

6.12.2009

תכנון וניתוח אלגוריתמים
הרצאה 7
אלגוריתמים הסתברותיים

המרצה: בועז פת-שמיר
רשימות: יונתן נתיב

בעיית המזכירה (Secretary problem)

מקרה לדוגמא:

נניח שאנו רוצים לקחת אדם לעבודה. בחדר ההמתנה יש N מועמדים. תוך כדי הראיון אנו קובעים את ציון ההתאמה של המועמד.

התנאים:

- 1) הציונים יכולים להיות כל מספר – לא חסום, בלי התפלגות א-פריורית.
 - 2) בסוף הראיון אנו מחוייבים להגיד לאדם אם הוא התקבל (לא מראינים את השאר), או לא התקבל (אין אפשרות להחזיר אותו).
 - 3) הצלחה מוגדרת אם בחרנו את המועמד הטוב ביותר.
- בעיה זו נקראת בעיית המזכירה.

הגדרת הבעיה:

- קלט: n מועמדים
- שיטה: קוראים לכל מועמד, מוצאים את ערכו, מחליטים האם לקבלו או לא. אם המועמד התקבל נגמר המשחק.
- מטרה: לשכור את הטוב ביותר.

טענה: אין אלגוריתם דטרמיניסטי אשר תמיד פותר את בעיית המזכירה בהצלחה.

הוכחה:

בהינתן אלג' A כלשהו, נראה כי קיים קלט עליו A נכשל. נקרא ערכים ל $n-1$ מועמדים ראשונים כלשהם. אם A בחר לשכור מועמד מבין $n-1$ המועמדים הראשונים נגדיר את המועמד ה n כבעל ציון התאמה גבוה מכל שאר המועמדים. אם A לא בחר מועמד עד כה ניתן למועמד ה n ציון נמוך מהמקסימום עד כה. בכל מקרה האלג' נכשל.

ראינו שאלגוריתם דטרמיניסטי לא עובד. הדבר המפתיע שיש אלגוריתם שכן עובד, ע"י איסוף סטטיסטיקה על הראיונות שבוצעו.

אלגוריתם פשוט שמצליח בהסתברות $\frac{1}{4}$ לבחור את המועמד הטוב ביותר:

1. האלגוריתם מסדר המועמדים בפרמוטציה אקראית.
2. מראיין את החצי הראשון של המועמדים וזוכר את הציון הגבוה ביותר. נסמנו h .
3. מראיין את החצי השני ושוכר את הראשון שציונו גבוה מ h .

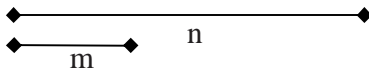
ניתוח האלגוריתם:

- בהסתברות $\frac{1}{2}$, הטוב ביותר בחצי השני
- בהסתברות $\frac{n/2}{n-1}$ השני הכי טוב בחצי הראשון
- כששני התנאים מתקיימים (בהסתברות $\sim 1/4$), האלגוריתם הצליח.

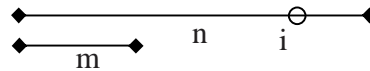
נראה כיצד ניתן לשפר את האלגוריתם הפשוט

אופטימיזציה:

1. כמקודם.
2. מראיין m מועמדים, זוכר את הטוב.
3. כמקודם.



בפשטות, אנחנו הופכים את הבעיה ליותר כללית, מראינים m מועמדים במקום $\frac{n}{2}$. השאלה היא מה ה- m שכדאי לקחת, ואת זה נחליט ע"י ניתוח מדויק של הסתברות ההצלחה.



נבצע ניתוח בהנחה שהמקסימום נמצא במקום i (נשים לב שהצלחה תתקבל אם האיבר השני בגובהו מתוך I האיברים נמצא ב m האיברים הראשונים).

$$\Pr[\text{success}] \stackrel{\text{bayes}}{=} \sum_{i=m+1}^n \Pr[\text{success} | \text{best in position } i] * \Pr[\text{best in position } i] =$$

no chance for success when $i < m$

$$\sum_{i=m+1}^n \Pr[\text{second best of } 1..i \text{ in pos } 1..m] \cdot \Pr[\text{best in position } i]$$

$$\Pr[\text{best in position } i] = \frac{1}{n}$$

$$\Pr[\text{second best of } 1..i \text{ in pos } 1..m] = \frac{m}{i-1}$$

$$\Pr[\text{success}] = \frac{1}{n} \sum_{i=m+1}^n \frac{1}{n} \frac{m}{i-1} = \frac{1}{n} \sum_{i=m+1}^n \frac{m}{i-1} = \frac{m}{n} \sum_{i=m+1}^n \frac{1}{i-1} =$$

$$= \frac{m}{n} (H_n - H_m) \approx \frac{m}{n} (\ln(n) - \ln(m)) = \frac{m}{n} \ln\left(\frac{n}{m}\right)$$

נרצה לעשות אופטימיזציה – למקסם את התוצאה שקיבלנו. לשם כך נציב $x = \frac{m}{n}$,

$$f(x) = x \cdot \ln\left(\frac{1}{x}\right)$$

$$f'(x) = \ln\left(\frac{1}{x}\right) + x \cdot \frac{1}{x} \cdot \frac{-1}{x^2} = \ln\frac{1}{x} - 1$$

$$f''(x) = \frac{1}{x} \cdot \frac{1}{x^2} < 0 \quad \forall x > 0$$

נשווה נגזרת ל-0 למצוא מקסימום:

$$\ln\left(\frac{1}{x}\right) = 1$$

$$x = \frac{1}{e}$$

$$\text{ואז הסתברות ההצלחה: } \frac{1}{e} \ln\left(\frac{1}{e}\right) = \frac{1}{e} \text{ בערך } 137\%$$

בעיית אספון הקלפים coupon collector

- קלט: N סוגי קלפים
- שיטה: בכל ניסוי מקבלים קלף באקראי כשהסתברות לסוג i היא 1/n לכל i
- מטרה: צריך לאסוף כל הסוגים
- שאלה: כמה ניסויים יש לבצע עד אשר מחזיקים את כל הקלפים?

הבעיה במילים: קונים קורנפלקס ואוספים קלף כל פעם. כדי לקבל פרס, יש לאסוף את כל סוגי הקלפים. כמה קופסאות צריך לקנות כדי להשלים את הסדרה?

תשובה I תוחלת:

נגדיר משתנה אקראי X_j $1 \leq j \leq n$ – אורך הזמן מהרגע שיש j סוגים שונים ועד שיש j+1 סוגים שונים.

נשאל: מהי התוחלת $E[X_j]$?

הסיכוי שקלף חדש אינו אחד מ j הסוגים שכבר התקבלו הוא: $\frac{n-j}{n}$

זהו משתנה גיאומטרי עם $p = \frac{n-j}{n}$ ולכן התוחלת שלו $\frac{n}{n-j}$.

$$E\left[\begin{array}{l} \#cards_to \\ complete_all \end{array}\right] = E\left[\sum_{j=0}^{n-1} x_j\right] = \sum_{j=0}^{n-1} E[x_j] : \text{לפי ליניאריות של התוחלת}$$

$$\sum_{j=0}^{n-1} E[x_j] = \sum_{j=0}^{n-1} \frac{n}{n-j} = nH_n$$

תשובה II "הסתברות גבוהה":

נרצה לאמר כעבור כמה זמן נסיים בהסתברות כלשהי להצלחה ← מה ההסתברות שלא נקבל קלף מסוג I מסוים ב t ניסיונות?

$$\left(1 - \frac{1}{n}\right)^t = \left[\left(1 - \frac{1}{n}\right)^n\right]^{\frac{t}{n}} < e^{-\frac{t}{n}}$$

נדרוש: $P\left[\begin{array}{l} card_ \#i_ missing \\ in_t_tries \end{array}\right] < \frac{1}{n^c}$ עבור $c > 0$ כלשהו.

הסתברות פולינומים קטנה, מספיק:

$$e^{-\frac{t}{n}} < n^{-c}$$

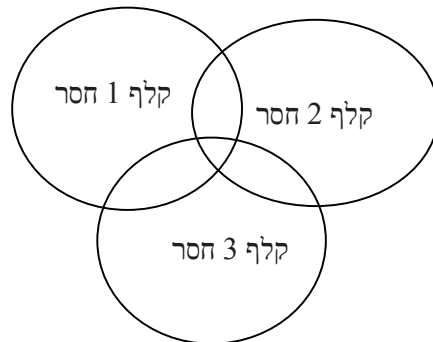
$$\frac{-t}{n} < -c \ln n$$

$$t > c \cdot n \cdot \ln n$$

$$P\left[\begin{array}{l} all_cards \\ in_t_tries \end{array}\right] = 1 - P\left[\begin{array}{l} card_ \#i_ missing \\ in_t_tries \end{array}\right]$$

נשתמש בחסם "ברברי" – חסם האיחוד (union bound):

במקום לסכום ההסתברות על איחוד השטחים (סיכמנו את השטחים החופפים מספר פעמים)



$$1 - P\left[\begin{array}{l} card_ \#i_ missing \\ in_t_tries \end{array}\right] \geq 1 - \sum_{i=1}^n P\left[\begin{array}{l} card_ \#i \\ is_ missing \end{array}\right] > 1 - ne^{-\frac{t}{n}} > 1 - n^{-c+1}$$

for: $t = cn \ln n$

הגדרה: מאורע X קורה בהסתברות גבוהה אם $P[x]=1-o(1)$ (ההסתברות לכישלון פחותה מקבוע).
 בד"כ, לצרכי חסם האיחוד: $c > 0$ $P[x]=1-O(n^{-c})$

הטלת מטבע הוגן עד לקבלת "עץ"

הגדרה: בהטלת מטבע הוגן נרצה בהסתברות גבוהה לקבל עץ. מס' ההטלות הדרוש כדי שבהסתברות לפחות $1-\delta$ נקבל "עץ".

$$P\left[\begin{array}{l} \text{no_pally} \\ \text{in_t_throws} \end{array}\right] = \left(\frac{1}{2}\right)^t$$

$$t > -\log_2 \delta$$

$$t > \log_2 \frac{1}{\delta} \quad \text{ונקבל } \left(\frac{1}{2}\right)^t < \delta$$

לקבלת $\delta = n^{-c}$ צריך להטיל $t > \log_2 n^c = c \log_2 n$

$$\delta = (1-p)^t \quad \text{ואז} \quad P\left[\begin{array}{l} t \\ \text{failures} \end{array}\right] = (1-p)^t$$

אם הסתברות הצלחה בניסיון בודד היא p אז

$$t > \frac{\log \delta}{\log(1-p)} \quad \text{כלומר:} \quad \log \delta < t \log(1-p)$$

- הניתוח מתבסס על חוסר תלות בניסויים.
 (אם נשנה את הניסוי למטבע שכל הטלה חוזרת על תוצאת ההטלה הראשונה התוחלת לא תשתנה אך הבטחה לסיכוי הצלחה מסוים היא לא פשוטה).

MAX3SAT

קלט: נוסחת 3CNF : m פסוקיות, כ"א 3 ליטרלים (משתנה או שלילתו).
נסמן n – מספר משתנים.

מטרה: לספק כמה שיותר פסוקיות.

ננסה להציב לכל המשתנים True או לכל המשתנים False, אחת משתי ההצבות הללו בהכרח תספק יותר ממחצית הפסוקיות (כל הפסוקיות השליליות הופכות לחיוביות בהצבה הנגדית).

נבדוק השמה אקראית: ניתן F/T בהסתברות $\frac{1}{2}$ לכל משתנה.
הסתברות שפסוקית מקבלת שקר – $\frac{1}{8}$ (הסתברות חצי עבור כל ליטרל בפסוקית).

מלינאריות של תוחלת: מספר הפסוקיות המסופקות = $m \cdot \frac{7}{8}$.

אי שוויון מרקוב (מורחב):

אינטואיציה – אם תוחלת גובהו של אדם הוא 1.75, אז פחות מ $\frac{1}{2}$ מהאנשים הם בגובה מעל 3.5 (אחרת התוחלת היתה מעל 1.75)

פורמלית:

יהי x משתנה אקראי דיסקרטי, וידוע ש $0 \leq X \leq \beta$.

רוצים לחשב $p = \Pr[x \leq L]$ ל נתון.

$$E[x] = \sum_{j=0}^B j \cdot \Pr[x = j] = \sum_{j=0}^L j \cdot \Pr[x = j] + \sum_{j=L+1}^B j \cdot \Pr[x = j]$$

$$\leq \sum_{j=0}^L L \cdot \Pr[x = j] + \sum_{j=L+1}^B B \cdot \Pr[x = j]$$

$$p = \Pr[x \leq L]$$

$$L \cdot p + B(1 - p)$$

$$\Rightarrow p \leq \frac{B - E[x]}{B - L}$$

מה ההסתברות שהשמה אקראית תספק לפחות $m \cdot \frac{7}{8}$ פסוקיות?

נשתמש באי שוויון מרקוב

X – מס' פסוקיות מסופקות

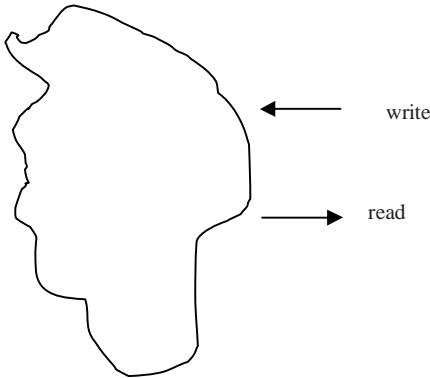
$$7/8m - E[x]$$

$$m - B$$

$$(7/8 - e)m - L$$

ולכן:

$$\Pr[\text{less than } \frac{7}{8}m \text{ sat.}] \leq \frac{m - \frac{7}{8}m}{m - (\frac{7}{8}m - \epsilon)} = \frac{1/8}{1/8 + \epsilon} = \frac{1}{1 + 8\epsilon} = 1 - \frac{8\epsilon}{1 + 8\epsilon}$$



מערכת קוורומים Quorum system

יש חווה של n שרתים לא אמינים. השרתים נותנים זכרון משותף. יש פעולות כתיבה וקריאה בכתיבה נכתוב יחד עם המידע את זמן הכתיבה. בקריאה ניקח את הערך המעודכן ביותר.

Quorum system

$|X|=n$ של שרתים

מערכת קוורומים: $S_i \in X, Q = \{S_1 \dots S_n\}$
 כך שלכל $S, S' \in Q, S \cap S' \neq \emptyset$ מתקיים

דוגמא טריוויאלית:

- הקוורום בנוי משרת יחיד $Q = \{\{1\}\}$
- יש קוורום אחד בלבד
- המערכת פגיעה לנפילות

דוגמא:

$$Q = \left\{ S : |S| > \frac{n}{2} \right\}$$

2 קבוצות שכל אחת מכילה לפחות חצי מהאיברים:

החסרון: בכל קריאה/כתיבה, יש להגיע למספר לינארי של מעבדים, ולכן עושים הרבה עבודה. $\Omega(n)$ - חסם תחתון של עבודה לגישה.

בכל כתיבה, תופסים איזשהו קוורום וכותבים לכל המעבדים בקוורום זה עם timestamp. כדי לקרוא, בוחרים קוורום וקוראים מכל המעבדים בו.

בהנתן מערכת קוורומים (מ"ק) נגדיר:

עמידות: #מעבדים מינימאלי שבלעדיהם אין קוורום חוקי (עמידות לתקלות).

עומס: עומס לשרת (תחת בחירה אקראית) – כמה עבודה (מס' הגישות) כל שרת כזה עושה?

$$\text{נניח } p \in \{S_1 \dots S_k\} \text{ ו- } |Q|=m$$

אז העומס על שרת p הוא k/m . (ניגשים לכל המעבדים בקוורום והשרת נמצא ב k מתוך m הקוורומים) עומס המערכת – עומס מקסימלי לשרת (אם בוחרים קוורום באקראי, מה תוחלת מקסימום הגישות למחשב). בדוגמא השניה, עומס המערכת הוא בערך $1/2$, כי בכל פעולה יש סיכוי $1/2$ שניגש לשרת מסויים (שנמצא ב $1/2$ מהקוורומים).

זמן: תוחלת גודל הקוורום (מכיוון שכותבים/קוראים מכולם).

דוגמאות:

$$Q = \{ \{1\} \} \text{ עמידות } = 1, \text{ עומס } = 1, \text{ זמן } = 1.$$

$$Q = \{ \{1,2,3\} \} \text{ עמידות } = 3, \text{ עומס } = 1, \text{ זמן } = 3.$$

$$Q = \{ S : |S| = (n+1)/2 \} \text{ עמידות } = n/2, \text{ עומס } = 1/2, \text{ זמן } = n/2.$$

נרצה לסדר את תתי הקבוצות באופן כזה שיביא למקסימום את העמידות ולמינימום את העומס וזמן.

תכונה:

תהי Q מ"ק עם עומס L (התוחלת של הכי גרוע), קוורום מינימאלי בגודל C , ו- n שרתים סה"כ.

$$L \geq \max\left(\frac{1}{c}, \frac{c}{n}\right) \text{ אז}$$

הוכחה:

כיוון שכל קוורום, גודלו $C \leq$,

בכל גישה מתבצעות לפחות C פניות לשרתים, לכן $L \geq \frac{c}{n}$.

בנוסף נסתכל על קוורום בגודל מינימלי:

בכל גישה לפחות אחד משרתיו משתתף \Leftarrow עומס שרתיו הממוצע $\leq C$. ■

$$\frac{1}{\sqrt{n}} \leq \text{עומס השרת: מסקנה}$$

בשביל לקבל מערכת אופטימאלית, נרצה ששני הביטויים יהיו זהים ולכן נבחר $c = \sqrt{n}$, ואז $L \geq 1/\sqrt{n}$.
מערכת אופטימלית

$$n = C^2$$

נסדר השרתים בריבוע $C \times C$.

הקוורומים: שורות/עמודות.

שורות לכתובה, עמודות לקריאה.

מערכת קוורומים ממש: כל קוורום: שורה + עמודה. כל פעם שנרצה לכתוב למשל, נבחר שורה ונכתוב בכולה.

$$L = \frac{2\sqrt{n} - 1}{n} \approx \frac{2}{\sqrt{n}}$$

פגיעות: \sqrt{n} (כל עמודה או שורה שנפלו מהמערכת).
עד כאן המערכת דטרמיניסטית.

מערכת הסתברותית:

נדרוש שלכל $S, S' \in Q$, $Pr[S \cap S' \neq \emptyset] > 1 - \delta$, ל $\delta < 1$ נתון.

נבחר קוורומים אקראיים בגודל C .

$$Pr[S \cap S' = \emptyset] = \left(\frac{n - |S|}{n}\right)^{|S'|} = \left(1 - \frac{c}{n}\right)^c = \left[\left(1 - \frac{c}{n}\right)^{\frac{n}{c}}\right]^{\frac{c^2}{n}} \approx e^{-c^2/n}$$

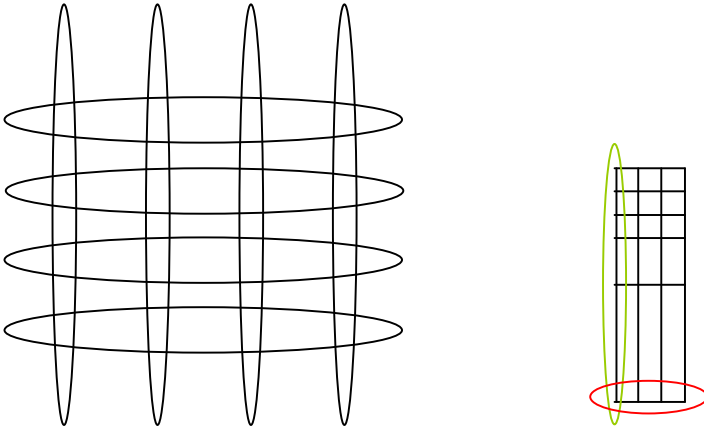
$$\delta > e^{-c^2/n}$$

$$\ln \delta > -\frac{C^2}{n}$$

$$\ln \frac{1}{\delta} < \frac{C^2}{n}$$

$$\sqrt{n \ln \frac{1}{\delta}} < C$$

אם הסתברות הכשלון פולינומית קטנה, אז מספיק $C > O(\sqrt{n}, \sqrt{\log n})$



פרדוקס יום ההולדת Birthday Paradox

מה הסיכוי שבקבוצה אקראית של t אנשים יש שנים עם אותו יום הולדת? בשנה יש n ימים וכולם שווים הסתברות.

$$E \left[\begin{array}{l} \# \text{ Pairs with} \\ \text{same birthday} \end{array} \right] = \sum_{i=1}^t \sum_{j=i+1}^t P[i, j \text{ have same BD}] = \sum_{i=1}^t \sum_{j=i+1}^t \frac{1}{n} = \frac{t(t+1)}{2} \frac{1}{n} \approx \frac{t^2}{2n}$$

תוחלת: $\frac{t^2}{2n}$ כלומר אם $t = \sqrt{2n}$ התוחלת בערך 1.

שתי קבוצות אקראיות מתוך n איברים יחתכו בהסתברות קבועה כאשר גודלן $\theta(\sqrt{n})$.

למשל, אם יש קבוצה שלא יודעים מה גודלה אבל ניתן להוציא ממנה איברים. נוציא איברים עד שנקבל את אותו האיבר פעמיים. אם ניסינו T פעמים, גודל הקבוצה הוא כנראה כ- T^2 .

מה תוחלת מספר ימי ההולדת המשותפים ל t אנשים?

באמצעות משתני אינדיקטור:

$$X_{ij} = \begin{cases} 1 - \text{Person } i \text{ and } j \text{ have same birthday date} & \text{נגדיר} \\ 0 - \text{otherwise} \end{cases}$$

נבחר $1 \leq i < j < t$

$$X = \sum_{ij} X_{ij}$$

$$Pr[X_{ij}=1] = 1/365$$

הסיכוי שאין התנגשות?

הראשון: 1

השני: $\frac{364}{365}$

השלישי: $\frac{363}{365}$

$$Pr[\text{no collision}] = \prod_{i=0}^{t-1} \left(1 - \frac{i}{365}\right) = \prod_{i=0}^{t-1} \left[\left(1 - \frac{i}{365}\right)^{\frac{365-i}{365}} \right]^{\frac{i}{365}} \approx \prod_{i=0}^{t-1} e^{-\frac{i}{365}} = e^{-\left(\sum_{i=0}^{t-1} \frac{i}{365}\right)} = e^{-\left(\frac{t(t-1)}{2 \cdot 365}\right)}$$

ולכן כאשר $t = \sqrt{30}$ הסיכוי קבוע.