנה האינטרנט ואת הקשרים בו ככלי לקבלת אינפורמציה על תוכנו.

באמצע שנות ה-90 היה Yahoo האתר הראשון שדרג את הדפים באינטרנט בצורה ידנית. כיום, האינטרנט כבר גדול מידי ולכן שיטה זו הופכת להיות יקרה ולא אפקטיבית. השאלה הגדולה היא כיצד לעשות את הדרוג בצורה אוטומטית? אפשרות אחת היא להמציא תוכנה הקוראת את דפי האינטרנט ומדרגת אותם לפי איכותם. כיום תוכנה כזו טרם הומצאה עקב המורכבות שבעניין. לא נפסלת האפשרות כי תוכנה כזו תומצא בעתיד. אפשרות אחרת היא לחשוב על מדד שייתן אינדיקציה על איכות הדפים באינטרנט ותוך כדי שהדבר יוכל להיעשות בצורה יחסית אוטומטית על ידי מחשב.

עבודתו של קליינברג היא מהראשונות המציעות איך לעשות את דרוג הדפים לא רק לפי תוכנם אלא על ידי בחינת דפי האינטרנט המצביעים עליהם, שהם מהווים כמעין "ממליצים". לצורך בדיקת איכותו של דף אינטרנט לא מספיק לבדוק את כמות הממליצים. שכן, אם זהו הקריטריון לבחינת חשיבותו של דף האינטרנט ניתן לבנות בצורה מלאכותית הרבה דפי אינטרנט שיצביעו על דף זה. דרך נוספת שניתן לדרג בה את הדפים היא לפי מספר הפעמים בהם מופיעה מילת החיפוש בדף. לדוגמא, אנו רוצים אינפורמציה באינטרנט על חלב. נבצע חיפוש באינטרנט ומנוע החיפוש יחזיר לנו דפים המדורגים לפי מספר הפעמים שמופיעה בהם המילה "חלב". גם שיטה זו אינה אפקטיבית במידה מספקת שכן אם העוסקים ב"קידום אתרים" מגלים אותה הם יוכלו לקחת דף אינטרנט שברצונם לקדם ולקדם אותו על ידי הוספת המילה המתאימה מספר רב של פעמים בתוכן הדף. בהמשך העבודה נראה מה מציע לעשות קליינברג.

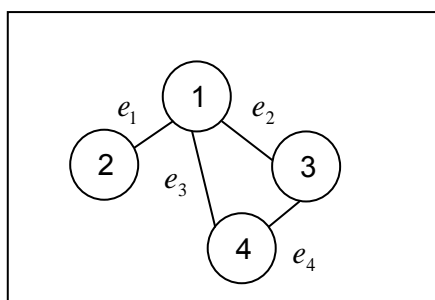
## מהו גרף

גרף הינו ייצוג מופשט של קבוצה של אובייקטים, כאשר כל תת-קבוצה של אובייקטים בקבוצה עשויים להיות מקושרים זה לזה. האובייקטים הניתנים לקישור מכונים קודקודים או צמתים (באנגלית: vertex), וקבוצת הקודקודים מסומנת באות  $V$ .

הקישורים בין הקודקודים מכונים צלעות או קשתות (באנגלית: edge), וקבוצת הצלעות מסומנת באות  $E$ . מתקיים כי קבוצת הצלעות מקיימת:  $E \subseteq V \times V$ , כלומר: כל צלע הינה זוג הקודקודים, אותם היא מקשרת.

גרף, אשר קבוצת הקודקודים שלו היא  $V$  וקבוצת הצלעות שלו היא  $E$  מסומן באופן הבא:  $G = (V, E)$ .

לדוגמא, נתבונן בגרף הבא:



**תורת הגרפים הספקטלית** היא תורה המאפשרת להתאים לכל גרף מטריצה וללמוד מהערכים העצמיים של המטריצה את תכונות הגרף.

---

קובץ זה נועד אך ורק לשימוש האישי של מורים למתמטיקה, פיזיקה, כימיה וביולוגיה ולהוראה בכיתותיהם. אין לעשות שימוש כלשהו בקובץ זה לכל מטרה אחרת, ובכלל זה: שימוש מסחרי, פרסום באתר אחר (למעט אתר בית הספר בו מלמד המורה), העמדה לרשות הציבור או הפצה בדרך אחרת כלשהי של קובץ זה או חלק ממנו.

## הצגת גרף בעזרת מטריצה

ניתן לייצג גרף בדוגמא על ידי כמה מטריצות:

### 1. מטריצת קודקודים-קשתות:

$$\begin{array}{c} e_1 \quad e_2 \quad e_3 \quad e_4 \\ \begin{array}{c} 1 \\ 2 \\ 3 \\ 4 \end{array} \begin{pmatrix} 1 & 1 & 1 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 1 \end{pmatrix} \end{array}$$

מטריצה זו מוגדרת על הקודקודים והקשתות על ידי ההגדרה הבאה:

$$a_{i,j} = \begin{cases} 1 & \text{אם קודקוד } j \text{ עובר דרך צלע } i \\ 0 & \text{אחרת} \end{cases}$$

### 2. מטריצת הסמיכויות:

מוגדרת על הקודקודים.

$$a_{i,j} = \begin{cases} 1 & \text{אם } v_i \text{ סמוך ל } v_j \\ 0 & \text{אחרת} \end{cases}$$

בדוגמא שלנו המטריצה היא:

$$\begin{array}{c} 1 \quad 2 \quad 3 \quad 4 \\ \begin{array}{c} 1 \\ 2 \\ 3 \end{array} \begin{pmatrix} 0 & 1 & 1 & 1 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 \\ 1 & 0 & 1 & 0 \end{pmatrix} \end{array}$$

קובץ זה נועד אך ורק לשימוש האישי של מורים למתמטיקה, פיזיקה, כימיה וביולוגיה ולהוראה בכיתותיהם. אין לעשות שימוש כלשהו בקובץ זה לכל מטרה אחרת, ובכלל זה: שימוש מסחרי, פרסום באתר אחר (למעט אתר בית הספר בו מלמד המורה), העמדה לרשות הציבור או הפצה בדרך אחרת כלשהי של קובץ זה או חלק ממנו.

**3. מטריצת הלפלטיאן:**

מטריצה זו מוגדרת על הקודקודים בלבד.

$$l_{i,j} = \begin{cases} \deg(v_i) & \text{אם } i = j \\ -1 & \text{אם } i \neq j \text{ וגם } v_i \text{ סמוך ל } v_j \\ 0 & \text{אחרת} \end{cases}$$

הצורה המנורמלת של מטריצה זו היא :

$$l_{i,j} = \begin{cases} 1 & \text{אם } i = j \text{ וגם } v_i \neq 0 \\ -\frac{1}{\sqrt{\deg(v_i) \cdot \deg(v_j)}} & \text{אם } i \neq j \text{ וגם } v_i \text{ סמוך ל } v_j \\ 0 & \text{אחרת} \end{cases}$$

כאשר הגרף הוא רגולרי, כלומר הדרגות של כל הקודקודים שוות, שתי המטריצות שקולות.

נבנה את המטריצות עבור הגרף שבדוגמא.

המטריצה הרגילה:

$$\begin{matrix} & \begin{matrix} 1 & 2 & 3 & 4 \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \end{matrix} & \begin{pmatrix} 3 & -1 & -1 & -1 \\ -1 & 1 & 0 & 0 \\ -1 & 0 & 2 & -1 \\ -1 & 0 & -1 & 2 \end{pmatrix} \end{matrix}$$

המטריצה המנורמלת:

$$\begin{matrix} & \begin{matrix} 1 & 2 & 3 & 4 \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \end{matrix} & \begin{pmatrix} 1 & -\frac{1}{\sqrt{3}} & -\frac{1}{\sqrt{6}} & -\frac{1}{\sqrt{6}} \\ -\frac{1}{\sqrt{3}} & 1 & 0 & 0 \\ \frac{1}{\sqrt{6}} & 0 & \frac{1}{2} & -\frac{1}{\sqrt{6}} \\ -\frac{1}{\sqrt{6}} & 0 & -\frac{1}{\sqrt{6}} & \frac{1}{2} \end{pmatrix} \end{matrix}$$

קובץ זה לועד אף ורק לשימוש האישי של מורים למתמטיקה, פיזיקה, כימיה וביולוגיה ולהוראה בכיתותיהם. אין לעשות שימוש כלשהו בקובץ זה לכל מטרה אחרת, ובכלל זה: שימוש מסחרי, פרסום באתר אחר (למעט אתר בית הספר בו מומד המורה), העמדה לרשות הציבור או הפצה בדרך אחרת כלשהי של קובץ זה או חלק ממנו.

2

3

4

נשים לב שבגרפים לא מכוונים המטריצות הנ"ל הן סימטריות.

צורה נוספת להסתכל על מטריצת הלפלסיאן היא להסתכל על מטריצה המתקבלת מההפרש בין מטריצת הדרגות למטריצת הסמיכויות של גרף  $G$ . על מטריצת הדרגות היא מטריצה סקלרית שבה מופיעות על האלכסון דרגות הקודקודים.

לדוגמא, בגרף הקודם:

$$\begin{pmatrix} 3 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 2 \end{pmatrix} - \begin{pmatrix} 0 & 1 & 1 & 1 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 \\ 1 & 0 & 1 & 0 \end{pmatrix} = \begin{pmatrix} 3 & -1 & -1 & -1 \\ -1 & 1 & 0 & 0 \\ -1 & 0 & 2 & -1 \\ -1 & 0 & -1 & 2 \end{pmatrix} \quad \text{סימון:}$$

$$L_G = D_G - A_G$$

הערכים העצמיים והוקטורים העצמיים של מטריצות אלו מספקים אינפורמציה על הגרף. עבור גרפים רגולריים, המידע שניתן להסיק מאחת ההצגות של המטריצה זהה למידע שניתן להסיק מההצגות האחרות. לעומת זאת, בגרפים שאינם רגולריים, בחירת ההצגה היא דבר המשנה מאוד.

הסבר לכך שבגרף רגולרי הייצוג לא משנה: ניתן לעשות טרנספורמציה פשוטה בין הערכים העצמיים של מטריצת הסמיכויות לערכים העצמיים של מטריצת הלפלסיאן.

במקרה של גרף  $d$  רגולרי הטרנספורמציה היא:  $\lambda \rightarrow d - \lambda$

$$Lap = d \cdot I - \begin{matrix} \Downarrow & \Downarrow \end{matrix}$$

הע"ע מוכפל ב 1 -      נוסף  $d$  לע"ע

הו"ע לא משתנה      הו"ע לא משתנה

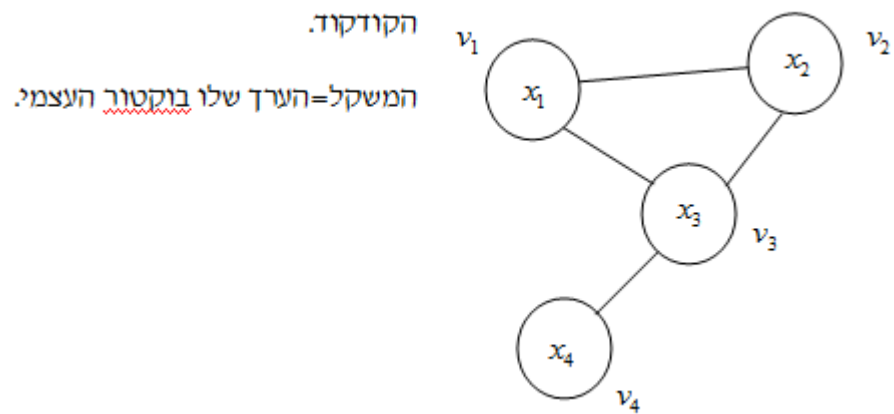
הערה: לאורך כל העבודה נשתמש במשפט מאלגברה לינארית שאם המטריצה היא ממשית וסימטרית אז הערכים העצמיים שלה ממשיים, והווקטורים העצמיים מהווים בסיס ל  $R^n$ . כלומר הם פורסים ובלתי תלויים.

קובץ זה נועד אך ורק לשימוש האישי של מורים למתמטיקה, פיזיקה, כימיה וביולוגיה ולהוראה בכיתותיהם. אין לעשות שימוש כלשהו בקובץ זה לכל מטרה אחרת, ובכלל זה: שימוש מסחרי, פרסום באתר אחר (למעט אתר בית הספר בו מלמד המורה), העמדה לרשות הציבור או הפצה בדרך אחרת כלשהי של קובץ זה או חלק ממנו.

## אינטרפרטציה חדשה למושג הערך העצמי במטריצות

אנו רגילים מאלגברה לינארית כי עבור מטריצה  $A$ ,  $x$  הוא וקטור עצמי עם ערך עצמי  $\lambda$ , אם מתקיים  $Ax = \lambda x$ .

נסתכל על אינטרפרטציה חדשה למושג הערך העצמי,



$$\begin{matrix} & 1 & 2 & 3 & 4 \\ \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \end{matrix} & \left( \begin{array}{cccc} & & & \\ & & & \\ & & & \\ & & & \end{array} \right) & \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix} = \lambda \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix} \\ & & & \Downarrow & \\ & & & \text{וקטור עצמי} & 
 \end{matrix}$$

$\lambda x_1 =$  סכום המשקלים של הקודקוד  $v_1$ .

לדוגמא: מסכמים את המשקל של השכנים (המחוברים אליו בקשת) של  $v_3$  ומקבלים  $\lambda \cdot v_3$ .

קובץ זה נועד אך ורק לשימוש האישי של מורים למתמטיקה, פיזיקה, כימיה וביולוגיה ולהוראה בכיתותיהם. אין לעשות שימוש כלשהו בקובץ זה לכל מטרה אחרת, ובכלל זה: שימוש מסחרי, פרסום באתר אחר (למעט אתר בית הספר בו מלמד המורה), העמדה לרשות הציבור או הפצה בדרך אחרת כלשהי של קובץ זה או חלק ממנו.

## רקע תיאורטי Spectral Graph Theory

משפט 1: 0 הוא תמיד ערך עצמי של מטריצת הלפלסיאן .

הוכחה:  $(1,1,\dots,1,1)$  הוא תמיד וקטור עצמי כי כל רכיב בוקטור  $Ax$  הוא למעשה סכום האיברים בכל שורה במטריצה  $A$  . הדבר קורה מכיוון ש  $x$  הוא וקטור האחדים. סכום זה הוא בדיוק 0 כי הוא מורכב מ:

הדרגה של כל קודקוד  $+(-1)^{m^*}$  הקודקודים הסמוכים לו  $=0$ .

משפט 2: מספר הפעמים שבהם 0 הוא ערך עצמי של מטריצת הלפלסיאן הוא כמספר רכיבי הקשירות של הגרף.

הוכחה: נניח שלגרף  $G, K$  רכיבי קשירות בגודל  $N_1, \dots, N_k$  . נחלק את המטריצה לפי רכיבי קשירות. נסדר את המטריצה כך ש  $N_1$  השורות הראשונות של  $L$  יהיו שייכות לרכיב הראשון.  $N_2$  השורות הבאות יהיו שייכות לרכיב השני וכך הלאה.... נקבל מטריצה מהצורה:

$$L = \begin{pmatrix} L_{1,1} & 0 & \dots & 0 \\ 0 & L_{2,2} & \dots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & \dots & L_{k,k} \end{pmatrix}$$

לפי משפט באלגברה לינארית: אם מחליפים את סדר השורות ובאותו האופן את סדר העמודות הערכים העצמיים לא משתנים כי פשוט שינינו את סדר הקודקודים במטריצה. המטריצה  $L_{i,i}$  היא תת-מטריצה  $N_i \times N_i$  המייצגת את מטריצת הלפלסיאן של הרכיב ה- $i$  של  $G$ .

נגדיר את הוקטור  $x_i$  כוקטור עם קואורדינטות 1 בכל הקואורדינאטות המתאימות לרכיב ה- $i$  ו 0 בשאר הקואורדינאטות. נעשה זאת עבור  $i=1, \dots, k$  .

כל הוקטורים  $x_1, x_2, \dots$  הם בת"ל כי המכפלה הפנימית של כל שניים מהם היא 0 מה שאומר שהם ניצבים. לכן יש ריבוי לערך העצמי כי יש לערך עצמי זה כמה וקטורים בלתי תלויים.

קובץ זה נועד אך ורק לשימושם האישי של מרצים למתמטיקה פיזיקה, כימיה וביולוגיה ולהוראה בכיתותיהם. אין לעשות שימוש כלשהו בקובץ זה לכל מטרה אחרת, ובכלל זה: שימוש מסחרי, פרסום באתר אחר (למעט אתר בית הספר בו מלמד המורה), העמדה לרשות הציבור או הפצה בדרך אחרת כלשהי של קובץ זה או חלק ממנו.

נסתכל על וקטור מהצורה :

ניתן לראות כי :  $Ly=0y$ .

אינטואיציה: התת גרף המיוצג על ידי  $L_{1,1}$  הוא בעל רכיב קשירות אחד ולכן הוקטור  $(1,1,\dots,1,1)$  הוא וקטור עצמי שלו עם ערך עצמי 0.

התרומה של שאר הרכיבים  $N_{2,2}, \dots, N_{k,k}$  היא 0 כי הגדרנו 0 בשאר הקואורדינטות. לכן קיבלנו כי  $(0, \dots, 0, 1, \dots, 1)$  הוא וקטור עצמי עם ערך עצמי 0.

גודל רכיב הקשירות הראשון.

באותו האופן, ניתן להגדיר עוד  $K-1$  וקטורים עצמיים נוספים, על ידי שניקח 1 ברכיבים של אותו רכיב קשירות ו 0 בשאר.

בצורה זו נקבל  $K$  בחירות בלתי תלויות של הוקטור  $y$ .

ולכן קיבלנו בסיס לתת-מרחב המתאים לערך העצמי 0 של מטריצת הלפלסיאן.

משפט 3: במטריצת הסמיכויות הערך העצמי המקסימלי בגרף  $d$  רגולרי הוא  $d$ .

הוכחה: גרף  $d$  רגולרי הוא גרף שבו לכל קודקוד יש  $d$  שכנים. כלומר בכל שורה של מטריצת

הסמיכויות יש  $d$  אחדים. מכיוון שהמטריצה סימטרית גם בכל עמודה יש  $d$  אחדים. נניח  $\lambda_1$  הוא הערך העצמי מקסימלי.

הכיוון  $\lambda_1 \geq d$  :

נסתכל על וקטור האחדים.  $A(1, \dots, 1) = (d, \dots, d) = d(1, \dots, 1)$ . מכאן  $x = (1, \dots, 1)$  עם ערך עצמי  $d$  ולכן הערך העצמי המקסימלי מקיים  $\lambda_1 \geq d$ .

הכיוון  $\lambda_1 \leq d$  :

נניח  $y = (y_1, \dots, y_n)$  וקטור עצמי ונניח  $y_i \geq y_j, \forall i$ . נרצה להראות שהערך העצמי קטן או שווה  $d$ . נפעיל על הוקטור  $y$  את המטריצה  $A$ .  $Ay = z$ . מכיוון שהמטריצה  $A$  היא מטריצת הסמיכויות, משמעות הפעולה היא כי מסכמים את משקלי הקודקודים הסמוכים ל  $y_i$  וזה יהיה הערך של  $z_i$  וכו'.

נניח  $Ay = z = (z_1, \dots, z_n)$ .  $\sum_{i|(1,i) \in E} y_i = dy_1$ . כלומר  $z_1 = \sum_{i|(1,i) \in E} y_i \leq \sum_{i|(1,i) \in E} y_i = dy_1$  ולכן  $\frac{z_1}{y_1} \leq d$ .

אי שיווין זה יהיה נכון לכל קואורדינטה  $i$  ולכן הערך העצמי יהיה קטן או שווה ל  $d$ .

משני כיוונים אלו נובע כי הערך העצמי המקסימלי הוא בדיוק  $d$ .

---

קובץ זה נועד אך ורק לשימוש האישי של מורים למתמטיקה, פיזיקה, כימיה וביולוגיה ולהוראה בכיתותיהם. אין לעשות שימוש כלשהו בקובץ זה לכל מטרה אחרת, ובכלל זה: שימוש מסחרי, פרסום באתר אחר (למעט אתר בית הספר בו מלמד המורה), העמדה לרשות הציבור או הפצה בדרך אחרת כלשהי של קובץ זה או חלק ממנו.

משפט 4: במטריצת הסמיכויות הע"ע המקסימלי: א. לא יותר גדול מהדרגה המקסימלית.

ב. לא יותר קטן מהדרגה הממוצעת.

הוכחת 4 א': גרף לא רגולרי הוא גרף שבו לא לכל קודקודים יש את אותה הדרגה.

נרצה להוכיח  $\lambda_1 \leq d_{\max}$  כאשר  $\lambda_1$  הוא הערך העצמי המקסימלי.

נניח  $y = (y_1, \dots, y_n)$  הוא וקטור עצמי ונניח  $y_i \geq 0, \forall i$ . נרצה להראות שהערך העצמי קטן או

שווה ל  $d_{\max}$ .

נפעיל על הוקטור  $y$  את המטריצה  $A$ .  $Ay = z$ . מכיוון שהמטריצה  $A$  היא מטריצת הסמיכויות,

משמעות הפעולה היא כי מסכמים את משקלי הקודקודים הסמוכים ל  $y_i$  וזה יהיה  $z_i$  וכו'...

נניח,  $Ay = z = (z_1, \dots, z_n)$ .

$$\frac{z_1}{x_1} \leq d_{\max} \quad \text{ולכן} \quad z_1 \leq d_{\max} \cdot y_1 \quad \text{כלומר} \quad z_1 = \sum_{i|(1,i) \in E} y_i \leq \sum_{i|(1,i) \in E} y_i \leq d_{\max} \cdot y_1$$

אי שוויון זה נכון לכל קואורדינאטה  $i$  ולכן הערך העצמי יהיה קטן או שווה  $d_{\max}$ .

הוכחת 4 ב': נסמן את הדרגה הממוצעת על ידי  $\bar{d} = \frac{\sum_{i=1}^n d_i}{n}$  כאשר  $d_i$  היא דרגת הקודקודים.

נרצה להוכיח כי  $\lambda_1 \geq \bar{d}$  כאשר  $\lambda_1$  נוכיח זאת בשני חלקים שמהם תנבע הטענה.

**חלק 1'**:  $\lambda_1 \geq y^T Ay$ . הטענה נכונה לכל וקטור יחידה  $\|y\| = 1$ .

**חלק 2'**: נרצה למצוא וקטור יחידה  $\|y\| = 1$  שעבורו מתקיים  $y^T Ay \geq \bar{d}$ .

**הוכחת חלק 1'**: נניח כי  $x_1, \dots, x_n$  הם הוקטורים העצמיים ו  $\lambda_1, \dots, \lambda_n$  הם הערכים העצמיים

מתאימים להם.  $x_1, \dots, x_n$  מהווים בסיס ל  $R^n$  ונניח שכל הוקטורים העצמיים הם מנורמלים

$$\|x_1\|^2 = \|x_2\|^2 = \dots = \|x_n\|^2 = 1 \quad \text{כלומר}$$

נסתכל על הביטוי  $y^T Ay$  ונבדוק מהו ערכו עבור וקטור יחידה  $y$ . ניתן להציג את  $y$  כצירוף

לינארי של אברי הבסיס  $x_1, \dots, x_n$ :  $y = \alpha_1 x_1 + \dots + \alpha_n x_n$ .  $\|y\| = 1$  ולכן לפי משפט פיתגורס ה- $n$

$$\sum_{i=1}^n \alpha_i^2 = 1$$

---

קובץ זה נועד אך ורק לשימוש האישי של מורים למתמטיקה, פיזיקה, כימיה וביולוגיה ולהוראה בכיתותיהם. אין לעשות שימוש כלשהו בקובץ זה לכל מטרה אחרת, ובכלל זה: שימוש מסחרי, פרסום באתר אחר (למעט אתר בית הספר בו מלמד המורה), העמדה לרשות הציבור או הפצה בדרך אחרת כלשהי של קובץ זה או חלק ממנו.

$$y^T Ay = (\alpha_1 x_1^T + \dots + \alpha_n x_n^T) \cdot A(\alpha_1 x_1 + \dots + \alpha_n x_n) = (\alpha_1 x_1^T + \dots + \alpha_n x_n^T) \cdot (\alpha_1 \lambda_1 x_1 + \dots + \alpha_n \lambda_n x_n) = \alpha_1^2 \lambda_1 + \dots + \alpha_n^2 \lambda_n$$

הסבר: הוקטורים העצמיים ניצבים כי  $x_i \cdot x_j = 0$  עבור  $i \neq j$ .

$$\lambda_1 \geq y^T Ay \text{ ולכן קיבלנו כי } \alpha_1^2 \lambda_1 + \dots + \alpha_n^2 \lambda_n \leq \alpha_1^2 \lambda_1 + \dots + \alpha_n^2 \lambda_1 = \lambda_1 (\alpha_1^2 + \dots + \alpha_n^2) = \lambda_1$$

$$y = \frac{1}{\sqrt{n}} (1, \dots, 1) \text{ כלומר } y \text{ להיות וקטור האחדים המנורמל.}$$

A היא מטריצת הסמיכויות. לכן מספר האחדים בה הוא כסכום דרגות הקודקודים.

נחשב את ערך הביטוי  $y^T Ay$ .

$$y^T Ay = \frac{1}{\sqrt{n}} (1, \dots, 1) \cdot A \frac{1}{\sqrt{n}} \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} = \frac{1}{\sqrt{n}} \cdot \frac{1}{\sqrt{n}} (1, \dots, 1) \cdot A \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} = \frac{1}{n} \cdot \text{סכום דרגות הדוקדוקים}$$

הסבר:

$$(1, \dots, 1) \cdot A = (d_1, \dots, d_n) \text{ כאשר } d_i \text{ היא דרגת הקודקוד ה-} i \text{ בגרף.}$$

$$(d_1, \dots, d_n) \cdot \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} = d_1 + \dots + d_n = \text{סכום דרגת הקודקודים}$$

הערה: הערך הכי גדול של הביטוי  $y^T Ay$  הוא הערך העצמי המקסימלי וזה קורה כאשר נשים במקום  $y$  את הוקטור העצמי המתאים לו.

**משפט 5:** בגרף דו-צדדי: גרף שקודקודיו מחולקים לשתי קבוצות זרות U, V כך שכל קשת מחברת קודקוד מ-U עם קודקוד מ-V, הערכים העצמיים באים בזוגות. כלומר אם  $\lambda$  ערך עצמי אז גם  $-\lambda$  – ערך עצמי.

**הוכחה:** נניח קודקודי הקבוצה U הם:  $p_1, \dots, p_t$ .

וקודקודי הקבוצה V הם:  $p_{t+1}, \dots, p_n$ .

טענה: אם  $(x_1, \dots, x_t, x_{t+1}, \dots, x_n)$  וקטור עצמי עם ערך עצמי  $\lambda$

אז  $(-x_1, \dots, -x_t, x_{t+1}, \dots, x_n)$  וקטור עצמי עם ערך עצמי  $-\lambda$ .

הוכחה: נראה שהדבר מתקיים עבור נציג מהקבוצה U,  $p_1$  ועבור נציג מהקבוצה V,  $p_{t+1}$ .

$$\text{מההנחה, } \lambda x_1 = \sum_{p_i}^i x_i \text{ כאשר } p_i \text{ מחובר בקשת ל } p_1.$$

קובץ זה נועד אך ורק לשימוש האישי של מורים למתמטיקה, פיזיקה, כימיה וביולוגיה ולהוראה בכיתותיהם. אין לעשות שימוש כלשהו בקובץ זה לכל מטרה אחרת, ובכלל זה: שימוש מסחרי, פרסום באתר אחר (למעט אתר בית הספר בו מלמד המורה), העמדה לרשות הציבור או הפצה בדרך אחרת כלשהי של קובץ זה או חלק ממנו.

$$\text{ולכן גם } -\lambda \cdot (-x_1) = \sum_{p_i}^i x_i \cdot p_i$$

כל משקלי הקודקודים מצד ימין של השוויון שייכים לקודקודים מהקבוצה V ולכן משקלם לא משתנה.

$$\text{בנוסף, } \lambda x_{p+1} = \sum_{p_j}^j x_j \cdot p_j \text{ כאשר } p_j \text{ מחובר בקשת ל } p_i$$

$$\text{ולכן גם } -\lambda x_{p+1} = \sum_{p_j}^j -x_j \cdot p_j \text{ כאשר } p_j \text{ מחובר בקשת ל } p_i$$

כל משקלי הקודקודים מצד ימין של השוויון שייכים לקודקודים מהקבוצה U ולכן משקלם מחליף סימן.

משפט 6: הערכים העצמיים של מטריצת הלפלאן הם גדולים שווים ל 0.

הלפלאן היא מטריצה מוגדרת חיובית (Positive-definite matrix). כלומר סימטרית וללא ערכים עצמיים שליליים.

הוכחה: בגרף רגולרי קל לראות זאת כי:

$$\text{מטריצת הסמיכויות } Lap = d \cdot I -$$



במטריצת הסמיכויות אין ערך עצמי הגדול מ d ולכן במטריצת הלפלאן אין ערכים עצמיים שליליים.

בגרף לא רגולרי,

הראינו כי  $\lambda_n \leq x^T Ax \leq \lambda_1$  כאשר  $\lambda_1$  הוא הערך העצמי המקסימלי ו  $\lambda_n$  הוא הערך העצמי המינימלי.

נראה כי לכל וקטור  $x$ ,  $x^T Lx \geq 0$  כאשר L היא מטריצת הלפלאן. מכאן ינבע כי כל הערכים העצמיים גדולים או שווים ל 0.

ניזכר כי מטריצת הלפלאן היא מהצורה:

$$\text{כאשר } d_i \text{ היא דרגת הקודקוד ה } i \cdot \begin{pmatrix} d_1 & \dots & \\ \vdots & \ddots & \vdots \\ & \dots & d_n \end{pmatrix}$$

בכל מקום אחר, יש -1 אם ישנה קשת המחברת בין הקודקודים ו 0 אחרת.

קובץ זה נועד אך ורק לשימוש האישי של מורים למתמטיקה, פיזיקה, כימיה וביולוגיה ולהוראה בכיתותיהם. אין לעשות שימוש כלשהו בקובץ זה לכל מטרה אחרת, ובכלל זה: שימוש מסחרי, פרסום באתר אחר (למעט אתר בית הספר בו מלמד המורה), העמדה לרשות הציבור או הפצה בדרך אחרת כלשהי של קובץ זה או חלק ממנו.



מכיוון שהסכום מורכב ממחבורים שהם גדולים או שווים ל 0 הוא גם כן גדול או שווה ל 0.  
מכאן, למטריצת הלפלסיאן אין ערכים עצמיים שליליים.

משפט 7:  $AA^T$  היא מטריצה סימטרית כאשר  $A = (a_{ij})$ ,  $A^T = (a_{ji})$ .

הוכחה: ניזכר בהגדרת כפל מטריצות:

נניח,  $A = (a_{ij})$ ,  $B = (b_{ij})$ , מטריצות  $n \times n$ ,  $C = AB$ .

$$C_{ik} = \sum_{j=1}^n a_{ij} b_{jk} \quad \text{אזי:}$$

$$C_{ki} = \sum_{j=1}^n a_{kj} b_{ji}$$

נרצה שיתקיים:  $C_{ik} = C_{ki}$  כאשר:  $(b_{ij}) = (a_{ji})$ . נראה שזה אכן מתקיים.

$$C_{ik} = \sum_{j=1}^n a_{ij} b_{jk} = \sum_{j=1}^n a_{ij} a_{kj}$$

$$C_{ki} = \sum_{j=1}^n a_{kj} b_{ji} = \sum_{j=1}^n a_{kj} a_{ij}$$

ומכאן  $C_{ik} = C_{ki}$  ולכן  $C = AA^T$  היא מטריצה סימטרית כנדרש.

באותה צורה ניתן להראות כי  $A^T A$  היא מטריצה סימטרית.

---

קובץ זה נועד אך ורק לשימוש האישי של מורים למתמטיקה, פיזיקה, כימיה וביולוגיה ולהוראה בכיתותיהם. אין לעשות שימוש כלשהו בקובץ זה לכל מטרה אחרת, ובכלל זה: שימוש מסחרי, פרסום באתר אחר (למעט אתר בית הספר בו מלמד המורה), העמדה לרשות הציבור או הפצה בדרך אחרת כלשהי של קובץ זה או חלק ממנו.

## מציאת ערך עצמי מקסימלי של מטריצה בעזרת ה Power Method :

**שיטה 1:** בעזרת ליכסון. מציאת הפולינום האופייני והשורשים שלו. בעזרת שיטה זו מוצאים את כל הערכים העצמיים.

**שיטה 2: (The Power Method).** אם מעוניינים רק בערך העצמי המקסימלי, ישנה צורה יותר מהירה מהשיטה הקודמת. לא מובטח כי השיטה תעבוד בוודאות, אבל היא תעבוד בהסתברות גבוהה. כלומר זה עובד באופן מעשי. בעזרת שיטה זו מוצאים גם את הוקטור העצמי. נראה כיצד היא עובדת:

$$(A)(v) = \lambda v$$

ניקח וקטור שרירותי כלשהו  $r$ . נפעיל עליו את המטריצה  $A$  ונקבל:  $Ar = r_1$ . כעת, על מנת להימנע מערכים גדולים בקואורדינטות, ננרמל את הוקטור  $r_1$ .

$$r_1^* = \frac{r_1}{\|r_1\|}$$

נמשיך הלאה ונקבל:  $Ar = r_1$

$$Ar_1^* = r_2$$

$$Ar_2^* = r_3 \quad \text{כאשר} \quad r_n^* = \frac{r_n}{\|r_n\|}$$

נניח הוקטורים העצמיים הם:  $v_1, v_2, \dots, v_n$  ונניח שהם מסודרים בסדר עולה. מכיוון שהם מהווים בסיס למרחב הוקטורי  $R^n$ ,

כל וקטור ניתן להצגה כצרוף לינארי שלהם. בפרט, הוקטור  $r$ .

$$r = \alpha_1 v_1 + \alpha_2 v_2 + \dots + \alpha_n v_n \quad \text{ולכן:}$$

---

קובץ זה נועד אך ורק לשימוש האישי של מורים למתמטיקה, פיזיקה, כימיה וביולוגיה ולהוראה בכיתותיהם. אין לעשות שימוש כלשהו בקובץ זה לכל מטרה אחרת, ובכלל זה: שימוש מסחרי, פרסום באתר אחר (למעט אתר בית הספר בו מלמד המורה), העמדה לרשות הציבור או הפצה בדרך אחרת כלשהי של קובץ זה או חלק ממנו.

$$Ar = A(\alpha_1 v_1 + \alpha_2 v_2 + \dots + \alpha_n v_n) = \alpha_1 Av_1 + \alpha_2 Av_2 + \dots + \alpha_n Av_n$$

$$r_1 = \alpha_1 \lambda_1 v_1 + \alpha_2 \lambda_2 v_2 + \dots$$

$$r_2 = \alpha_1 \lambda_1^2 v_1 + \alpha_2 \lambda_2^2 v_2 + \dots \quad \text{נמשיך כך ונקבל:}$$

$r_k = \alpha_1 \lambda_1^k v_1 + \alpha_2 \lambda_2^k v_2 + \dots$  מתכנס לוקטור המתאים לערך העצמי הכי גדול אם נחזור על התהליך מספיק פעמים.

אם  $\lambda_1$  הכי גדול אז כל שאר הרכיבים נהיים זניחים. מכאן מוצאים את הוקטור העצמי.

ככל שחוזרים על התהליך יותר פעמים, מתקרבים יותר לוקטור העצמי.

אחרי שמצאנו את הוקטור העצמי נציב אותו במשוואה ההתחלתית ונמצא את הערך העצמי.

$$A(v) = \lambda v$$

הערות:

• כאשר בוחרים באקראי את הוקטור  $r$ , חשוב שנבחר אותו כך שלא יהיה ניצב לוקטור העצמי

המתאים לערך העצמי הגדול ביותר  $v_1$ . מכיוון ש  $\alpha_i = \langle r, v_i \rangle$ , אם  $r$  ניצב ל  $v_1$  אז

$$\alpha_1 = \langle r, v_1 \rangle = 0$$

וההתכנסות תהיה לערך העצמי הבא בתור מבחינת גודל. הסיכוי שנעשה בחירה לא מתאימה של הוקטור  $r$  הוא מאוד קטן.

• כמה מהר נתכנס לוקטור העצמי? השאלה היא מהם שאר הערכים העצמיים. אם הערכים העצמיים השייכים לאותם וקטור עצמי הם רחוקים, ההתכנסות תהיה מהירה. ואם לא, זה ייקח יותר זמן.

• בעיה: בגרף דו-צדדי. בין הערכים העצמיים של מטריצת הסמיכויות מופיעים בזוגות  $d, -d$ .

זו אכן בעיה כי אם נעלה בחזקה זוגית, הערך המוחלט שלהם יהיה אותו הדבר.

---

קובץ זה נועד אך ורק לשימוש האישי של מורים למתמטיקה, פיזיקה, כימיה וביולוגיה ולהוראה בכיתותיהם. אין לעשות שימוש כלשהו בקובץ זה לכל מטרה אחרת, ובכלל זה: שימוש מסחרי, פרסום באתר אחר (למעט אתר בית הספר בו מלמד המורה), העמדה לרשות הציבור או הפצה בדרך אחרת כלשהי של קובץ זה או חלק ממנו.

## הרעיון של קליינברג

הרעיון של קליינברג הוא להחשיב כדפים ממליצים רק את הדפים שהם איכותיים בעצמם. אך לשם כך שוב נשאלת השאלה מהו דף ממליץ איכותי? אחת האפשרויות היא להחשיב כממליץ איכותי דף ששאר הדפים שהוא ממליץ עליהם הם איכותיים גם כן. למעשה נכנסנו לתהליך הבנוי משלבים שחוזרים על עצמם כאשר בכל שלב עלינו לבדוק את איכות דף האינטרנט. תהליך זה נקרא תהליך איטראטיבי ( Iterative ). תהליך שבו חוזרים על אותה פעולה שוב ושוב כאשר כל פעולה בודדת נקראת איטראציה. לכן התהליך שקליינברג מציע הוא תהליך איטראטיבי.

קליינברג בעבודתו מתייחס לשני סוגי דפים.

Authority = דף איכותי שנקרא לו דף סמכות. דף בעל איכות טובה עם תוכן מעניין על פי נושא חיפוש מסויים.

Hub = דף רכז. דף שממנו קל להגיע לדפים איכותיים, הוא מהווה ריכוז של הרבה הצבעות כלומר הוא מצביע על הרבה אתרים טובים. לכן הרכזים עוזרים לגלות דפים איכותיים. לדוגמא, אם אנו מחפשים דפים בעלי איכות טובה עם אינפורמציה על כלבים אז דף רכז יכול להיות פורום של מגדלי כלבים. דף רכז הוא מעין צומת מרכזית שכדאי לעבור דרכה על מנת להגיע למקומות חשובים.

מפאת גודלה של הרשת, לא היינו רוצים להפעיל את אלגוריתם הדרוג על כל הרשת כולה. לכן בדרך כלל יש שלב איתחול הבוחר חלק מהרשת. במאמר לא נאמר בדיוק איך עושים זאת אך בדרך כלל שלב זה מורכב משני חלקים. החלק הראשון הוא סינון הדפים שהכילו את מילת החיפוש בעצמה. החלק השני הוא הרחבת הקבוצה זו על ידי הוספת הדפים שמצביעים עליהם דפים מהחלק הראשון. נרצה לדמות את חלק האינטרנט שבחרנו לגרף דו-צדדי, גרף שקודקודיו מחולקים לשתי קבוצות, דפים איכותיים ודפי רכז. כאשר ישנן קישורים בין שתי קבוצות של דפים אלו אך לא בתוך הקבוצות עצמן. ניתן להסתכל על הדפים כעל קודקודים ועל הקישורים כעל קשתות מכוונות. האינטואיציה לרעיון של קליינברג היא כי אם מופיע באינטרנט מבנה של גרף דו-צדדי צפוף מעין זה, כנראה שמסתתרים כאן דפים בעלי תוכן איכותי. לצורך אישוש ההשערה אפשר פשוט להסתכל מדגמית על כמה מהם ולראות אם אלו דפים איכותיים או לא.

---

קובץ זה נועד אך ורק לשימוש האישי של מורים למתמטיקה, פיזיקה, כימיה וביולוגיה ולהוראה בכיתותיהם. אין לעשות שימוש כלשהו בקובץ זה לכל מטרה אחרת, ובכלל זה: שימוש מסחרי, פרסום באתר אחר (למעט אתר בית הספר בו מלמד המורה), העמדה לרשות הציבור או הפצה בדרך אחרת כלשהי של קובץ זה או חלק ממנו.

במקרה וחברה לקידום אתרים מגלה את השיטה היא יכולה לבנות מבנה מסוג זה סביב דף אינטרנט שברצונה לקדם. לכל תת גרף כזה נוכל להתאים מטריצה הנקראת מטריצת הסמיכויות. בעזרת הווקטורים העצמיים של מטריצה זו נוכל ליצור דרוג לקבוצת הדפים המיוצגים על ידי הגרף. המאמר מתמודד עם השאלה איך בכלל צריך לגשת לעניין דרוג הדפים. עדיין לא מדובר על איך להפוך את האלגוריתם ליעיל, אלא קודם צריך להבין היטב את השאלה ואת הדרישות מהתשובה בצורה ברורה.

### שאלות ומקורות מוסמכים

שאלתה היא המשפט שהמשתמש כותב בשורת מנוע החיפוש. ישנם שלשה סוגי שאלות. סוג א': נושא חיפוש ספציפי. זו שאלה ספציפית לדוגמא: מהו מזג האוויר מחר בירושלים? סוג ב': נושא חיפוש רחב. לדוגמא: מצא מידע אודות התכנות בשפת ג'אווה. סוג ג': בקשה למציאת דפים דומים. לדוגמא: מצא דפים דומים ל `java.sun.com`. נתמקד בשני הסוגים הראשונים. בכל אחד מהם, הקשיים בחיפוש הם שונים. בחיפוש ספציפי, יש מעט מאוד דפים שיכלו את מילות המפתח המבוקשות ולכן ייתכן שיהיה קשה למצוא את האינפורמציה המבוקשת. לעומת זאת, בחיפוש רחב יהיו אלפים רבים של דפים שיתאימו להגדרת החיפוש. לכן על מנת לאפשר חיפוש אפקטיבי, חייבת להימצא דרך לסנן את הדפים המוסמכים ביותר מתוך המגוון הרב הקיים. המושג "דף איכותי" קשור לשאלת נושא החיפוש הרחב. מושג זה יהווה את מקור ההתמקדות המרכזי שלנו במאמר. המכשול המרכזי בעניין הוא כיצד נקבע אם דף אינטרנט מסויים הוא דף איכותי עבור שאלתה מסויימת או לא?

### ניתוח הרשת המתוקשרת

ניתוח הרשת המתוקשרת יעזור לנו להתמודד עם הקשיים שהעלנו עד כה. למעשה, הקישורים באינטרנט טומנים בתוכם שיפוט אנושי סמוי והטענה שלנו היא שזהו בדיוק השיפוט אותו אנו צריכים כדי להגדיר מהו "דף איכותי". לדוגמא, נניח יוצרו של דף  $p$  באינטרנט כולל בו קישור לדף  $q$ , משמעות הדבר היא כי דף  $q$  חשוב לו במידה מסויימת ולכן הוא מהווה עבורו דף איכותי. יותר מכך, הקשרים נותנים לנו את ההזדמנות למצוא דפים איכותיים פוטנציאליים אך ורק על ידי הדפים המצביעים עליהם. העניין מאפשר לנו לעקוף את הבעיה שהועלתה מקודם שכן הרבה דפים אינם מספקים די אינפורמציה על עצמם. ישנם מספר חסרונות בשימוש בקישורים לצורך מטרה זו.

- א. קישורים באינטרנט נוצרו מסיבות שונות ולא דווקא לצורך הקניית איכות לדף.
- ב. ישנם קישורים שמטרתם היחידה היא ניווט באינטרנט וקישורים לצורך פרסומות.
- ב. הקושי למצוא את האיזון המתאים בין רלוונטיות ופופולאריות. כאשר שני מושגים אלו תורמים לרעיון האיכות שאנו מדברים עליו. בעבודה זו, נציע מודל המבוסס על קשרים לצורך הקניית המושג דף איכותי ונראה כיצד הוא מוביל לשיטה המזהה דפים רלוונטיים איכותיים

---

קובץ זה נועד אך ורק לשימוש האישי של מורים למתמטיקה, פיזיקה, כימיה וביולוגיה ולהוראה בכיתותיהם. אין לעשות שימוש כלשהו בקובץ זה לכל מטרה אחרת, ובכלל זה: שימוש מסחרי, פרסום באתר אחר (למעט אתר בית הספר בו מלמד המורה), העמדה לרשות הציבור או הפצה בדרך אחרת כלשהי של קובץ זה או חלק ממנו.

באינטרנט עבור שאילתא כלשהי. המודל מבוסס על הקשר הקיים בין מקורות הסמכות לנושא מסויים לבין אותם דפים המקושרים להרבה מקורות סמכות קרובים. נקרא לדפים מהסוג השני אלו-רכזים. נראה כי סוג מסויים של הקבלה טבעית זו קיימת בין רכזים לדפים איכותיים קיימת בגרף המוגדר על ידי מבנה הרשת המתקשרת. ננצל זאת על מנת לפתח אלגוריתם שיזהה את הדפים משני הסוגים בו זמנית. האלגוריתם מופעל על תת גרף מרוכז של האינטרנט, אותו נבנה מהפלט שנקבל ממנוע החיפוש. בניית תת גרף זה תייצר אוסף קטן של דפים שיכילו ככל הנראה את הדפים הרלוונטיים ביותר עבור החיפוש.

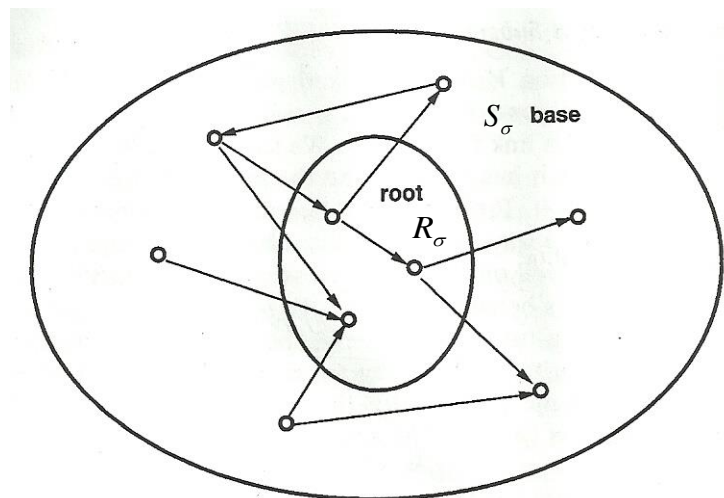
### בניית גרף מרוכז באינטרנט

נוכל להסתכל על אוסף  $V$  של דפים מקושרים באינטרנט כעל גרף מכוון.  $G=(V,E)$  כך שהקודקודים יהיו הדפים וכל קשת מכוונת  $(p,q) \in E$  מראה על קיומו של קישור מ  $p$  ל  $q$ . נגדיר את דרגת היציאה של קודקוד  $p$  להיות מספר הקודקודים שיש לו קישורים אליהם ואת דרגת הכניסה של קודקוד  $p$  להיות מספר הקודקודים שיש להם קישורים אליו. מתוך הגרף  $G$  נוכל לבדוד אזורים קטנים או תתי גרפים באופן הבא: אם  $W \subseteq V$  היא תת קבוצה של דפים, נסמן ב  $G[W]$  את הגרף המושרה מ  $W$ . קודקודיו יהיו הדפים ב  $W$  והקשתות שלו יהיו כל הקישורים בין הדפים ב  $W$ . נניח כי נתונה שאילתא בעלת נושא חיפוש רחב המצויינת בעזרת מחרוזת  $\sigma$ . נרצה לזהות את הדפים המוסמכים על ידי ניתוח מבנה הקישור. ראשית, אנו חייבים לזהות את תת הגרף באינטרנט עליו יפעל האלגוריתם. מטרתנו היא למקד את המאמץ החישובי בדפים הרלוונטיים. לכן, לדוגמא, נוכל להגביל את הניתוח שלנו בקבוצה  $Q_\sigma$  של כל הדפים המכילים את המחרוזת  $\sigma$ . אך יש לכך שני חסרונות חשובים. ראשית, קבוצה זו יכולה להכיל מעל מיליון דפים ולכן תדרוש עלות חישובית רבה. שנית, קודם לכן ציינו כי חלק מהמקורות המוסמכים החשובים ביותר יכולים להיות לא מוכללים בקבוצה זו. באופן אידיאלי, נרצה להתמקד באוסף  $S_\sigma$  של דפים שהוא בעל התכונות הבאות:

- א.  $S_\sigma$  קטנה יחסית.
  - ב.  $S_\sigma$  עשירה בדפים רלוונטיים.
  - ג.  $S_\sigma$  מכילה את רוב (או הרבה) ממקורות הסמכות החזקים.
- על ידי כך שנשמור את  $S_\sigma$  להיות קטנה נוכל להרשות לעצמינו את העלות החישובית על ידי שימוש באלגוריתמים לא טריוויאליים. על ידי הבטחת העושר בדפים הרלוונטיים אנו הופכים את התהליך למציאת מקורות סמכות טובים לקל יותר מכיוון שסביר להניח כי יש אליהם הפניות רבות בתוך  $S_\sigma$ .
- כיצד נוכל למצוא אוסף כזה של דפים? עבור פרמטר  $t$  (בדרך כלל הוא בערך 200), קודם לכן ניקח את  $t$  הדפים המדורגים ביותר עבור המחרוזת  $\sigma$  על ידי מנוע חיפוש כגון AltaVista או Hotbot.

קובץ זה נועד אך ורק לשימוש האישי של מורים למתמטיקה, פיזיקה, כימיה וביולוגיה ולהוראה בכיתותיהם. אין לעשות שימוש כלשהו בקובץ זה לכל מטרה אחרת, ובכלל זה: שימוש מסחרי, פרסום באתר אחר (למעט אתר בית הספר בו מלמד המורה), העמדה לרשות הציבור או הפצה בדרך אחרת כלשהי של קובץ זה או חלק ממנו.

מנועי חיפוש אלו היו קיימים באותה התקופה, לפני קיומו של Google. נתייחס ל  $t$  דפים אלו כאל קבוצה שורשית  $R_\sigma$ . קבוצה זו מקיימת את תנאים  $\alpha$  ו  $\beta$  שצויינו קודם אך היא רחוקה מאוד מלקיים את תנאי  $\gamma$ . על מנת לראות זאת, נשים לב כי  $t$  הדפים שבראש הרשימה המוחזרים על ידי מנוע החיפוש כולם יכללו את המחרוזת  $\sigma$  ולכן גם את  $R_\sigma$ . קבוצה זו היא בוודאות תת קבוצה של כל האוסף  $Q_\sigma$  של כל הדפים המכילים את  $\sigma$ . עם זאת טענו כי אפילו  $Q_\sigma$  בדרך כלל אינה מקיימת את תנאי  $\gamma$ . מעניין לציין כי בדרך כלל יש מעט מאוד קישורים בין דפים ב  $R_\sigma$ , מה שמשאיר אותה למעשה חסרת צורה. לדוגמא, בניסויים שלנו, הקבוצה השורשית עבור השאילתא "java" מורכבת מ 15 קישורים בין דפים מתחומים שונים. נוכל להשתמש בקבוצה השורשית  $R_\sigma$  על מנת לייצר את קבוצת הדפים  $S_\sigma$  שתענה על הדרישות שלנו. אם ניקח מקור סמכות חזק עבור שאילתת חיפוש רחב, ייתכן והוא לא יהיה כלול בקבוצה  $R_\sigma$ , אך סביר להניח שלפחות אחד הדפים ב  $R_\sigma$  יצביע עליו על ידי קישור. לכן נוכל להגדיל את מספר מקורות הסמכות החזקים על ידי הרחבת הקבוצה  $R_\sigma$ . בצורה זו נשיג את הקבוצה  $S_\sigma$  המתקבלת על ידי הרחבת הקבוצה  $R_\sigma$  בעזרת הוספת כל הדפים שישנם קישורים ב  $R_\sigma$  המצביעים עליהם. על מנת להבין זאת נסתכל על הציור הבא:



נרצה לעשות הגבלה על ידי כך שנאפשר לכל דף ב  $R_\sigma$  להביא לכל היותר  $d$  דפים המצביעים על דפים שיהפכו ל  $S_\sigma$ . פרט זה הוא חשוב מכיוון שאם נוסיף לקבוצה דפים ללא הגבלה נפגע בתכונה  $\alpha$  הדורשת שהקבוצה  $S_\sigma$  תהיה קטנה יחסית.

ישנה שיטה שהיא שימושית מאוד לצורך קיזוז השפעת הקישורים שמטרתם היחידה היא ניווט בתוך האינטרנט. נניח ישנו  $G[S_\sigma]$  כמו שתיארנו קודם. נבדיל בין שני סוגים של קישורים בתוך  $G[S_\sigma]$ . קישור חיצוני, הוא קישור המקשר בין דפים עם שמות תחום שונים וקישור פנימי,

קובץ זה נועד אך ורק לשימוש האישי של מורים למתמטיקה, פיזיקה, כימיה וביולוגיה ולהוראה בכיתותיהם. אין לעשות שימוש כלשהו בקובץ זה לכל מטרה אחרת, ובכלל זה: שימוש מסחרי, פרסום באתר אחר (למעט אתר בית הספר בו מלמד המורה), העמדה לרשות הציבור או הפצה בדרך אחרת כלשהי של קובץ זה או חלק ממנו.

המקשר בין דפים עם אותו שם תחום. קישור חיצוני הוא קישור היוצא ומצביע על דף של ארגון אחר. קישור פנימי הוא יוצא ומצביע על דף בתוך הארגון עצמו. הכוונה ב"שם תחום" כאן היא הרמה הראשונה של מחרוזת ה URL הקשורה לעמוד (Domain). כלומר, החלק הראשון של כתובת האינטרנט של הדף. בדרך כלל, מטרתם של הקישורים הפנימיים היא ניווט בתוך התשתית של האתר ולכן הם מספקים הרבה פחות אינפורמציה על הדפים האיכותיים שהם מצביעים אליהם. לכן נמחק את כל הקישורים הפנימיים של הגרף  $G[S_\sigma]$  ונשאיר רק את הקשתות המתייחסות לקישורים החיצוניים. התוצאה תהיה הגרף  $G_\sigma$ .

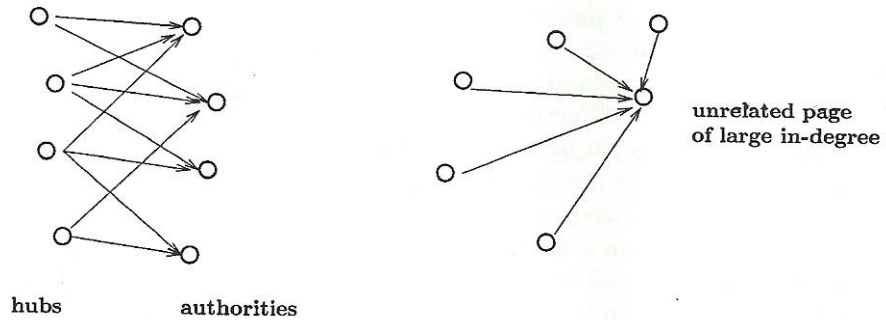
### חישוב רכזים ודפים איכותיים

ראינו כיצד ניתן ליצור תת גרף קטן  $G_\sigma$  שהוא מתמקד יחסית בנושא השאילתא. יש בו הרבה דפים רלוונטיים ומקורות סמכות חזקים. עכשיו נראה כיצד לחלץ את מקורות הסמכות מתוך אוסף כל הדפים. נעשה זאת אך ורק על ידי ניתוח הקשרים של  $G_\sigma$ . ניתן לטעון כי השיטה הפשוטה ביותר תהיה לסדר את הדפים על ידי דרגת הכניסה שלהם-מספר הקישורים אשר מצביעים עליהם ב  $G_\sigma$ . קודם לכן, דחינו את הרעיון כי היה מדובר בכל הדפים המכילים את המחרוזת  $\sigma$ . כעת יש לנו אוסף קטן של דפים רלוונטיים שבנינו במפורש, המכיל בתוכו את כל מקורות הסמכות שאנו רוצים למצוא. לכן מקורות סמכות אלו הם גם שייכים ל  $G_\sigma$  וגם יש להם הפניות בצורה רבה על ידי דפים מתוך  $G_\sigma$ . אכן, גישת הדרוג המתבצעת אך ורק על ידי דרוג דרגת הכניסה עובדת טוב יותר בהקשר של  $G_\sigma$  מאשר במסגרת שהצענו קודם לכן. במקרים מסויימים, זה יכול לייצר באופן אחיד תוצאות ברמה גבוהה. בכל זאת, לשיטה עדיין כמה בעיות רציניות. לדוגמא, עבור השאילתא "java" הדפים בעלי דרגת הכניסה הגבוהה ביותר מורכבים מ [www.gamelan.com](http://www.gamelan.com) ו [java.sun.com](http://java.sun.com) יחד עם דפי פרסומות לחופשות באיים הקאריביים ודף הבית של Amazon Books. תערובת זו מציגה סוג של בעיה שיכולה לצוץ בעקבות שיטה של דרוג פשוט. בזמן ששני הדפים הראשונים יכולים להיחשב כתוצאות טובות, השאר אינם רלוונטיים לנושא השאילתא. יש להם דרוג גבוה אך הם חסרות אחדות נושאית. הקושי הבסיסי שנגרם מכך הוא המתח הטמון בגרף  $G_\sigma$  בין מקורות הסמכות החזקים לבין דפים שהם פופולאריים בצורה אוניברסאלית. אנו מצפים מהסוג השני של הדפים דרגת כניסה גבוהה ללא תלות בנושא השאילתא הבסיסי. ניתן לתהות האם עקיפת בעיות אלו לא תדרוש לעשות שימוש עתידי בתוכן המילולי של של קבוצת הבסיס  $S_\sigma$  יותר מאשר במבנה הקישורי של  $G_\sigma$ . כעת נראה שזה לא המקרה. כלומר למעשה אפשרי לדלות אינפורמציה בצורה אפקטיבית יותר מהקשרים מאשר מהתוכן המילולי. נתחיל מהעובדה הבאה. דפים סמכותיים הרלוונטיים לשאילתא ההתחלתית, לא רק שהם אמורים להיות עם דרגת כניסה גבוהה, אלא מכיוון שהם כולם מקורות סמכות בעלי נושא משותף, אמורה להיות גם חפיפה רבה בקבוצות הדפים המצביעים עליהם. לכן, בנוסף לדפים בעלי הסמכות הגבוהה נצפה למצוא דפים רכזים. דפים רכזים אלו דפים שיש בהם

---

קובץ זה נועד אך ורק לשימוש האישי של מורים למתמטיקה, פיזיקה, כימיה וביולוגיה ולהוראה בכיתותיהם. אין לעשות שימוש כלשהו בקובץ זה לכל מטרה אחרת, ובכלל זה: שימוש מסחרי, פרסום באתר אחר (למעט אתר בית הספר בו מלמד המורה), העמדה לרשות הציבור או הפצה בדרך אחרת כלשהי של קובץ זה או חלק ממנו.

קישורים לדפים איכותיים רבים. דפים רכזים אלו מקבצים יחד דפים איכותיים בעלי נושא משותף ומאפשרים לנו לזרוק את הדפים הלא רלוונטיים בעלי הדרוג הגבוה. נסתכל על הציור הבא לצורך המחשה.



דפים רכזים ודפים איכותיים מקיימים יחס של חיזוק הדדי. דף רכז טוב הוא דף המצביע על הרבה דפים איכותיים. דף איכותי טוב הוא דף שמצביעים עליו הרבה דפים רכזים טובים. ניתן להגיד בברור כי על מנת לזהות דפים רכזים ודפים איכותיים של תת גרף  $G_\sigma$  אנו צריכים שיטה שתשבור מעגליות זו.

### האלגוריתם האיטרטיבי

נשתמש בקשר בין הדפים הרכזים לדפים האיכותיים בעזרת אלגוריתם איטרטיבי המתחזק ומעדכן את המשקלים המספריים של כל דף.

בהינתן גרף מכוון, נחלק משקלים לקודקודים. לדף  $p$  נתאים משקל סמכות  $x^{(p)}$  ומשקל רכזות  $y^{(p)}$ . נרצה שהמשקלים יהיו מנורמלים כלומר  $\sum_{p \in S_\sigma} (x^{(p)})^2 = 1$ ,  $\sum_{p \in S_\sigma} (y^{(p)})^2 = 1$ . כעת, נבצע את

התהליך האיטראטיבי הבא: תחילה נחלק משקלי רכזות לקודקודים כלומר את  $y^{(p)}$ . למצב זה נקרא  $y_0$ . כעת, נבצע לפי התור את הפונקציות הבאות:

$$Q: y^{(p)} \leftarrow \sum_{q:(p,q) \in E} x^{(q)}, F: x^{(p)} \leftarrow \sum_{q:(q,p) \in E} y^{(q)}$$

כלומר, נתחיל מוקטור משקלי רכזות אקראי מנורמל  $y_0 = z$  ואז סדר המשקלים יהיה:

$$y_0, x_0, y_1, x_1, y_2, x_2, \dots$$

למעשה, אם מטריצת הסמיכויות של הגרף היא  $A$  אז הפעולות הן:

$$Q: y_i \leftarrow Ax_i, F: x_i \leftarrow A^T y_i$$

קובץ זה נועד אך ורק לשימוש האישי של מורים למתמטיקה, פיזיקה, כימיה וביולוגיה ולהוראה בכיתותיהם. אין לעשות שימוש כלשהו בקובץ זה לכל מטרה אחרת, ובכלל זה: שימוש מסחרי, פרסום באתר אחר (למעט אתר בית הספר בו מלמד המורה), העמדה לרשות הציבור או הפצה בדרך אחרת כלשהי של קובץ זה או חלק ממנו.

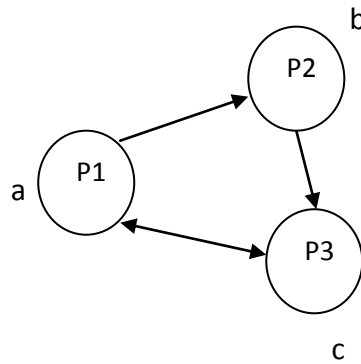
נסביר זאת באמצעות דוגמא:

נסתכל בגרף המכוון הבא, נניח כי הקודקודים הם:  $p_1, p_2, p_3$ .

$Y =$  משקל הרכזות. הסכום של אילו שהוא מצביע עליהם.

$X =$  משקל האיכות. הסכום של אילו שמצביעים עליו.

נניח כי משקלי הרכזות הם:  $(a, b, c)$ .



במקרה זה, מטריצת הסמיכויות היא:  $A = \begin{pmatrix} 0 & 1 & 1 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{pmatrix}$  והמטריצה המשחולפת היא

$A^T = \begin{pmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 1 & 1 & 0 \end{pmatrix}$ . נשים לב כי מטריצת הסמיכויות היא לא סימטרית מכיוון שהגרף הוא מכוון.

טענה: עבור כל מטריצות  $A^T A$  ו  $AA^T$  הן סימטריות (אך לאו דווקא שוות).

אם נעדיכן את משקלי הסמכות נקבל את הוקטור  $(c, a, a+b)$  וזה בדיוק:

---

קובץ זה נועד אך ורק לשימוש האישי של מורים למתמטיקה, פיזיקה, כימיה וביולוגיה ולהוראה בכיתותיהם. אין לעשות שימוש כלשהו בקובץ זה לכל מטרה אחרת, ובכלל זה: שימוש מסחרי, פרסום באתר אחר (למעט אתר בית הספר בו מלמד המורה), העמדה לרשות הציבור או הפצה בדרך אחרת כלשהי של קובץ זה או חלק ממנו.

את וקטור משקלי הרכזות אז באמת היינו מקבלים  $A^T \begin{pmatrix} a \\ b \\ c \end{pmatrix} = \begin{pmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 1 & 1 & 0 \end{pmatrix} \begin{pmatrix} a \\ b \\ c \end{pmatrix}$ . כנ"ל אם וקטור משקלי הסמכות היה  $(a, b, c)$  והיינו רוצים לעדכן

את וקטור משקלי הרכזות אז באמת היינו מקבלים  $A \begin{pmatrix} a \\ b \\ c \end{pmatrix} = (b+c, c, a)$  כנדרש.

לפי הגדרת הפעולות שדיברנו עליהם:  $x_i = A^T y_i$  כעת,  
 $y_{i+1} = Ax_i$

$$x_{i+1} = A^T y_{i+1} = (A^T A)x_i$$

$$x_{i+1} = (A^T A)x_i$$

מכאן,  $x_{i+2} = (A^T A)x_{i+1} = (A^T A)(A^T A)x_i = (A^T A)^2 x_i$

$$x_{i+k} = (A^T A)^k x_i$$

$$x_k = (A^T A)^k x_0$$

כמו כן,  $x_0 = A^T y_{i+1} \equiv Ax_z^T$

$$x_k = (A^T A)^k Ax_{i+1}^T = AA^T y_{i+1}$$

$$y_{i+2} = AA^T y_{i+1}$$

$$y_{i+1} = AA^T y_i$$

$$y_{i+k} = (AA^T)^k y_i$$

$$y_k = (AA^T)^k y_0 = (AA^T)^k z$$

הערה: למעשה בכל שלב אנו מנרמלים את הוקטור המתקבל, אם כי לא ציינו זאת במפורש.

נסמן:  $AA^T = B$ . סימטרית לפי הטענה. קבלנו כי:  $y_k = B^k y_0$ .

כל התהליך שעשינו מקביל Power Method רק שהוא מתבצע בשני שלבים. לכן לפי ה Power Method נקבל כי  $y_k$  מתכנס לוקטור העצמי  $y^*$  של  $B$  המתאים לערך העצמי הגדול ביותר של  $B$ .

השאלה היא למה זו בדיוק נקודת שבת?

התהליך מורכב משני שלבים:

קובץ זה נועד אך ורק לשימוש האישי של מורים למתמטיקה, פיזיקה, כימיה וביולוגיה ולהוראה בכיתותיהם. אין לעשות שימוש כלשהו בקובץ זה לכל מטרה אחרת, ובכלל זה: שימוש מסחרי, פרסום באתר אחר (למעט אתר בית הספר בו מלמד המורה), העמדה לרשות הציבור או הפצה בדרך אחרת כלשהי של קובץ זה או חלק ממנו.

אם נמצא זוג כך ש:  $Ax = \lambda y$  וגם  $A^T y = \lambda' x$ , נוכל לנרמל על ידי חלוקה ב  $\lambda$  ו  $\lambda'$  בהתאמה ואז נקבל כמעין נקודת שבת כי:  $AA^T y = A \cdot \lambda' x = \lambda' Ax = \lambda' \lambda y$  כלומר  $AA^T y = \lambda' \lambda y$  ו

$$y = \frac{1}{\lambda' \lambda} B y$$

לכן  $y$  היא נקודת השבת של המטריצה  $B = \frac{1}{\lambda' \lambda} B$ .

המסקנה היא כי כל התהליך מקביל ל Power Method.

כמו כן ניתן להגדיר מטריצה  $C = A^T A$  ובאותה הדרך להגיע למצב  $x_k = C^k x_0$  ואז

$$A^T Ax = \lambda \lambda' x, Cx = \lambda \lambda' x, Cx = x, \frac{1}{\lambda \lambda'} Cx = x$$

הגענו לתוצאה כי ל  $B$  ו  $C$  אותם ערכים עצמיים אך וקטורים עצמיים שונים.

מסקנה: בעזרת תהליך ה Power Method הגענו לוקטור המשקלים המתאים לערך העצמי המקסימלי. וקטור זה מהווה נקודת שבת ולכן נקבל דרוג לקודקודים.

הערה: האלגוריתם שהוצג כאן הוא רק אחד מתוך כמה שיקולים בעזרתם מנועי החיפוש מבצעים את הדרוג. הדרוג הסופי נבנה בעזרת איזון בין כמה שיקולים.

לפי התוצאות המוצגות במאמר ההתכנסות של Iterate היא די מהירה. למעשה, הממצאים הם כי מספיקות כ-20 אינטראציות כדי להתכנס להגיע למצב שבו כמעט ולא זזים מנקודת השבת. אחרי שהתכנסנו לנקודת השבת ניתן להגדיר את המספר  $c$ . מספר זה ייצג את כמות הדפים הראשונים בדרוג שנרצה שיוצגו על ידי מנוע החיפוש.

---

קובץ זה נועד אך ורק לשימוש האישי של מורים למתמטיקה, פיזיקה, כימיה וביולוגיה ולהוראה בכיתותיהם. אין לעשות שימוש כלשהו בקובץ זה לכל מטרה אחרת, ובכלל זה: שימוש מסחרי, פרסום באתר אחר (למעט אתר בית הספר בו מלמד המורה), העמדה לרשות הציבור או הפצה בדרך אחרת כלשהי של קובץ זה או חלק ממנו.

## רשימה ביבליוגרפית

- JON M. KEINBERG/Authoritative Sources in a Hyperlinked Environment/*Journal of the ACM, Vol. 46, No. 5, September 1999, pp. 604-632.*
- FAN R. K. CHUNG/ Spectral Graph Theory, 1994.