

קבוצות זרות – בעיית Union/Find

חומר קריאה לשיעור זה:

Chapter 22- Data Structures for Disjoint Sets (440 - 462)

קבוצות זרות – בעיית Union/Findנתון עולם של איברים U , $|U|=n$. לצורך הנוחות נניח $U = \{1,2,\dots,n\}$.**מבנה נתונים לשמירת קבוצות זרות תומך בפעולות הבאות:**1. $\text{Makeset}(i)$ – מחזיר קבוצה חדשה בעלת איבר בודד i .2. $\text{Find}(i)$ – מחזיר את הקבוצה לה שייך האיבר i .3. $\text{Union}(p,q)$ – מאחד את הקבוצות p ו- q , כלומר מחזיר קבוצה חדשה המכילה את איחוד האיברים בקבוצות p,q . (הקבוצות p,q המקוריות חדלות מלהתקיים.)לדוגמא, נניח ברגע נתון הגענו למצב: $\{1,3\}$ $\{5\}$ $\{2,4,6,7\}$ $p = \text{Find}(6)$ $q = \text{Find}(5)$ $r = \text{Union}(p,q)$ $s = \text{Find}(6)$ הקבוצה שהמשתנה r שומר היא $\{5,2,4,6,7\}$ המשתנים p ו- q אינם מחזיקים יותר קבוצות.והמשתנה s מחזיק את אותה קבוצה כמו r .

דוגמא היפותטית לשימוש

נתונות n ערים $\{1, \dots, n\}$. בתחילה כל הערים מנותקות. כל זמן מה בונים כביש ישיר בין שתי ערים. יש לאפשר את הפעולות הבאות:

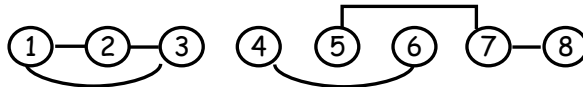
Add-Road(x, y) -- הוסף כביש ישיר בין שתי ערים x ו- y .

Check-Connectivity(x, y) -- קבע האם שתי הערים x ו- y מחוברות במסלול כלשהו.

פתרון: ניצור n קבוצות: $\{1\}, \{2\}, \dots, \{n\}$ ע"י: for ($i=1; i \leq n; i++$) Makeset(i);

Add-Road(x, y) ממומשת ע"י Union(Find(x), Find(y)).

Check-Connectivity(x, y) מחזירה "אמת" אם ורק אם Find(x)=Find(y).



| כביש שהוסף | {1} | {2} | {3} | {4} | {5} | {6} | {7} | {8} |
|------------|---------|-----|-----|-------|---------|-----|-----|-----|
| (1,2) | {1,2} | | {3} | {4} | {5} | {6} | {7} | {8} |
| (1,3) | {1,2,3} | | | {4} | {5} | {6} | {7} | {8} |
| (2,3) | {1,2,3} | | | {4} | {5} | {6} | {7} | {8} |
| (5,7) | {1,2,3} | | | {4} | {5,7} | {6} | | {8} |
| (4,6) | {1,2,3} | | | {4,6} | {5,7} | | | {8} |
| (7,8) | {1,2,3} | | | {4,6} | {5,7,8} | | | |

פתרון נאיבי ראשון

נשתמש במערך A מטיפוס SET בגודל n ובמונה counter המאותחל ל-0. בתא $A[i]$ נשמור את שם הקבוצה אליה שייך האיבר i . לדוגמא: $\{1,3\}$ $\{5\}$ $\{2,4,6,7\}$. מיוצגות ע"י המערך הבא (במימוש זה SET הוא פשוט int):

| | | | | | | |
|---|----|---|----|---|----|----|
| 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| 4 | 12 | 4 | 12 | 5 | 12 | 12 |

 Makeset(i) ממומשת ע"י $A[i]=++\text{counter}$.

Find(i) ממומשת ע"י $A[i]$.

Union(p, q) ממומשת ע"י הגדלת ה-counter באחד, מעבר על המערך A וכתיבת ערך ה-counter בכל מקום שבו כתוב p או q . הפונקציה מחזירה את הערך של counter.

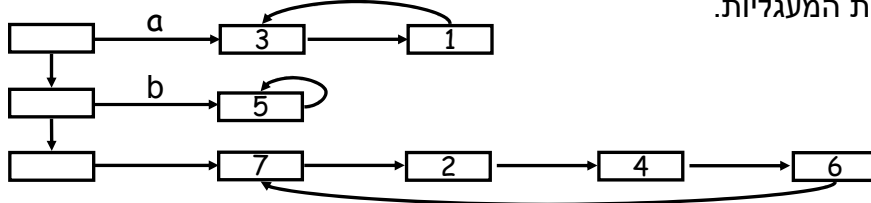
דוגמא: לאחר Union(p, q) כאשר $q=5, p=4$, מאוחדות הקבוצות p ו- q ומקבלים:

| | | | | | | |
|----|----|----|----|----|----|----|
| 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| 13 | 12 | 13 | 12 | 13 | 12 | 12 |

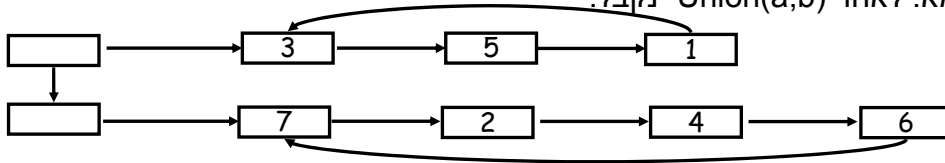
סיבוכיות הזמן של Find ו-Makeset היא $O(1)$ ושל Union היא $O(n)$.

פתרון נאיבי שני

נייצג כל קבוצה כרשימה מעגלית. נחזיק רשימה מקושרת של מצביעים אל הרשימות המעגליות.

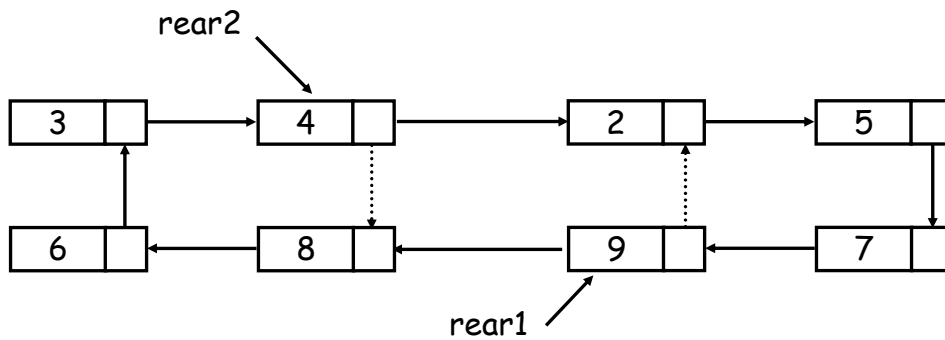


Union(p,q) -- ממושת ע"י איחוד הרשימות המוצבעות ע"י p ו-q. זמן $O(1)$.
 במימוש זה SET הוא מצביע לצומת מסוג NODE.
 דוגמא: לאחר Union(a,b) נקבל:



Find(i) ממושת ע"י מעבר על כל הרשימות עד שנמצא האיבר i. זמן $O(n)$.
 Makeset(i) ממושת ע"י הוספת רשימה עם איבר בודד. זמן $O(1)$.

תזכורת: איחוד רשימות מעגליות



הערות בהקשר לקידוד בשפת C

כותרות הפונקציות בכל המימושים הניתנים בשיעור זה:

SET Makeset(VALUE value)

SET Find(VALUE value)

SET Union(SET p,q)

typedef int VALUE;

typedef int SET;

במימוש הראשון ובמימושים נוספים

typedef NODE *SET;

במימוש השני ובמימושים נוספים

הערות

| Union | Find | Makeset | פתרון |
|--------|--------|---------|---------|
| $O(n)$ | $O(1)$ | $O(1)$ | נאיבי 1 |
| $O(1)$ | $O(n)$ | $O(1)$ | נאיבי 2 |

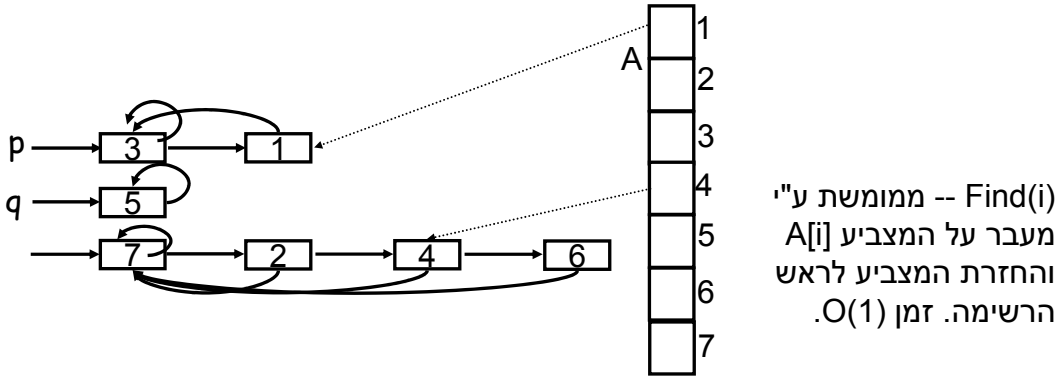
הערה נוספת: בהגדרת הבעיה שמות האיברים נבחרו להיות מספרים שלמים. ניתן להשתמש בשמות כלליים (כלומר מחרוזות תווים). לשם כך ניתן להשתמש במבנה נתונים הממיר ביעילות בין שמות כלליים ומספרים שלמים, למשל Hash table.

נפתח כעת פתרון ובו הזמן של פעולת Find הוא $O(1)$ והזמן המשוער (יוגדר בהמשך) של פעולת Union הוא $O(\log n)$.

ולסיום נפתח מבנה נתונים בו הזמן המשוער לפעולת Find הוא "כמעט קבוע" וזמן פעולת Union הוא $O(1)$.

פתרון שלישי

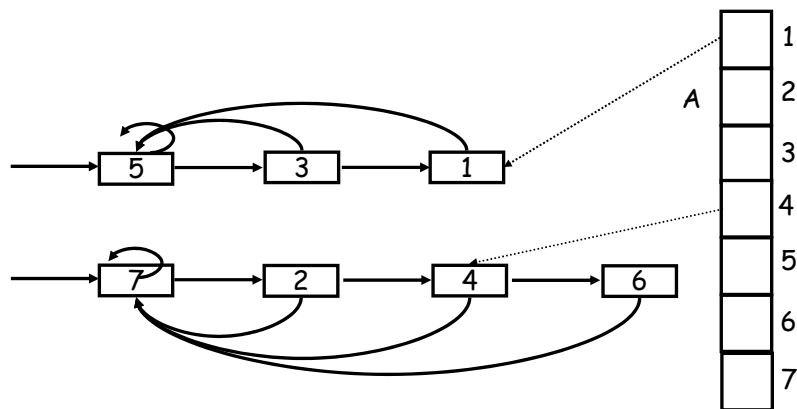
נייצג כל קבוצה כרשימה. כל איבר ברשימה מצביע גם לראש הרשימה וגם לאיבר הבא ברשימה. נחזיק מערך מיפוי איברים $A[i]$ ובו מצביע לאיבר i . קבוצה מוחזרת כמצביע לראש הרשימה המתאימה.



Union(p,q) -- ממומשת ע"י אחוד הרשימות המוצבעות ע"י p, q לרשימה אחת וכן עדכון מצביעי האיברים לראש הרשימה החדשה. הפונקציה מחזירה את ראש הרשימה החדשה. זמן $O(n)$.

דוגמא

דוגמא: לאחר Union(p,q) נקבל:

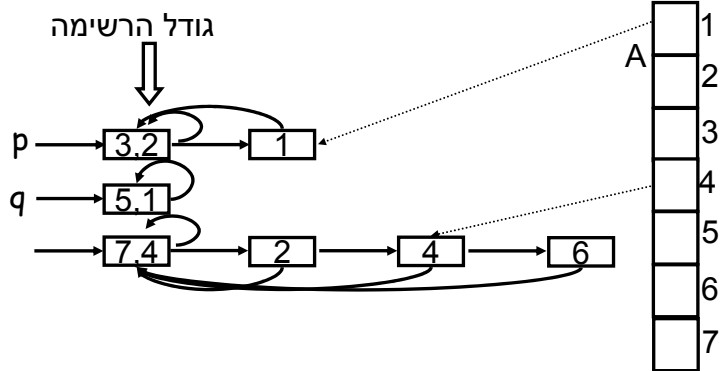


כאשר מאחדים שתי קבוצות יש לשנות באחת הקבוצות את כל המצביעים ליותרת החדשה. ברור שאת שינוי המצביעים עדיף לעשות בקבוצה הקטנה מבין השתיים (בניגוד לדרך שפעלנו בדוגמא זו).

פתרון שלישי משופר

שוב נייצג כל קבוצה כרשימה. בנוסף, בראש הרשימה נשמור את גודלה.

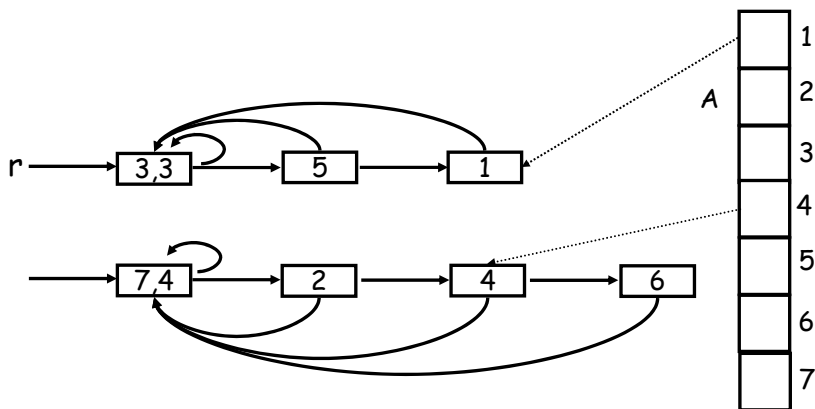
Find יתבצע כמו קודם ו- Union יתבצע ע"י הוספת הקבוצה הקטנה לתוך הקבוצה הגדולה. כלומר ראש הרשימה החדשה יהיה ראש הרשימה הגדולה יותר, וכל המצביעים יצביעו עליו.



כיצד נבצע $r = \text{Union}(p,q)$?

דוגמא

לאחר $r = \text{Union}(p,q)$ נקבל:



שימו לב שהקבוצה {5} שהייתה הקטנה מבין שתי הקבוצות הוספה לתוך הקבוצה הגדולה יותר {3,1} והקבוצה החדשה שנוצרה מוצבעת ע"י r.

ניתוח השיפור

טענה א: אם בכל איחוד מוסיפים את הקבוצה הקטנה לגדולה אזי כל איבר x משנה את הקבוצה אליה הוא שייך לכל היותר $\log_2 n$ פעמים כאשר n הוא מספר האיברים במבנה.

הוכחה: בכל פעם שאיבר x משנה את הקבוצה אליה הוא שייך הוא עובר לקבוצה חדשה הגדולה לפחות פי 2 מקבוצתו העכשווית. לאחר לכל היותר $\log_2 n$ מעברים הקבוצה הנוצרת מכילה את כל n האיברים, ולפיכך, שייכו של איבר x לא ישתנה יותר.

מסקנה: אם בכל איחוד מוסיפים את הקבוצה הקטנה לגדולה אזי הזמן הכולל לביצוע סדרה כלשהי של m פעולות Union/Find/Makeset הוא לכל היותר $O(m \log n)$.

הוכחה: n הוא מספר האיברים ב- U עבורם בוצע Makeset. לאור טענה א, סה"כ הזמן הנדרש לכל פעולות ה-Union הוא $O(n \log n)$ לכל היותר. מתקיים $m \leq n$. שאר הפעולות הן פעולות Find/Makeset שזמן כל אחת מהן הוא $O(1)$. לפיכך, סיבוכיות הזמן היא $O(m + n \log n)$. סיבוכיות זמן זו חסומה ע"י $O(m \log n)$.

לפיכך כל פעולה לוקחת זמן משוערך (Amortized time) של $O(\log n)$.

זמן משוערך

הגדרה: אם m פעולות מתבצעות בתוך זמן M , אזי הזמן המשוערך לכל פעולה מוגדר להיות M/m .

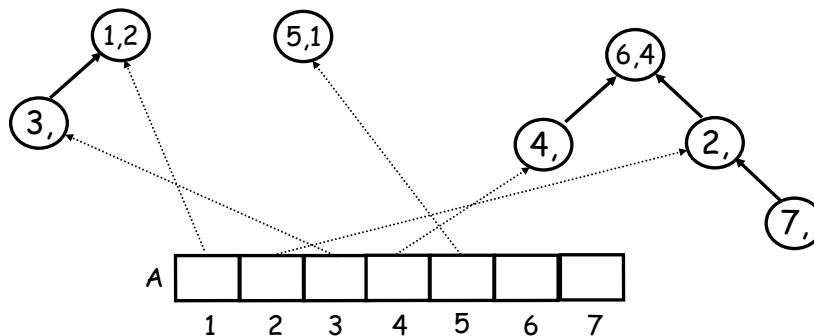
נימוקים להגדרה: כאשר ישנה סדרה ארוכה של פעולות שרובן קלות לביצוע (כגון find במימוש האחרון) וחלקן קשה לביצוע (כגון union באותו מימוש), אז ההשפעה של הפעולות הקשות מתחלקת על סדרת הפעולות כולה. יתכן גם מצב שזמן בצוע הפעולה משתנה כתוצאה מ"עזרה" הניתנת בזמן פעולה קודמת ואז הזמן המשוערך הוא מדד המתחשב בתרומה של פעולה אחת לרעותה.

מימוש נוסף: עצים הפוכים (Up trees)

לכל קבוצה ניצור עץ הפוך (בנים מצביעים להורה) ובו צומת לכל איבר בקבוצה. שורש העץ יכיל גם את מספר האיברי הקבוצה. בנוסף נחזיק מערך גישה לאיברים כמו במימוש השלישי.

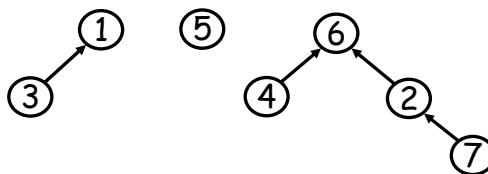
בזמן פעולת Union תולים את שורש העץ הקטן יותר מתחת לשורש העץ הגדול יותר ומעדכנים את גודל הקבוצה המאוחדת.

לדוגמה הקבוצות {1,3} {5} {2,4,6,7} מיוצגות ע"י:



מימוש עצים הפוכים בעזרת מערכים

בהנחה שמספר האיברים ידוע מראש, ניתן לייצג את כל הרשימות במערכים ללא שימוש במצביעים. ניתן לעשות זאת גם בכל המימושים הקודמים.



תאור גרפי:

| | | | | | | | |
|----------|----|---|---|---|----|----|---|
| size | 2 | | | | 1 | 4 | |
| parent | -- | 6 | 1 | 6 | -- | -- | 2 |
| elements | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

מימוש:

האיבר בשורש משמש מציין לקבוצה

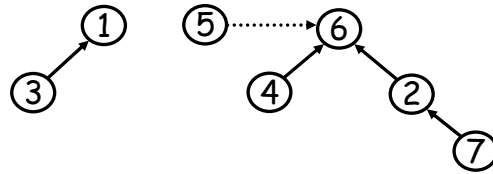
Find(i) -- ממומש ע"י סיור במעלה העץ. זמן $O(h)$ כאשר h הוא גובה העץ.

Union(p,q) -- ממומש ע"י אחוד העצים ששורשיהן p ו- q כך ששורש העץ הקטן יופנה לשורש העץ הגדול. זמן $O(1)$.

דוגמא

מצב התחלתי (כמו בשקף הקודם):

| | | | | | | | |
|----------|----|---|---|---|----|----|---|
| size | 2 | | | | 1 | 4 | |
| parent | -- | 6 | 1 | 6 | -- | -- | 2 |
| elements | 1 | 2 | 3 | 4 | 5 | 6 | 7 |



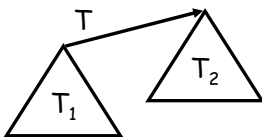
לאחר Union(5,6):

| | | | | | | | |
|----------|----|---|---|---|---|----|---|
| size | 2 | | | | 5 | | |
| parent | -- | 6 | 1 | 6 | 6 | -- | 2 |
| elements | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

זמן לפעולת Union הוא $O(1)$.

ניתוח גובה העץ

טענה: עבור עץ T בעל גובה h שנבנה מאיחודים של קבוצות קטנות לתוך גדולות, מתקיים $|T| \geq 2^h$.



הוכחה: באינדוקציה על h. עבור $h = 0$ ישנם 2^0 צמתים.

יהי T עץ בגובה h בעל מספר צמתים מינימלי. נבחן את פעולת האיחוד האחרונה:

• בפעולת האיחוד האחרונה חובר עץ T_1 לתוך עץ T_2 ולפיכך התקיים $|T_2| \geq |T_1|$.

• נסמן ב- h_1 וב- h_2 את גובה העצים T_1 ו- T_2 בהתאמה.

• במקרה א: $h_1 < h_2$. סתירה שכן אז T הוא עץ בעל גובה זהה ל- T_2 ואינו בעל מספר צמתים מינימלי.

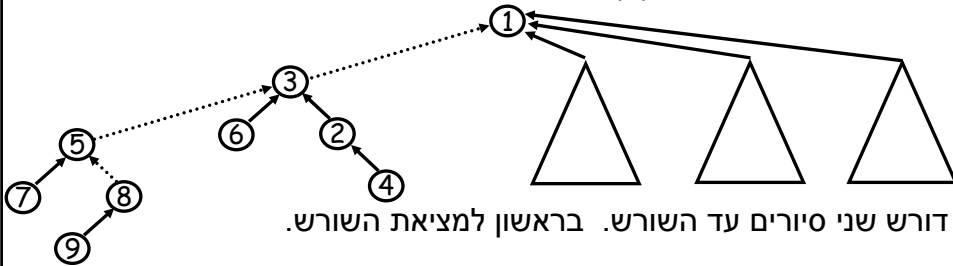
• במקרה ב: $h_1 \geq h_2$. במקרה זה מתקיים $h = h_1 + 1$ ולפיכך

• $|T| = |T_1| + |T_2| \geq 2 \cdot |T_1| \geq 2 \cdot 2^{h-1} = 2^h$ כנדרש.

מסקנה 1

שיפור נוסף: כיווץ מסלולים

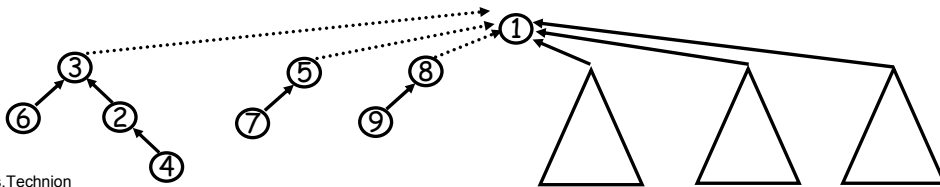
בזמן ביצוע $Find(i)$, עדכן את שדה ה-parent של כל הצמתים מ- i ועד השורש כך שיצביעו ישירות לשורש. דוגמא: $Find(8)$.



המימוש דורש שני סיורים עד השורש. בראשון למציאת השורש.

בשני לעדכון המצביעים לכיוון השורש.

תהליך העדכון:



ניתוח זמנים

משפט: במימוש באמצעות אוסף עצים של n איברים, סיבוכיות זמן בצו m פעולות $union/find/makeset$ חסומה ע"י $O(m \log^* n)$.

לפיכך הזמן המשוערך (amortized time) לכל פעולה חסום ע"י $O(\log^* n)$,

כאשר הפונקציה $\log^* n$ מוגדרת כדלהלן:

$$\log^{(i)}(n) = \overbrace{\log_2 \log_2 \dots \log_2}^i n$$

$$\log^*(n) = \min \{i \geq 0 \mid \log^{(i)}(n) \leq 1\}$$

ובמילים, $\log^* n$ הוא מספר הפעמים שצריך לקחת \log כדי לקבל מספר קטן או שווה לאחד. זוהי פונקציה מונוטונית העולה בקצב מאד איטי. לכל מספר n פרקטי $\log^* n \leq 5$. כלומר פעולת $find$ לוקחת זמן משוערך שהוא קבוע מבחינה מעשית.

$$\log^*(1) = 0 \quad \log^*(2) = 1 \quad \log^*(2^2) = 2 \quad \log^*(2^{2^2}) = 3 \quad \log^*(2^{2^{2^2}}) = 4$$

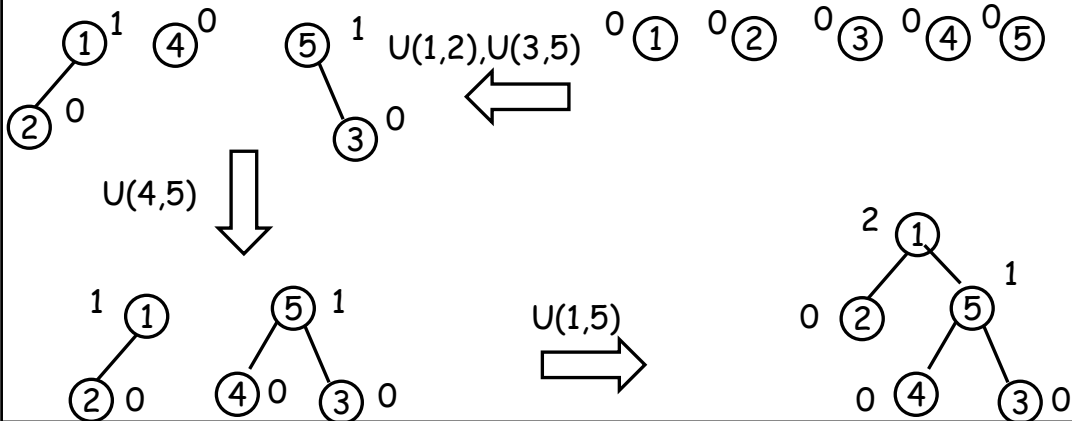
$$\log^*(2^{2^{2^{2^2}}}) = \log^*(2^{65,536}) = 5$$

הוכחת המשפט ותאור חסם הדוק עוד יותר באמצעות פונקציות אקרמן ניתן למצוא בספר

הלימוד בעמודים 450-453.

וריאנט של האלגוריתם – שימוש ב-rank

לכל איבר נוסף שדה rank, כך שהשורש של קבוצת יחידה מקבל $rank(x) = 0$.
 בזמן $Union(x,y)$: אם $rank(x) < rank(y)$ הוסף את x כבן נוסף של y.
 אם יש שוויון, הוסף שרירותית את x כבן נוסף של y, אך הגדל את $rank(y)$ באחד.



וריאנט של האלגוריתם – שימוש ב-rank

לכל איבר נוסף שדה rank, כך שהשורש של קבוצת יחידה מקבל $rank(x) = 0$.
 בזמן $Union(x,y)$: אם $rank(x) < rank(y)$ הוסף את x כבן נוסף של y.
 אם יש שוויון, הוסף שרירותית את x כבן נוסף של y, אך הגדל את $rank(y)$ באחד.

טענה 1: $rank(x) < rank(\text{parent}(x))$

הוכחה: מתקיים כאשר x הופך לבן של $\text{parent}(x)$.

אחרי זמן זה $rank(x)$ לא גדל, אך $rank(\text{parent}(x))$ לא קטן (יכול אפילו לגדול).



מתי גדל?

• בזמן Union כאשר לאבא מצטרף תת-עץ בעל אותו rank

• בזמן Find כאשר הבן מתנתק מהאבא והופך להיות בן של השורש

טענה 2:

בעץ שלשורשו יש $\text{rank}(r) = k$ יש לכל הפחות 2^k צמתים.
הוכחה: אינדוקציה על $\text{rank}(r)$.

$$\text{rank}(x) \leq \log n \quad \text{מסקנה}$$

טענה 3 יש לכל היותר $n/2^r$ צמתים x בעלי $\text{rank}(x) = r$
הוכחה: בזמן ש- x קיבל $\text{rank}(x) = r$ היו בעץ שלו לפחות 2^r צמתים.
 נסמן ב- x את כל הצמתים בעץ הנ"ל.
 במהלך האלגוריתם כל צומת מסומן פעם אחת לכל ערך של r .
 כיוון מסלולים לא משנה את rank .

חישוב זמן ה-Find:

נחלק את הצמתים לבלוקים, כך שבלוק i יכיל את כל צמתים v כך ש- $\log^* \text{rank}(v) = i$.

$$g(n) = 2^{\left\{ \begin{matrix} 2 \\ \vdots \\ 2 \end{matrix} \right\}^n} \quad \text{תהי } g \text{ הפונקציה ההפוכה ל-} \log^*$$

$$\text{Block}(j) = \{v : \log^* \text{rank}(v) = j\} = \{v : g(j-1) < \text{rank}(v) \leq g(j)\}$$

חלק מהמחיר של ה-Find נזקוף לבלוק (יסומן B),
 וחלק למסלול (יסומן P). נחסום כל חלק בנפרד.

פעולות B:

נסתכל על מסלול:

את המעבר על הצומת האחרון של כל בלוק נזקוף לחובת הבלוק.
וכן נזקוף לחובת הבלוקים את המעבר על בן של השורש.

כיון ש- $rank(v) \leq \log n$,

בכל מסלול מספר פעולות B חסום ע"י $n \log^* n \leq 2 \log^* \log n + 1$ (צמתים + השורש ובנו).

טענה מהרגע שאביו של צומת ואביו שייכים לבלוקים שונים,
לצומת לא יהיה אב ששייך לאותו בלוק כמוהו.

הוכחה כיון שהבן איננו שורש ה- $rank$ שלו לא יגדל, ולכן הבן יישאר באותו בלוק.
אם האב הוא שורש, ה- $rank$ שלו עשוי לגדול, ואז הוא עלול לעבור לבלוק
עם $rank$ עוד יותר גדול, ועוד יותר רחוק מהבלוק של הבן.

פעולות P:

כמה צמתים יש בבלוק ?

$$\begin{aligned}
 N(j) &= \sum_{g(j-1) < i \leq g(j)} \# \text{ nodes with rank } i \\
 &\leq \sum_{g(j-1) < i \leq g(j)} n/2^i = \frac{n}{2^{g(j-1)+1}} \sum_{0 \leq k \leq g(j)-g(j-1)-1} 1/2^k \\
 &< \frac{n}{2^{g(j-1)+1}} \sum_{0 \leq k} 1/2^k \\
 &= \frac{n}{2^{g(j-1)+1}} 2 = \frac{n}{2^{g(j-1)}} \\
 &= \frac{n}{g(j)}
 \end{aligned}$$

פעולות P (המשך)

כיון ש- $\text{rank}(x) < \text{rank}(\text{parent}(x))$ (טענה 1),
 בפעולת P על צומת x, $\text{parent}(x)$ הופך להיות השורש.
 כיון שבפעולת P x אינו שורש או בנו, $\text{rank}(\text{parent}(x))$ גדל ממש.
 יהי $x, \text{parent}(x) \in \text{Block}(j)$.
 כיון ש- $\text{rank}(\text{parent}(x))$ לא יכול לגדול יותר $g(j)-g(j-1)$ פעמים מבלי ש-
 $\text{parent}(x)$ יעבור לבלוק אחר,
 אף צומת בבלוק j לא יכול לקבל יותר מאשר $g(j)-g(j-1)$ זקיפות P.
 מרגע ש-x ו- $\text{parent}(x)$ לא נמצאים באותו הבלוק, לא ייזקפו ל- x פעולות P.

מכאן:

$$\text{Total cost of P} = \sum_{0 \leq j < \log^* n} N(j)(g(j) - g(j-1))$$

$$\leq \sum_{0 \leq j < \log^* n} \frac{n}{g(j)} (g(j) - g(j-1)) < \sum_{0 \leq j < \log^* n} n = n \log^* n$$

פונקציית אקרמן

נגדיר את פונקציית אקרמן

$$A(1, j) = 2^j$$

$$A(i, 1) = A(i-1, 2)$$

$$A(i, j) = A(i-1, A(i, j-1))$$

| | j=1 | j=2 | j=3 |
|-----|-----------|---|---|
| i=1 | 2^1 | 2^2 | 2^3 |
| i=2 | 2^2 | 2^{2^2} | $2^{2^{2^2}}$ |
| i=3 | 2^{2^2} | $2^{\left. \begin{matrix} 2^2 \\ 2^2 \end{matrix} \right\} 16}$ | $2^{\left. \begin{matrix} 2^2 \\ 2^2 \end{matrix} \right\} 2^{\left. \begin{matrix} 2^2 \\ 2^2 \end{matrix} \right\} 16}$ |

הפונקציה ההפוכה

$$\alpha(m, n) = \min \{i \geq 1 : A(i, \lfloor m/n \rfloor) > \log n\}$$

אם פונקציה גדלה מהר אז ההופכית שלה גדלה לאט.

$$\alpha(m, n) \leq 4 \text{ לכל ערך סביר של } m, n.$$

משפט Tarjan (בלי הוכחה)

הזמן המשוערך של m פעולות Union/Find על קבוצות זרות מעל n איברים הוא

$$\theta(m\alpha(m, n))$$