



האוניברסיטה  
העברית  
בירושלים  
THE HEBREW  
UNIVERSITY  
OF JERUSALEM

# אלגוריתמים בביו' חישובית

76558

שערוך פרמטרים ושיטות

היוריסטיות

תומי קפלן

14/1/2024

# ציון העימוד ככלל החלטה אופטימלי

- שתי ההיפותיזות שלנו:

$H_0$ : אין קשר בין שני הרצפים

$H_1$ : לשני הרצפים יש אב קדמון משותף

- תחת הנחת אי-תלות בין העמודות:

$$P_1(S, T) = \prod_i P_1(S_i, T_i), \quad P_0(S) = \prod_i P_0(S_i)$$

- לוג יחס הנראות מתפרק לסכום ציוני העמודות:

$$LLR = \log \frac{P_1(S, T)}{P_0(S) \cdot P_0(T)} = \log \frac{\prod P_1(S_i, T_i)}{\prod P_0(S_i) \cdot \prod P_0(T_i)} = \sum_i \log \frac{P_1(S_i, T_i)}{P_0(S_i) \cdot P_0(T_i)} \triangleq \sum_i \sigma(S_i, T_i)$$

# כיצד נשערך את $P_1$ ואת $P_0$ ?

$$P_1(a, b) = P_1(b, a)$$

• ובפרט, תחת אילוצי סימטריה:

$$P_0(a) = \sum_b P_1(a, b)$$

• ונוסחת ההסת' השלמה:

# שערוך פרמטרים פשוט

• נניח מטבע עם התפלגות ברנולי:  $\Omega = \{H, T\}$

$$P_\theta = \Omega \rightarrow [0,1]$$

• וסדרת תצפיות  $D = \{a_1, a_2, a_3, \dots, a_n\}$

• הגדרת הניראות:

$$L(\theta : D) = P_\theta(D) = \prod_i^n P_\theta(a_i)$$

# שערוך פרמטרים פשוט

$$D = \{H, H, T, H\}$$

• ספציפית, אם

$$P_\theta(H) = \theta, \quad P_\theta(T) = 1 - \theta$$

• ונגדיר

$$L(\theta : D) = \theta^3(1 - \theta)$$

• נקבל

$$L(\theta : D) = \theta^{N_H(D)}(1 - \theta)^{N_T(D)}$$

• ובאופן כללי

$$N_H(D) = \sum_i 1\{a_i = H\}$$

$$N_T(D) = \sum_i 1\{a_i = T\}$$

# סטטיסטיים מספיקים

- במקרה שלנו,  $N_H$  ו- $N_T$  הם נתונים מספיקים  
[Sufficient Statistics] לחישוב הנדאות

$\vec{S}(D)$

- כלומר, בהינתן וקטור הססט' המספיקים  
הנדאות כבר אינה תלויה בסדרה  $D$ :

$\forall_{D_1}, \forall_{D_2}, \forall_{\theta} :$

$$\vec{S}(D_1) = \vec{S}(D_2) \Rightarrow L(\theta : D_1) = L(\theta : D_2)$$

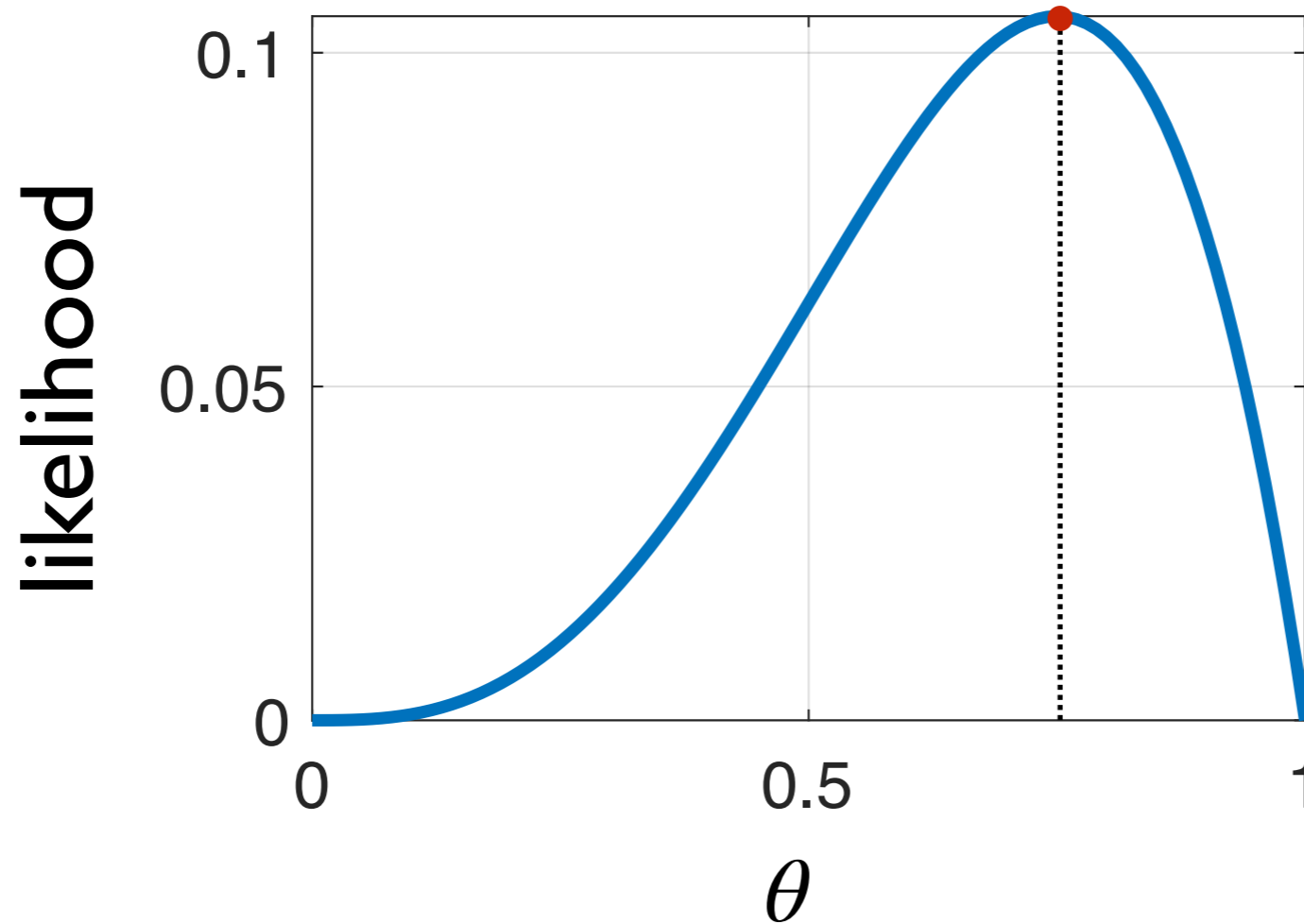
# Maximum Likelihood Estimator

$$L(\theta : D) = \theta^3(1 - \theta)$$

$$\hat{\theta}_{ML} = \arg \max_{\theta} L(\theta : D)$$

# Maximum Likelihood Estimator

$$L(\theta : D) = \theta^3(1 - \theta)$$



$$\hat{\theta}_{ML} = \arg \max_{\theta} L(\theta : D)$$

• איך נמצא? נגזור ונשווה לאפס. קלללל.



# מציאת ערך הניראות המקסימלי

$$L(\theta : D) = \theta^{N_H} \cdot (1 - \theta)^{N_T}$$

- אהה! לוג היא פונ' מונוטונית.
- המקסימום של  $L$  מושג באותה הנקודה (= ארגמקס) בה מושג המקסימום של  $\log L$

$$\log L(\theta : D) = N_H \cdot \log \theta + N_T \cdot \log(1 - \theta)$$

- נגזור ונקבל:  $LL' = \frac{N_H}{\theta} + \frac{N_T}{1 - \theta} \cdot (-1)$

# מציאת ערך הניראות המקסימלי

$$LL' = \frac{N_H}{\theta} - \frac{N_T}{1 - \theta}$$

• נשווה לאפס

$$\frac{N_H}{\theta} = \frac{N_T}{1 - \theta}$$

$$(1 - \theta) \cdot N_H = \theta \cdot N_T$$

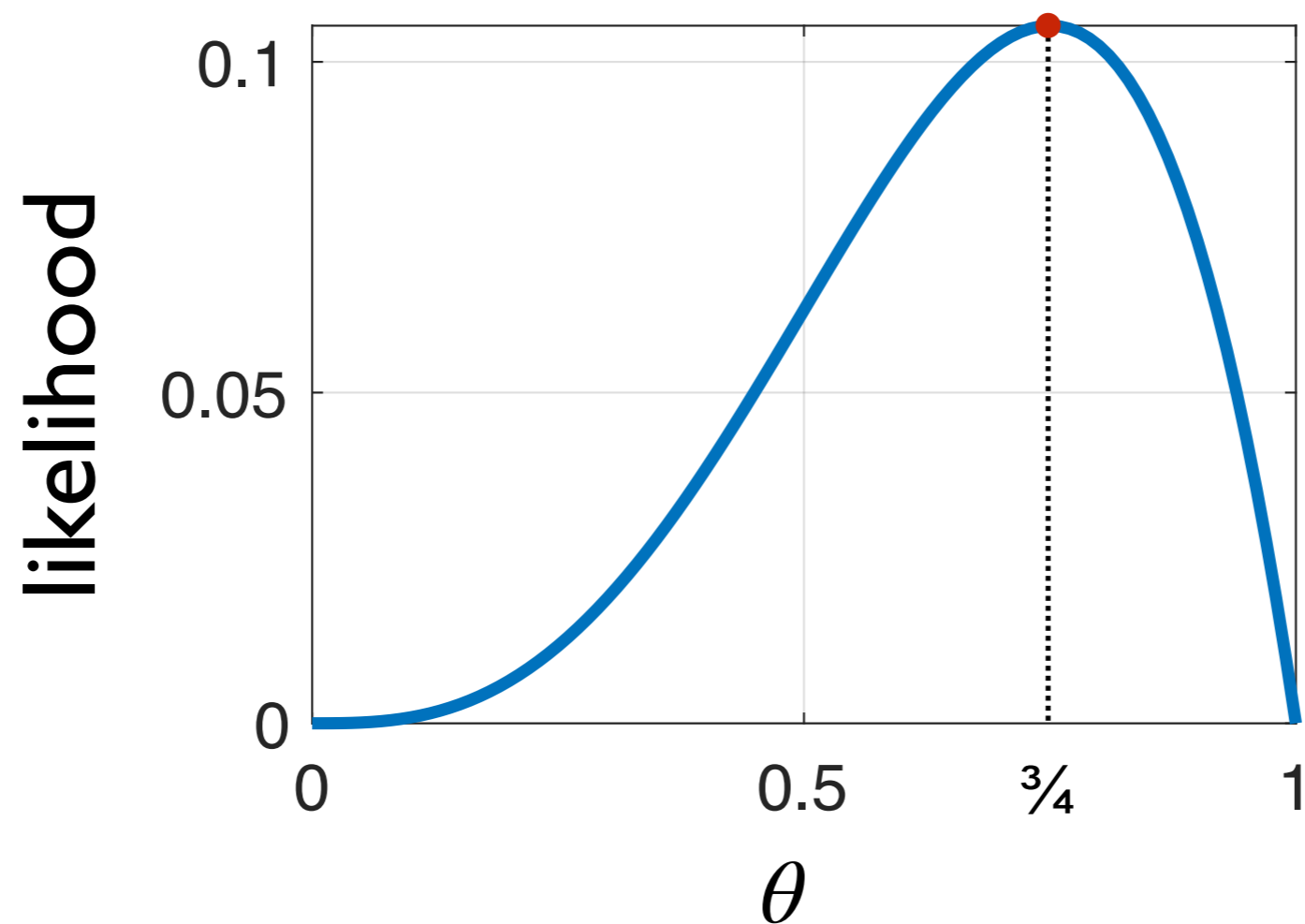
$$N_H - \theta \cdot N_H = \theta \cdot N_T$$

$$N_H = \theta \cdot (N_T + N_H)$$

• כלומר, אומד הניראות המקסימלי:  $\hat{\theta} = \frac{N_H}{N_T + N_H}$

# אומד ניראות מקסימלי

$$L(\theta : D) = \theta^3(1 - \theta)$$



$$\hat{\theta}_{ML} = \frac{N_H}{N_H + N_T} = \frac{3}{4}$$

• ועבור המטבע נקבל:

# נכליל להתפלגות מולטינומיאלית

$$\vec{\theta} = \langle \theta_1, \theta_2, \dots, \theta_k \rangle \in \mathbb{R}^k \quad \theta_i \geq 0, \quad \sum_i \theta_i = 1$$

• והנידאות:

$$L(\theta : D) = \prod_i^n \theta_{a_i} = \prod_i^k \theta_i^{N_i}$$

• לוג הנידאות:

$$\log L(\theta : D) = \sum_i^k N_i \log \theta_i$$

נקודת המקסימום  
בכפוף לאילוצים!

• מיהו אומד הנידאות המירבית?

# שיטת כופלי לגראנז'

$$\arg \max_{\vec{X} \in \mathbb{R}^k} f(\vec{X}) \quad \text{s.t.} \begin{cases} c_1(\vec{X}) = 0 \\ \vdots \\ c_l(\vec{X}) = 0 \end{cases} \quad \text{מציאת אופטימום בכפוף לאילוצים:}$$

$$J(\vec{X}, \vec{\lambda}) = f(\vec{X}) - \sum_i c_i(X) \cdot \lambda_i \quad \bullet \text{ נוסף } \lambda \text{ כופלי לגראנז':}$$

$$\forall_i c_i(X) = 0 \Rightarrow J(\vec{X}, \vec{\lambda}) = f(\vec{X}) \quad \bullet \text{ כאשר } X \text{ עומד בתנאים:}$$

# כופלי לגראנז' להתפלגות מולטינומית

$$\arg \max_{\vec{\theta} \in \mathbb{R}^k} f(\vec{\theta})$$

$$\sum_i \theta_i = 1 \quad \text{אילוצים}$$

$$J = \sum_i^k N_i \log \theta_i - \lambda \left( \sum_i \theta_i - 1 \right)$$

•  $\lambda$  כופלי לגראנז':

$$\frac{\partial J}{\partial \theta_i} = \frac{N_i}{\theta_i} - \lambda$$

$$\frac{\partial J}{\partial \theta_i} = 0 \Rightarrow \hat{\theta}_i = \frac{N_i}{\lambda}$$

$$\frac{\partial J}{\partial \lambda} = - \sum_i \theta_i - 1$$

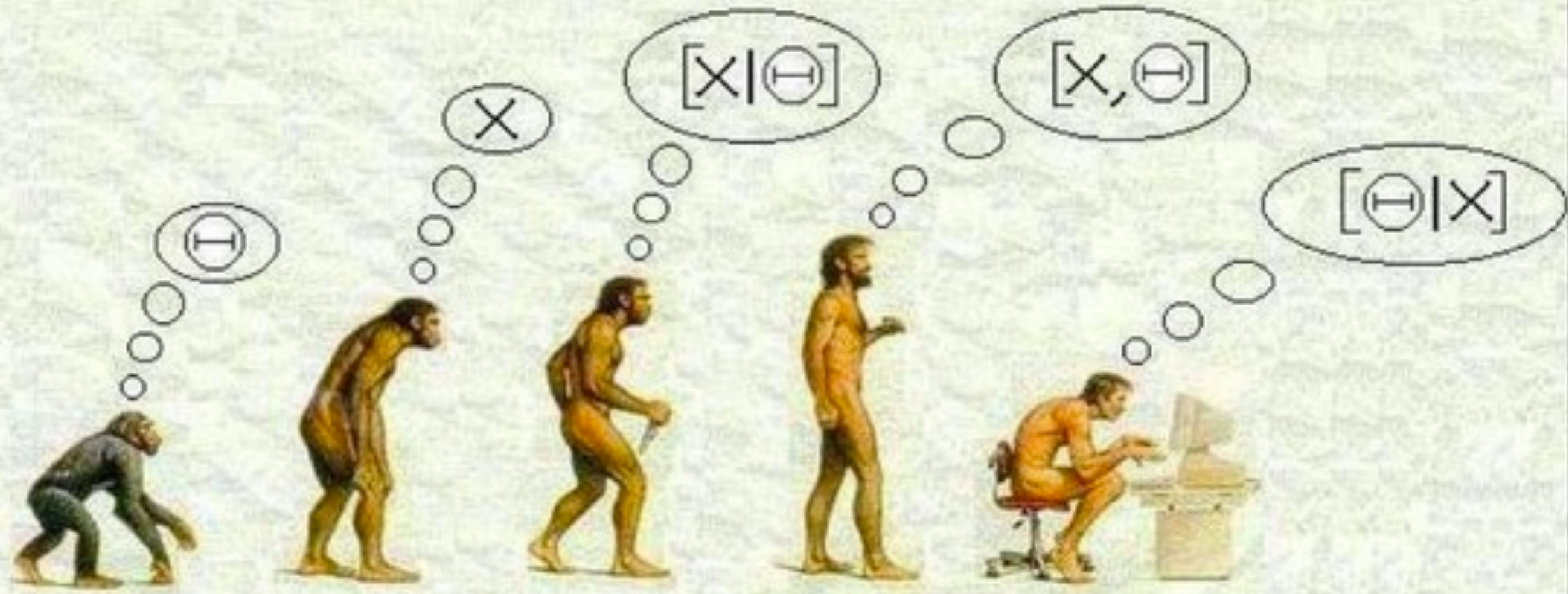
$$\frac{\partial J}{\partial \lambda} = 0 \Rightarrow \sum_i \hat{\theta}_i = 1$$

$$\hat{\theta}_i = \frac{N_i}{N} \quad \begin{array}{l} \text{אומד הנדאות} \\ \text{המידבית} \end{array}$$

$$\sum_i \hat{\theta}_i = \frac{\sum_i N_i}{\lambda} = 1 \Rightarrow \lambda = \sum_i N_i = N$$

# נושאים לדיון

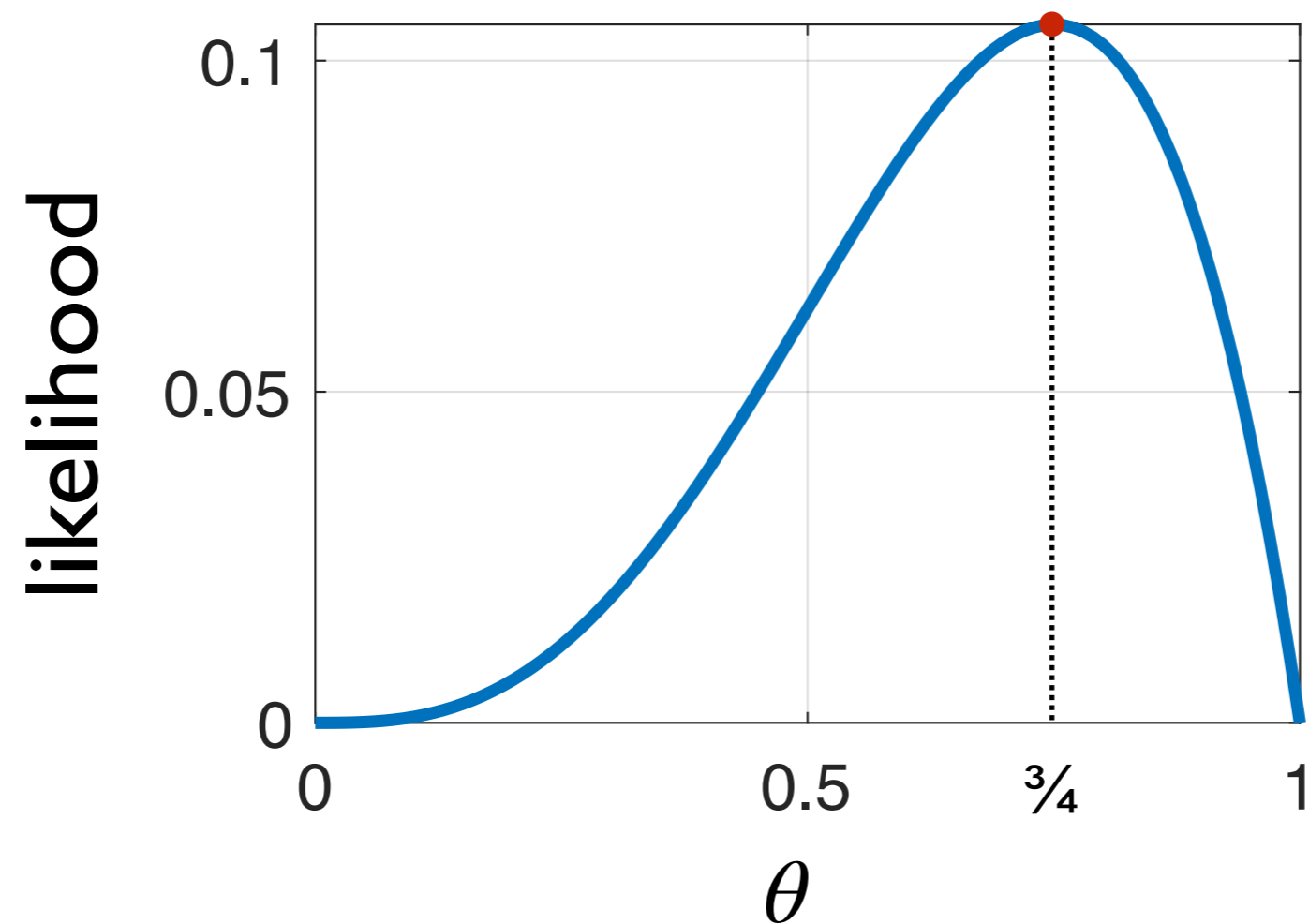
(YET ANOTHER) HISTORY OF LIFE AS WE KNOW IT...



**HOMO APRIORIUS**      **HOMO PRAGMATICUS**      **HOMO FREQUENTISTUS**      **HOMO SAPIENS**      **HOMO BAYESIANIS**

# שערוך בייזיאני

- במקום אומד  $MLE$ , נשערוך אומד שמתחשב בכל ערכי התטא האפשריים (כל אחד לפי הפוסטריוור שלו)





# שערוך בייזיאני

- במקום אומד MLE, נשערוך אומד שמתחשב בכל ערכי התטא האפשריים [כל אחד לפי הפוסטריוור שלו]

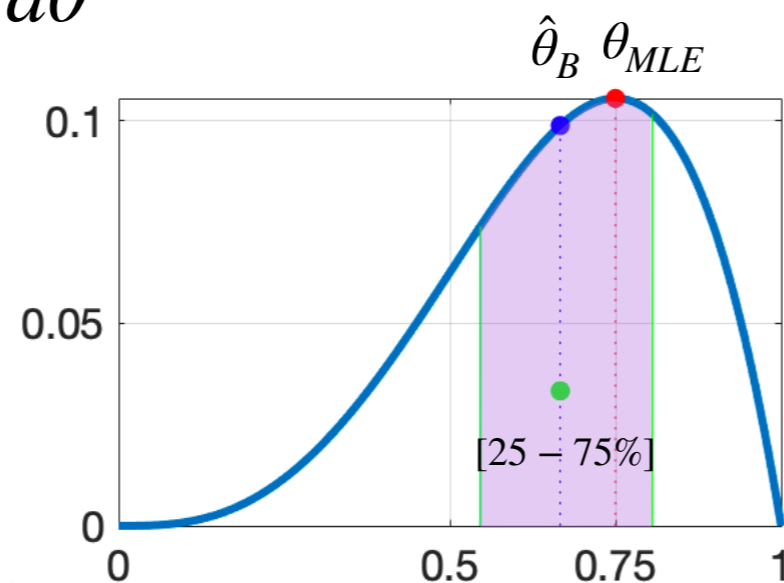
$$\hat{\theta} = \mathbb{E}(\theta | D) = \int \theta \cdot P(\theta | D) d\theta$$

$$P(\theta | D) = \frac{P(D | \theta) \cdot P(\theta)}{P(D)} \propto \underbrace{P(D | \theta)}_{\text{ניראות}} \cdot \underbrace{P(\theta)}_{\text{פריור}}$$

- לפי חוק ביים:

$$\hat{\theta} = \int \theta \cdot \frac{P(D | \theta)}{P(D)} \cancel{P(\theta)} d\theta$$

- עם פריור יוניפורמי:



# שערוך $P_1$ ו- $P_0$

- לפי למת נימון-פירסון, הסקוד האופטימלי לעימוד

$$\sigma(a, b) \triangleq \log \frac{P_1(a, b)}{P_0(a) \cdot P_0(b)} \quad (+c) \quad ?$$

רצפים צריך להיות

- הסיכוי למצוא  $a$  בחלבון מקרי  $P_0(a)$

- הסיכוי למצוא  $a$  מועמד מול  $b$   $P_1(a, b)$

[בעמדות מתאימות בחלבונים בעלי אב קדמון משותף]

- אומד MLE  
עבוד דאטה-בייס של חלבונים?  
 $\hat{P}_0(a) = \frac{N_a}{N}$

שערוך  $P_0$  ו-  $P_1$ 

- רגע רגע - אבל זה תלוי במרחק האבולוציוני ביניהם!

איזור הדמדומים

זהות מוחלטת

0%

15%

40%

62%

80%

100%

BLOSUM

 $\sigma$  שואפת ל- $P_0$  $\sigma$  אלכסונית

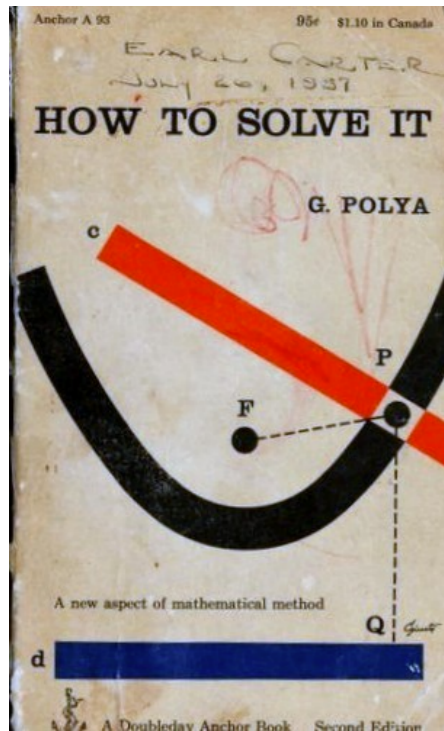
## Henikoff &amp; Henikoff [92']

1. נעמד זוגות חלבונים [לפי מרחק עריכה פשוט]
2. נמצא איזורים בעלי % זהות עד 62%
3. נחשב עליהם סטטיסטיקות:

$$\text{BLOSUM}_{62} \triangleq \left[ \log \frac{P_{62}(a, b)}{P_{62}(a) \cdot P_{62}(b)} \right]$$

4. והרי לנו  $\sigma$  לעימוד חלבונים

# היוריסטיקות



- פוליה '45 - "איך לפתור את זה?"
- גישה מעשית ויעילה לפתרון בעיות, ללא הבטחה למציאת הפתרון האופטימלי

- ארבעה שלבים:

1. הבנת הבעיה
2. גיבוש תכנית פעולה
3. פתרון!
4. ניתוח פתרון והכללה



האוניברסיטה  
העברית  
בירושלים  
THE HEBREW  
UNIVERSITY  
OF JERUSALEM

# אלגוריתמים בביו' חישובית

76558

שערוך פרמטרים ושיטות  
היוריסטיות

תומי קפלן  
16/1/2024

# FASTA [Pearson & Lipman, 88']

- היוריסטיקה למציאת עימוד לוקאלי בין שני רצפים.
- מבוססת על מציאת תתי-מילים [seeds] זהות.

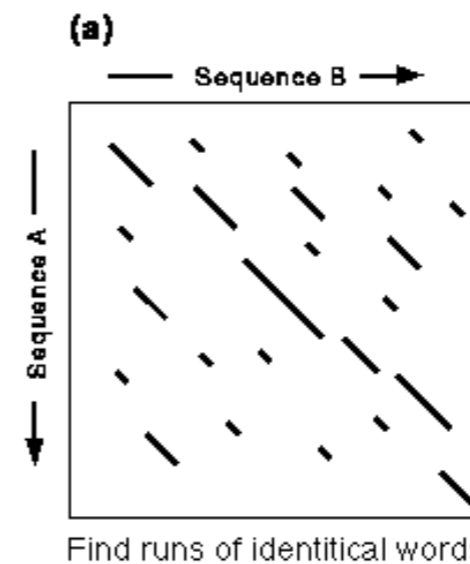
# BLAST [Altschul et al, 90']

- היוריסטיקה דומה:
- תתי-רצפים ארוכים יותר
- לא נדרשת התאמה מלאה

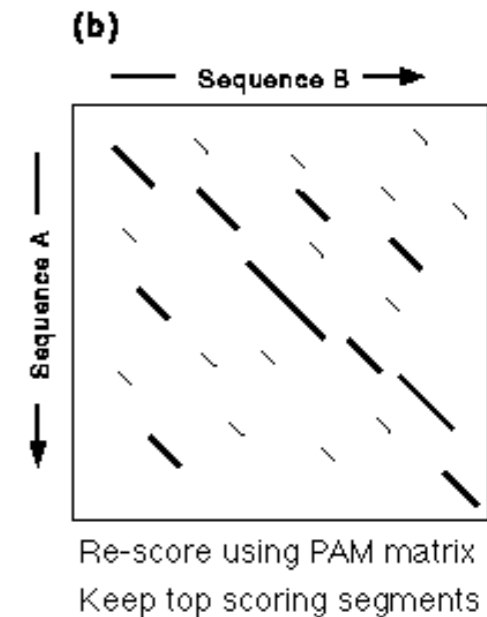
# FASTA

- Heuristic: local alignments of query vs. database
- “Good alignments have exact small matches”

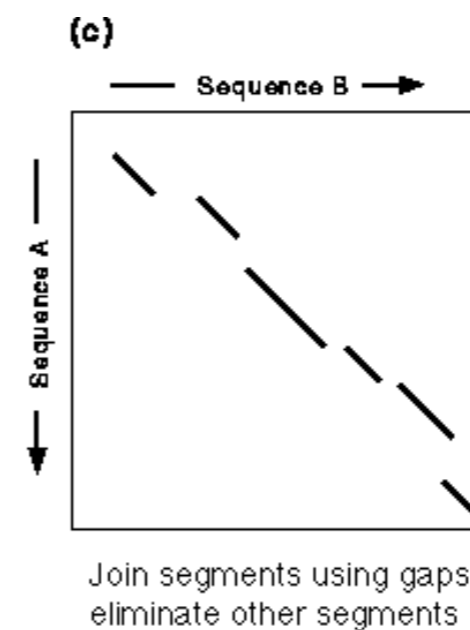
[a] Find promising matches  
[focus on top **diagonals**]



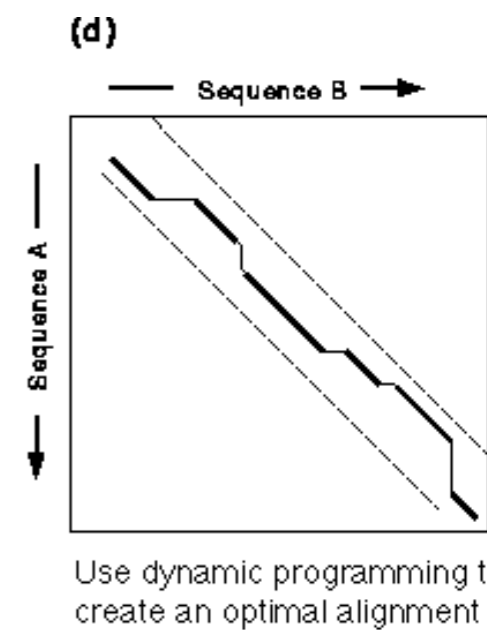
[b] Re-score [using PAM]



[c] Join gapped segments



[d] Finalize [local alignment]



# BLAST

- Heuristic: local alignments of query vs. database
  - “Good alignments have **multiple** small **hits**”
- [a] Break query into words
  - [b] Filter low-scoring words
  - [c] Store in hash table
  - [d] Run other sequences through hash table  
[allowing for mismatches]
  - [e] Find high-scoring pairs [HSPs]
  - [f] Extend