

# מבני נתונים ואלגוריתמים - הרצאה 12

6 בדצמבר 2011

## דחיסה

תחילה נעסוק בדחיסה שהיא Lossless - לא מאבדים מידע בדחיסה.

## טענה

כל דחיסה מגדילה חלק מהקבצים.

## הוכחה

נניח שקיימת דחיסה שמטקינה או משאירה שווה כל קובץ:

$$f: L \rightarrow K$$

ותמיד  $k \leq L$

ניקח את  $m$  להיות גודל הקובץ המינימלי שעובר דחיסה ממש.

קיים קובץ בגודל  $m$  שעובר לגודל  $m_1$  כאשר  $m_1 < m$ .

יש  $2^{m_1}$  קבצים אפשריים בגודל  $m_1$ . לאחר הדחיסה יש לפחות  $2^{m_1} + 1$  קבצים אפשריים בגודל  $m_1$ . לפי שובך היונים, לפחות קובץ אחד מופיע פעמיים בקבצים בגודל  $m_1$ . לכן קיימים שני קבצים שונים שנדחסים לאותו קובץ, בניגוד לעובדה שהדחיסה ללא הפסד.

## אי שוויון גיבס

מכאן והלאה -  $p_i, q_i$  הן התפלגויות:

$$\begin{aligned} q_i, p_i &> 0 \\ \sum_i p_i &= 1 \\ \sum_i q_i &= 1 \end{aligned}$$

אי שוויון גיבס אומר: לכל  $p, q$ :

$$-\sum_i p_i \log(p_i) \leq -\sum_i p_i \log(q_i)$$

נשים לב שצד שמאל הוא אנטרופיה:

$$H = -\sum_i p_i \log(p_i)$$

אפשר לכתוב אחרת:

$$\begin{aligned} -\sum_i p_i [\log(p_i) - \log(q_i)] &\leq 0 \\ -\sum_i p_i \log\left(\frac{p_i}{q_i}\right) &\leq 0 \end{aligned}$$

## הוכחה

$$\begin{aligned} -\sum_i p_i \log\left(\frac{p_i}{q_i}\right) &\leq \sum_i p_i \log\left(\frac{q_i}{p_i}\right) \\ &\leq \sum_i p_i \left(\frac{q_i}{p_i} - 1\right) \\ &= \sum_i q_i - \sum_i p_i = 0 \end{aligned}$$

השתמשנו בכך ש  $\ln x \leq x - 1$ :

$$\begin{aligned} f(x) &= x - 1 - \ln x \\ f'(x) &= 1 - \frac{1}{x} = 0 \\ x &= 1 \\ f(1) &= 0 \\ f(x) &\geq 0 \\ \ln x &\leq x - 1 \end{aligned}$$

לכן אי השוויון נכון.  
אם נגדיר

$$D(p, q) = \sum_i p_i \log\left(\frac{p_i}{q_i}\right) \geq 0$$

ומתקיים:

$$D(p, q) = 0 \iff p = q$$

לכן זו סמי-מטריקה (זו לא מטריקה כי אין סימטריות).  
 $D(p, q)$  נקרא מרחק קולבק-לייבלר.  
יש גם מקבילה רציפה למרחק זה:

$$D(p, q) = \int p(x) \log\left(\frac{p(x)}{q(x)}\right) dx$$

## אי שוויון Kraft - קוד רישא

קוד רישא הוא קוד שבו אף מילה אינה רישא של מילה אחרת, שבאלפבית שלו יש  $D$  אותיות.  
בהינתן קוד רישא יש לכל מילה אורך  $\ell_i$ .  
אי שוויון Kraft הוא:

$$\sum_i D^{-\ell_i} \leq 1$$

## הוכחה

נכתוב את קוד הרישא בעץ - מילים הן עלים.  
נסמן ב  $k$  את עומק העץ.  
נמלא את העץ עד עומק  $k$ , כלומר נוסיף לעלים בנים עד עומק  $k$ , ו"נסמן" את הבנים שהוספנו ושלא באמת שייכים לעץ המקורי.  
מתחת לכל מילה אורך  $\ell_i$  יש  $D^{k-\ell_i}$  עלים.  
יש בסה"כ  $D^k$  עלים בעץ (כי זה עץ מלא בעומק  $k$  ולכל צומת יש  $D$  בנים).  
נעבור על כל המילים וניקח את כל העלים שיש מתחת אליהן בעץ, וזה בהכרח לכל היותר מס' העלים בעץ:

$$\sum_i D^{k-\ell_i} \leq D^k$$

נחלק ב- $D^k$  ונקבל:

$$\sum_i D^{-\ell_i} \leq 1$$

משל.

## טענה

יהי  $\ell_i$  אורך המילים בקוד רישא, אז:

$$\sum_i p_i \ell_i \geq H$$

כאשר  $H$  היא האנטרופיה.

## הוכחה

נגדיר

$$z = \sum_i 2^{-\ell_i} \leq 1$$
$$q_i = \frac{2^{-\ell_i}}{z}$$

נשים לב ש- $q_i$  היא התפלגות:

$$q_i \geq 0$$
$$\sum_i q_i = \frac{\sum_i 2^{-\ell_i}}{z} = \frac{z}{z} = 1$$

נעביר אגפים ונקבל:

$$2^{-\ell_i} = z \cdot q_i$$
$$\ell_i = -\log_2 z - \log_2 q_i$$

אזי

$$\begin{aligned} \sum_i p_i \ell_i &= -\sum_i p_i \log_2 z - \sum_i p_i \log_2 q_i \\ &= -\log_2 z - \sum_i p_i \log_2 q_i \\ &\geq 0 - \sum_i p_i \log_2 p_i = H \end{aligned}$$

(המעבר לשורה האחרונה הוא לפי Gibbs Kraft).  
לכן בעצם אלגוריתם הדחיסה הכי טוב בקוד רישא יידחס עד לגבול  $H$ .  
נשים לב שאם אנו רוצים שהאלגוריתם יפעל כך, אנו צריכים שיתקיים:

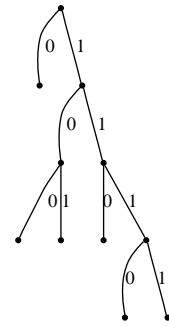
$$\log_2(z) = 0$$
$$p_i \sim q_i$$

כלומר התנאים שנרצה שיתקיימו הם:

$$\log_2(z) = 0$$
$$p_i \sim \frac{2^{-\ell_i}}{z}$$

## קוד הופמן

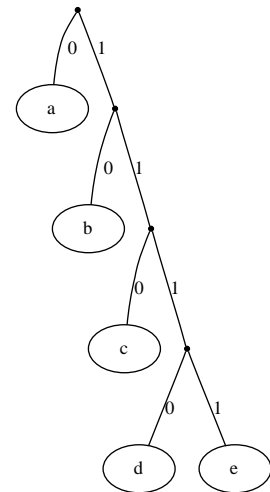
נשים לב שאם יש לנו עץ שלם (כלומר לא קיים צומת עם בן יחיד), יתקיים בהכרח  $\log_2 z = 0$ , כי  $\sum 2^{-\ell_i} = 1$



קוד הופמן תמיד יתן לנו עץ שלם, ולכן התנאי הראשון תמיד יתקיים בו. נניח לדוגמה שאלה המילים והשכיחויות:

שכיחות	מילה
$\frac{1}{2}$	a
$\frac{1}{4}$	b
$\frac{1}{8}$	c
$\frac{1}{16}$	d
$\frac{1}{16}$	e

העץ שנקבל מקוד הופמן הוא:



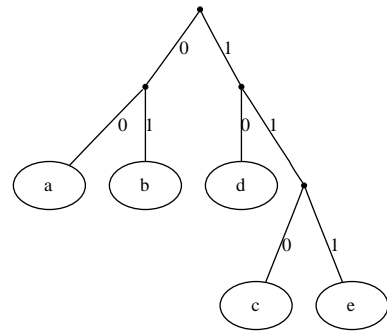
אז נקבל:

$\ell_i$	המילה בקוד הופמן	שכיחות	מילה
1	0	$\frac{1}{2}$	a
2	10	$\frac{1}{4}$	b
3	110	$\frac{1}{8}$	c
4	1110	$\frac{1}{16}$	d
4	1111	$\frac{1}{16}$	e

אזי במקרה הזה  $z = 2^{-1} + 2^{-2} + 2^{-3} + 2 \cdot 2^{-4} = 1$  וזה יתקיים בכל מקרה בקוד הופמן. אך התנאי השני לא בהכרח יתקיים, למרות שפה הוא כן מתקיים. אם השכיחויות לא יהיו חזקות של 2 או משהו קרוב אליהן, יכול להתקבל משהו שונה:

שכיחות	מילה
$\frac{1}{2}$	a
$\frac{1}{4}$	b
$\frac{1}{8}$	c
$\frac{1}{8}$	d
0	e

עץ הופמן שנקבל:



ואז:

מילה	שכיחות	קוד הופמן	אורך
a	$\frac{1}{3}$	00	2
b	$\frac{1}{3}$	01	2
c	$\frac{1}{3}$	110	3
d	$\frac{2}{3}$	10	2
e	0	111	3

ובמקרה זה התנאי השני לא מתקיים.