



האוניברסיטה
העברית
בירושלים
THE HEBREW
UNIVERSITY
OF JERUSALEM

אלגוריתמים בביו' חישובית

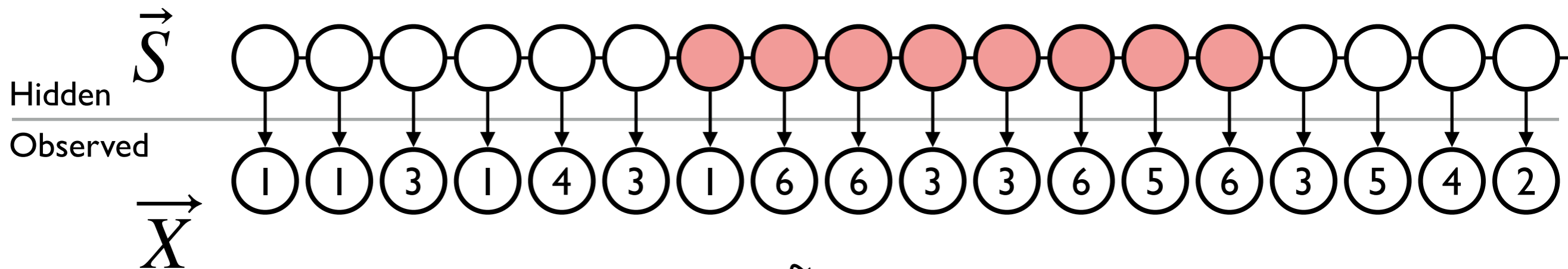
76558

זיהוי מוטיבים בעזרת
מודלים מרקוביים חבויים

תומי קפלן
30/1/2024

שרשראות מרקוב חבויות

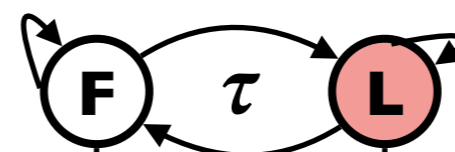
- משתני התהליך המרקוביים הם חבויים
- פולטים (בהתאם למצבם) משתנים אחרים (ניצפים)



$$S_i \in \{F, L\}$$

$$\tau = P(S_i | S_{i-1})$$

τ	F	L
F	99%	1%
L	5%	95%



$$\forall_k, e_F(k) = \frac{1}{6} \begin{matrix} 1/6 \\ 1/6 \\ 1/6 \\ 1/6 \\ 1/6 \\ 1/6 \end{matrix} \quad e_L(k) = \begin{cases} \frac{1}{10} & k = 1, \dots, 5 \\ \frac{1}{2} & k = 6 \end{cases}$$

$$X_i \in \{1, 2, 3, 4, 5, 6\}$$

מה עוד היינו רוצים לדעת?

$$P(S_i | X_1, \dots, X_n)$$

- פוסטריוור על מצב חבוי

האם סביר שעברנו במצב מסויים בזמן i

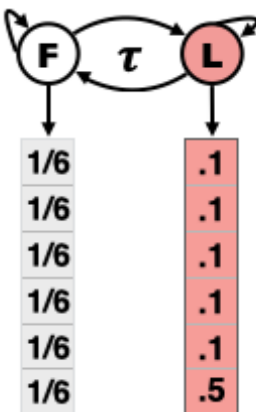
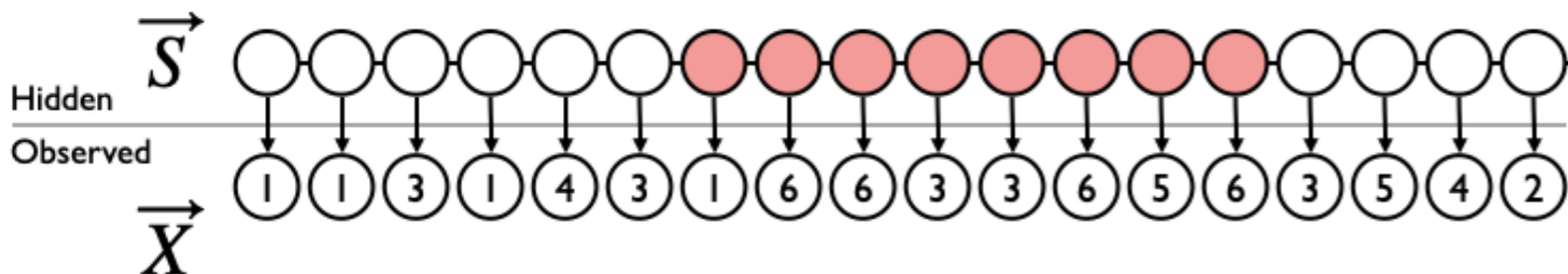
$$P(S_{i-1}, S_i | X_1, \dots, X_n)$$

- פוסטריוור על מעבר חבוי

$$\arg \max_{S_1, \dots, S_n} P(S_1, \dots, S_n, X_1, \dots, X_n)$$

- מסלול ויטרבי

MPE / decoding



חישוב ניראות: אלגוריתם Forward

$$F_k(i) = P(X_1, \dots, X_i, S_i = k)$$

• נגדיר:

$$F_l(i+1) = P(X_1, \dots, X_i, X_{i+1}, S_{i+1} = l)$$

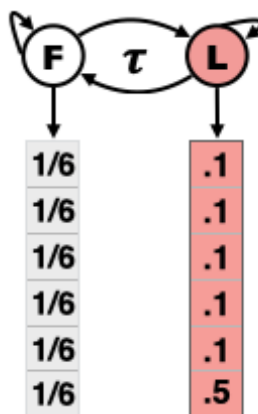
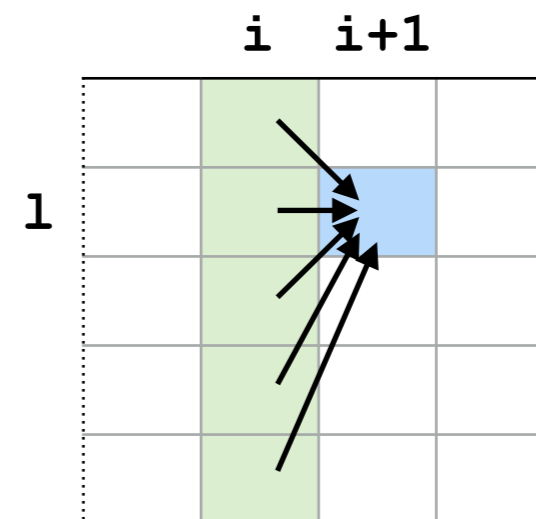
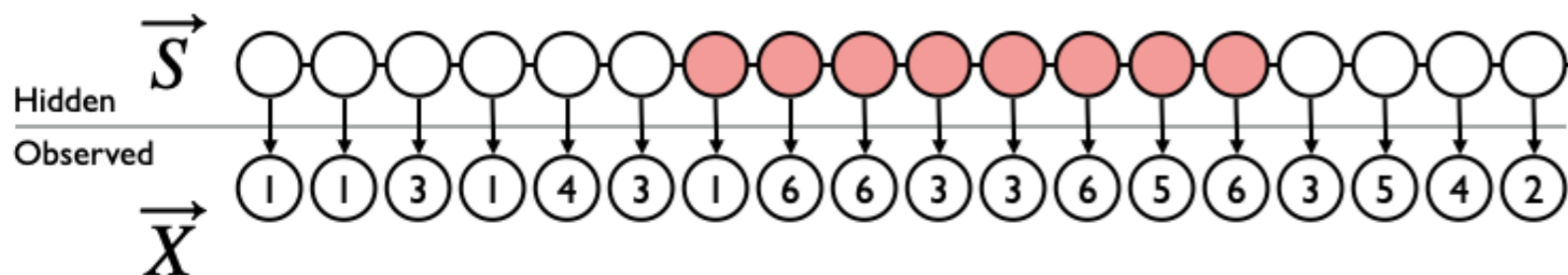
• נפתח רקורסיבית את

$$= \sum_k P(X_1, \dots, X_i, X_{i+1}, S_i = k, S_{i+1} = l)$$

$$= \sum_k P(X_1, \dots, X_i, S_i = k) \cdot P(S_{i+1} = l | X_1, \dots, X_i, S_i = k) \cdot P(X_{i+1} | X_1, \dots, X_i, S_i = k, S_{i+1} = l)$$

$$= \sum_k [P(X_1, \dots, X_i, S_i = k) \cdot P(S_{i+1} = l | S_i = k)] \cdot P(X_{i+1} | S_{i+1} = l)$$

$$= \sum_k [F_k(i) \cdot \tau_{k,l}] \cdot e_l(X_{i+1})$$



חישוב ניראות: אלגוריתם Backward

$$B_k(i) = P(X_{i+1}, \dots, X_n | S_i = k)$$

• נגדיר:

$$B_l(i-1) = P(X_i, X_{i+1}, \dots, X_n | S_{i-1} = l)$$

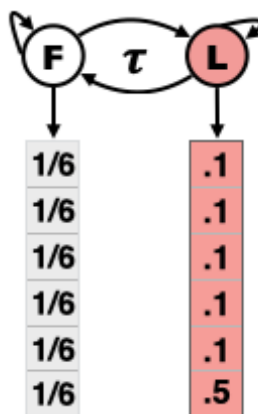
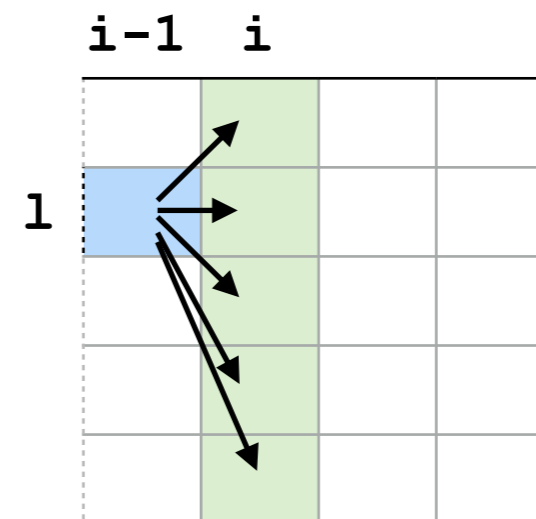
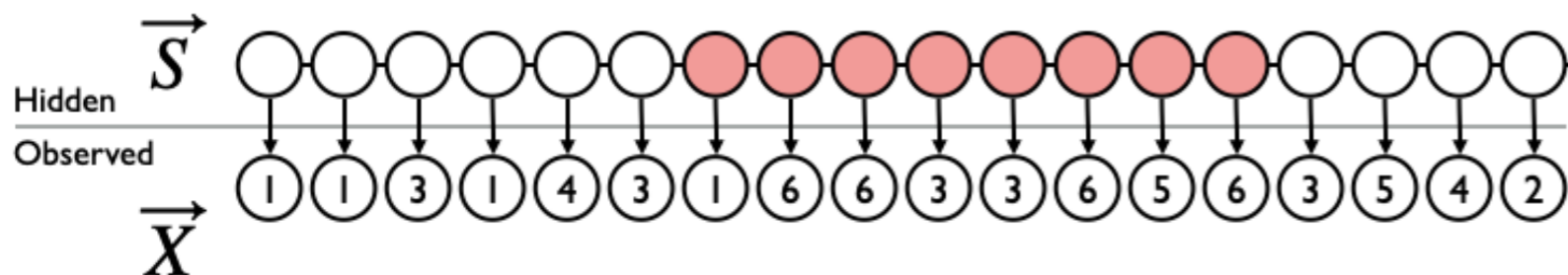
• נפתח רקורסיבית את

$$= \sum_k P(S_i = k, X_i, X_{i+1}, \dots, X_n | S_{i-1} = l)$$

$$= \sum_k P(S_i = k | S_{i-1} = l) \cdot P(X_i | S_i = k, S_{i-1} = l) \cdot P(X_{i+1}, \dots, X_n | X_i, S_i = k, S_{i-1} = l)$$

$$= \sum_k P(S_i = k | S_{i-1} = l) \cdot P(X_i | S_i = k) \cdot P(X_{i+1}, \dots, X_n | S_i = k)$$

$$= \sum_k \tau_{l,k} \cdot \mathbf{e}_k(X_i) \cdot B_k(i)$$



חישוב הפוסטריוור

$$P(\vec{X}, S_i = k) = P(X_1, \dots, X_n, S_i = k)$$

$$= P(X_1, \dots, X_i, S_i = k) \cdot P(X_{i+1}, \dots, X_n | X_1, \dots, X_i, S_i = k)$$

$$= P(X_1, \dots, X_i, S_i = k) \cdot P(X_{i+1}, \dots, X_n | S_i = k)$$

$$= F_k(i) \cdot B_k(i)$$

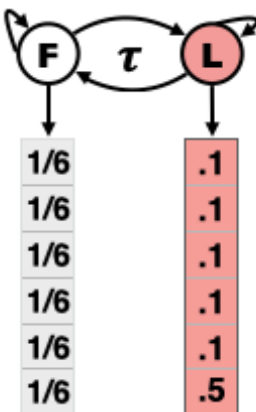
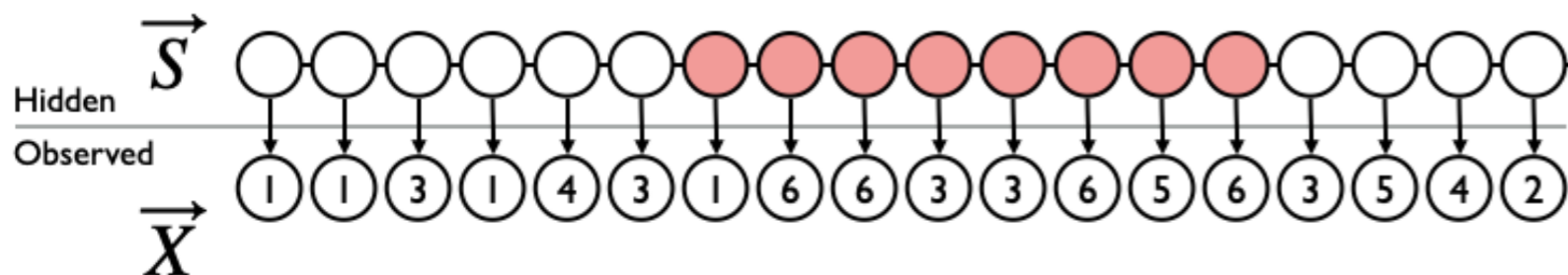
• ומכאן ש:

$$P(S_i = k | \vec{X}) = \frac{F_k(i) \cdot B_k(i)}{P(\vec{X})}$$

• ונוכל לחשב את הפוסטריוור:

$$P(S_{i-1} = l, S_i = k | \vec{X}) = \frac{F_{i-1}(l) \cdot \sum_{\ell k} e_{\ell k}(x_i) \cdot B_i(k)}{P(\vec{X})}$$

גודל ביג:
פוסטריוור
לש בייג.
טאן גוליה:
הטלוני:



מציאת מסלול ויטרבי

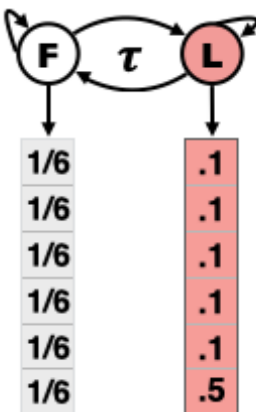
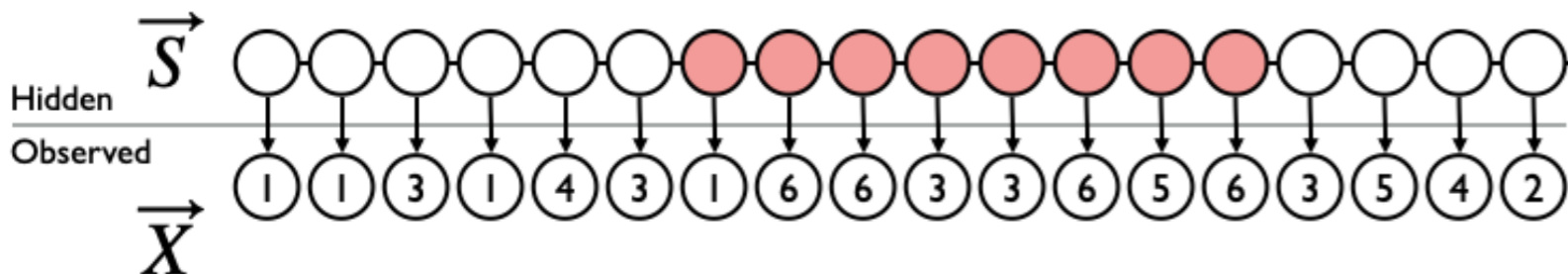
- איך נמצא את המסלול הסביר ביותר, בהנתן התצפיות?

$$\pi^* = \arg \max_{\vec{S}} P(S_1, \dots, S_n | X_1, \dots, X_n) = \arg \max_{\vec{S}} P(\vec{S}, \vec{X})$$

למה?

- אלגוריתם תכנון דינאמי. נשמור בטבלה V את ההסתברות למסלול הסביר ביותר (עד i) המסתיים במצב $S_i=k$

$$V_k(i) = \max_{S_1, \dots, S_{i-1}} P(S_1, \dots, S_{i-1}, S_i = k, X_1, \dots, X_i)$$



מציאת מסלול ויטרבי

$$V_k(n) = \max_{S_1, \dots, S_{n-1}} P(S_1, \dots, S_{n-1}, S_n = k, X_1, \dots, X_n)$$

• לפיכך

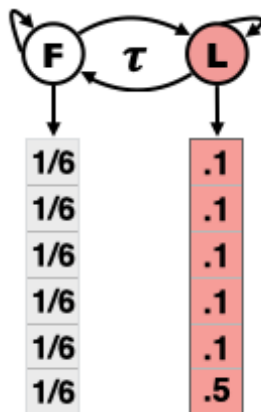
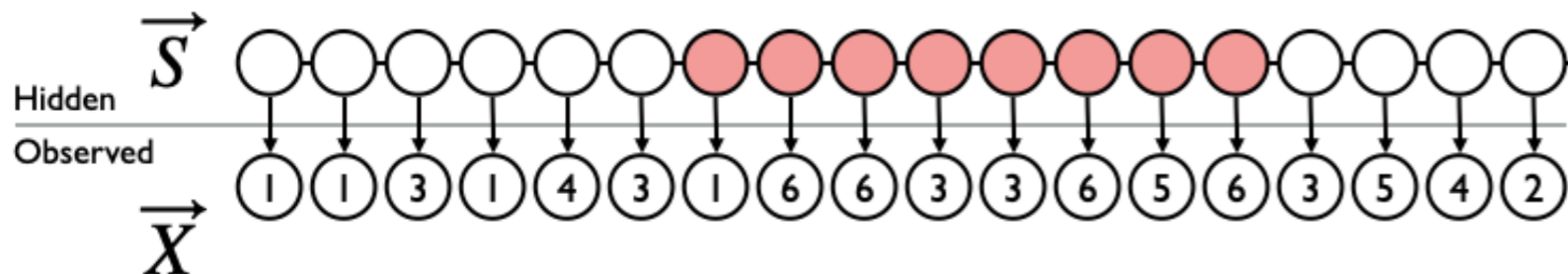
$$\max_k V_k(n) = \max_k \max_{S_1, \dots, S_n} P(S_1, \dots, S_n = k, X_1, \dots, X_n) = \max_{\vec{S}} P(\vec{S}, \vec{X})$$

• ולכן

• כלומר נמצא את המצב האחרון במסלול האופטימלי

$$\pi^*[n] = \arg \max_k V_k(n)$$

• באופן דומה נוכל לשמור טבלה נוספת, לשיחזור המסלול



מציאת מסלול ויטרבי

Init: $\forall_k V_k(0) = 0$

For $i = 1 \dots n$

For $k = 1 \dots K$

$$V_k(i) = \mathbf{e}_k(X_i) \cdot \max_l [v_l(i-1) \cdot \tau_{l,k}]$$

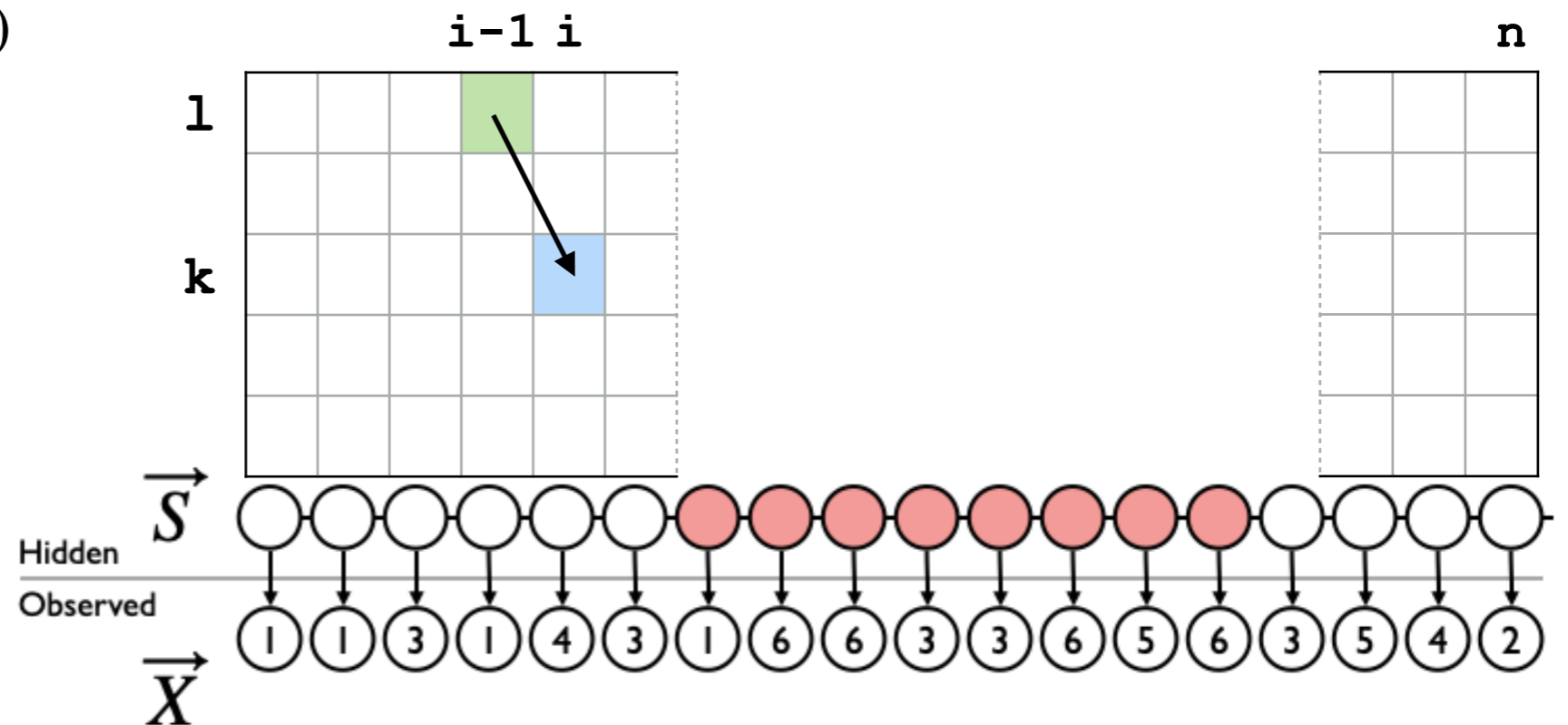
$$Ptr_k(i) = \arg \max_l [v_l(i-1) \cdot \tau_{l,k}]$$

Traceback:

$$\pi^*[n] = \arg \max_k V_k(n)$$

For $i = n \dots 2$

$$\pi^*[i-1] = Ptr_{\pi^*[i]}(i)$$



דגימה ממרחב הפוסטריור

- האם נדע לדגום מסלולים (שלמים) מהסתברות הפוסטריור?
 [שימו לב להבדל, לעומת הפוסטריור המקומי $P[S_i|X]$]

נפתח:

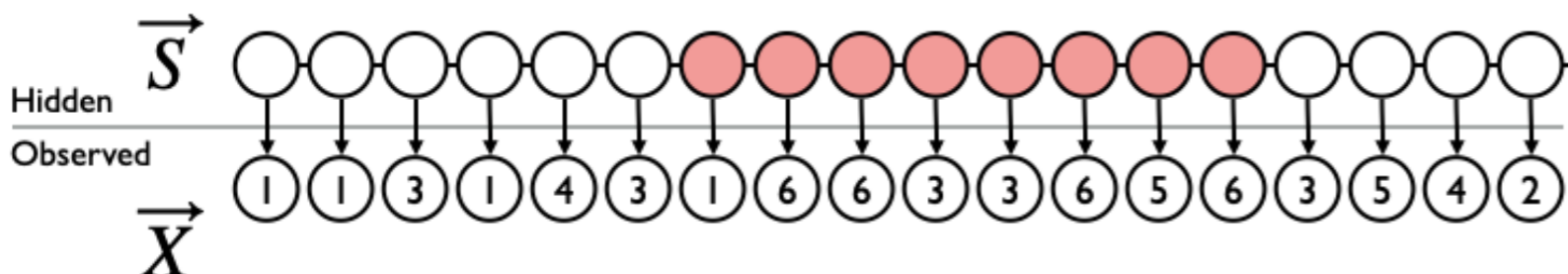
$$P(\vec{S} | \vec{X}) = \prod_i P(S_i | S_{i-1}, \vec{X})$$

$$= \prod_i P(S_i | S_{i-1}, X_1, \dots, X_i, \dots, X_n) = \prod_i \frac{P(X_i, \dots, X_n, S_{i-1}, S_i)}{P(X_i, \dots, X_n, S_{i-1})}$$

$$= \prod_i \frac{P(X_i | S_i) \cdot P(X_{i+1}, \dots, X_n | S_{i-1}, S_i) \cdot P(S_i | S_{i-1})}{P(X_i, \dots, X_n, S_{i-1})}$$

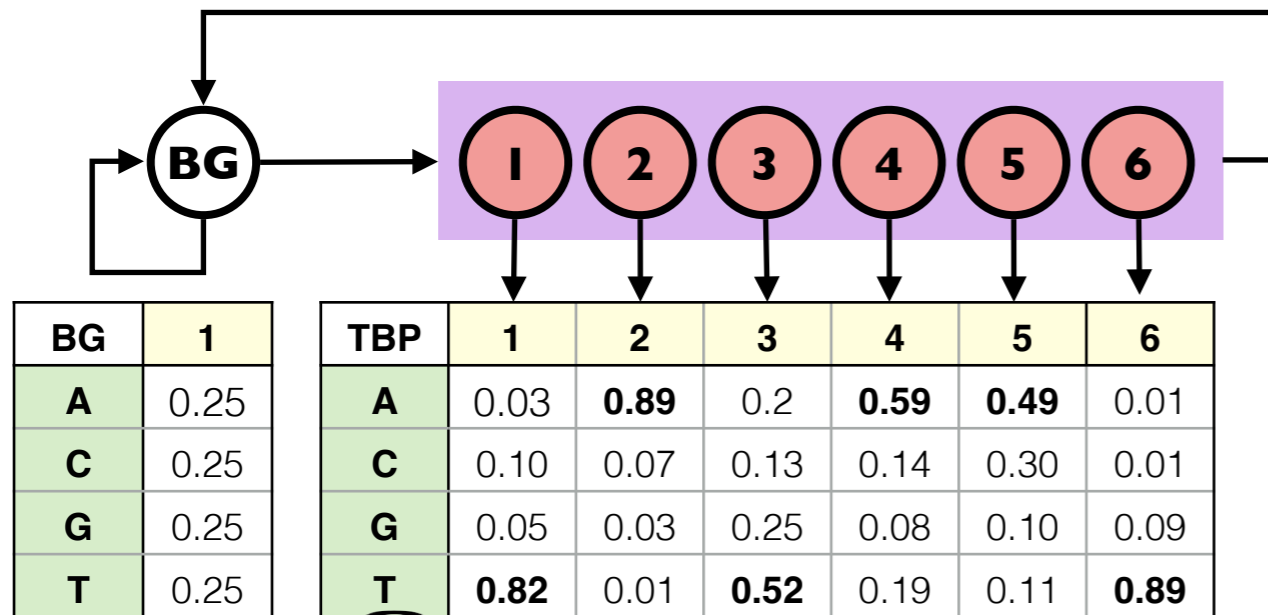
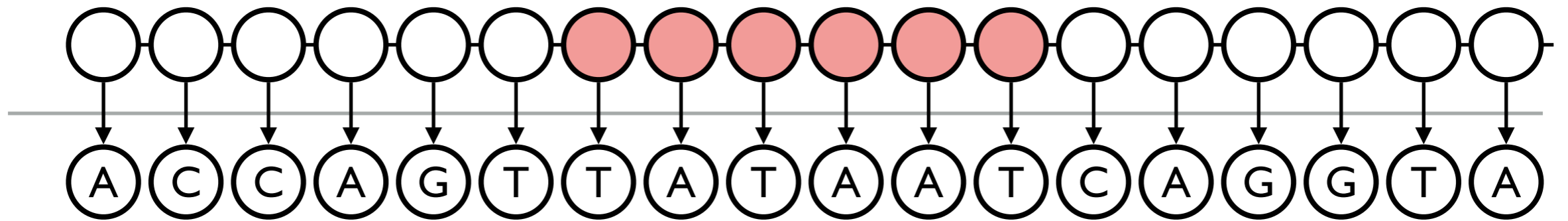
- כלומר נוכל לדגום עמדה עמדה:

$$P(S_i | S_{i-1}, \vec{X}) \propto e_{S_i}(X_i) \cdot B_{S_i}(i) \cdot \tau_{S_{i-1}, S_i}$$



זיהוי מוטיבים עם מודלים חבויים

- נניח שנתונים רצפי דנ"א, וברצוננו למצוא בהם מוטיב ידוע (למשל TATA-box, אתר הקישור למכונת השעתוק)



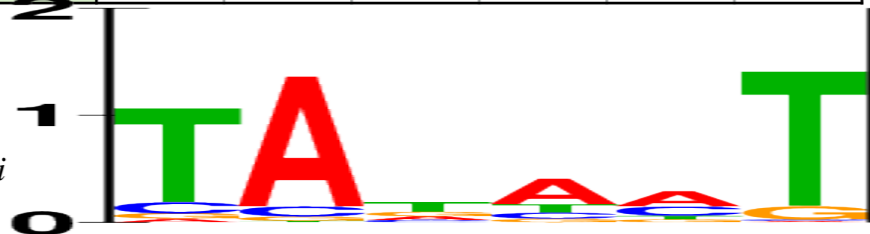
GCGGCACCCATCTCTTATAAATTCGCTTGATAGTAACGTTTCGAAAGCACAA

AAAGCGGTAGACTTCAGGCATAAAAGGATTAATTTTGGACAATCCCCGAT

CATCCACAGTTAAGTTATCCGAAGGGGGTTCGGCATCGGAAGCTTGAAATT

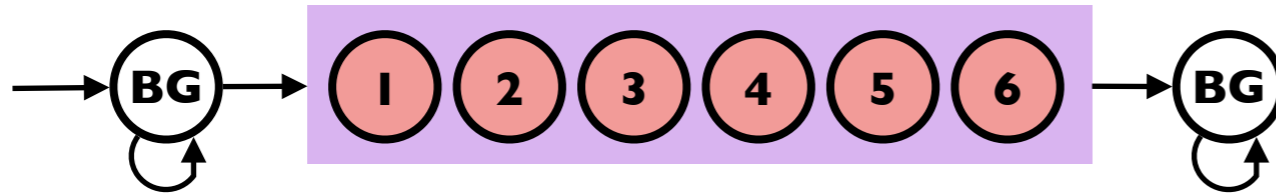


$$H(p) = - \sum_i p_i \log_2 p_i$$

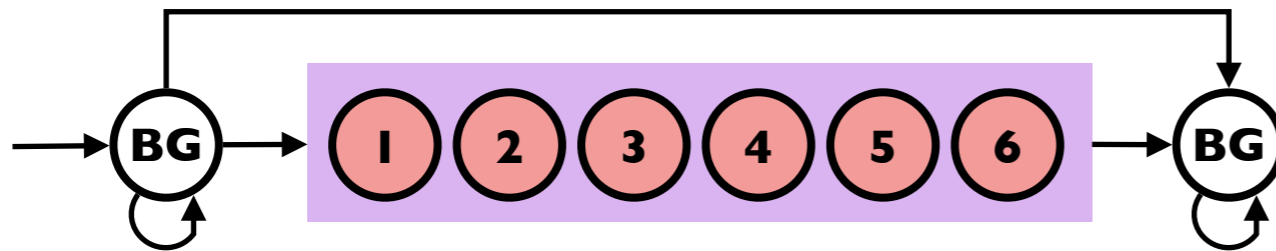


זיהוי מוטביבים עם מודלים חבויים

• מודל OOPS



• מודל ZOOPS



• מודל TCM

