

A beautiful game

Lecture 12

(Gene prediction) ①

The best bet for simpletons (※)

The best bet for simpletons starts out with two players who select words of length l in 0-1 alphabet. Player I selects a sequence A of l alphabets from $\{0,1\}^l$, and Player II, knowing what A is, selects another sequence B of length l . The players then flip a coin to obtain heads (1) or tails (0) in turns until either A or B appears as a block of l consecutive outcomes. If A comes first, then A wins the game. (B)

Example, $l=3$,
 $A \rightarrow 010$
 $B \rightarrow 110$ | $\langle 0,0,1,1,1, \dots \rangle$

(Fact 1) The game will terminate with probability 1,

assuming the coin is a fair one or at least both sides have positive probability. (?) 两面机率不等也可以!

(Fact 2) B has higher probability to win the game. (How?)
($l \geq 3$)

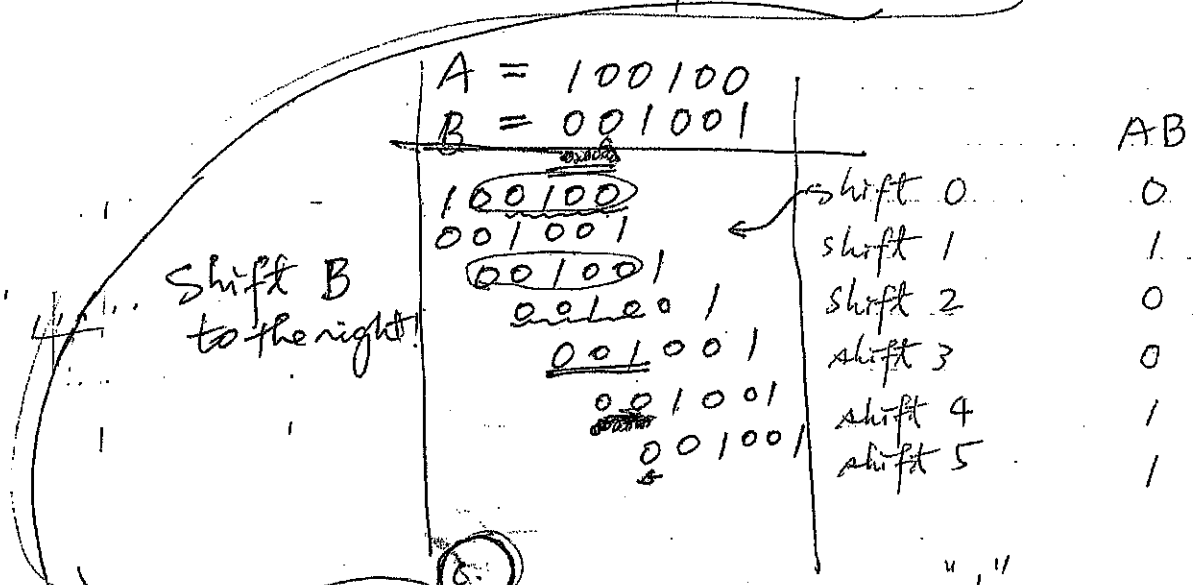
(Fact 3) The best bet for simpletons is a non-transitive game. (谁的運氣比較嗎?)

$A > B, B > C \Rightarrow A > C$.

好, 可以

(2)

Given two l -letter words A and B , the correlation of A and B , to be denoted by $AB = (c_0, c_1, \dots, c_{l-1})$, is an l -letter boolean word defined as follows:



the i -th bit of AB is defined to be "1" if the $(l-i)$ -prefix of B coincides with the $(l-i)$ -suffix of A . Otherwise, "0".

$AB = (0, 1, 0, 0, 1, 1) = (c_0, c_1, c_2, c_3, c_4, c_5)$

Definition The correlation polynomial of A and B is defined

as $K_{AB}(t) = c_0 + c_1 t + c_2 t^2 + \dots + c_{l-1} t^{l-1}$. We also denote

$K_{AB} = K_{AB}(\frac{1}{2})$. Fair coin!

$K_{AB} = (t + t^4 + t^5)(\frac{1}{2}) = \frac{1}{2} + \frac{1}{16} + \frac{1}{32} = \frac{19}{32}$

3

John Conway

The odds that B will win over A is

$$\frac{K_{AA} - K_{AB}}{K_{BB} - K_{BA}}$$

Martin Gardner, 1974 : (Remark)

I have no idea why it works. It just cranks out the answer as if by magic, like so many of Conway's other algorithms. (1993 找到最简洁的证明.)

A = 00

$K_{AB}(t) = (0, 0) \sim 0 + 0t$

B = 10

$K_{AA}(t) = (1, 1) \sim 1 + t$

如果 A 先选 00,
则 B 选 10 的胜率
最大了!

$K_{BB}(t) = (1, 0) \sim 1$

$K_{BA}(t) = (0, 1) \sim 0 + t$

A = 100
B = 010

$$\frac{K_{AA} - K_{AB}}{K_{BB} - K_{BA}} = \frac{1 + \frac{1}{2}}{1 - \frac{1}{2}} = \frac{\frac{3}{2}}{\frac{1}{2}} = \frac{3}{1}$$

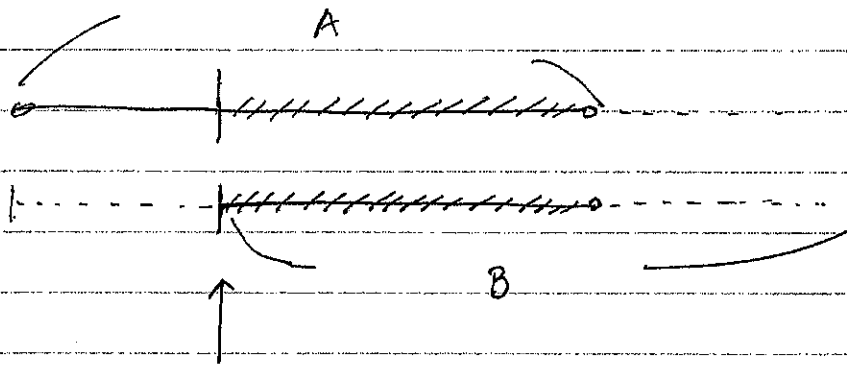
直观: 除了 00 是 A 赢之外
(10) 01, 11 都是 B 赢。

Who wins?

Proof. 1980, Li || 1981, Guibas and Odlyzko,

1993, Perzner

shortest one!



(*) 在切断之後，後半段 (A) 与 B 的前段可能是一样。
(基本上)

因此，我们探讨所谓的 Self-correlation or

auto-correlation: 自身循環移动後相似的情况。
(Optical codes)

(*) 如果 Letters 都一样，则无论如何移动都具有自身相似的情况。

(*) 所以我们对论的对象是 Letters 出现的机率给定的情况下所產生的 Words，例如 0, 1 出现的机率各 $\frac{1}{2}$ ；在不断地生成 (0, 1)-字串时，那种类型 (pattern) 的 Words 会比较容易出現，还是都一样！

Example

00, 01

010, 101

A B

) 誰的出現機率大？

前三次	A	B	
000			A 先贏
001			B 先贏
010	✓		
011			B —
100			A —
101		✓	
110			B —
111			B —

整體而言, B 贏的機率較大。

A : 010

B : 101

If Conway is correct, then what is the "策略" to win? (B over A!)

$$\Rightarrow \frac{K_{AA} - K_{AB}}{K_{BB} - K_{BA}} \geq 1$$

(*) Either $K_{AA} \gg K_{BB}$ or

$K_{AB} \ll K_{BA}$!

\Rightarrow Page 9

The Conway Equation

Finding Signals in DNA (continued)

Review of Correlation Polynomial

A', B' : l -letterwords

Correlation of A', B' : $AB' = (c_0, c_1, \dots, c_{l-1})$ where the i th bit of AB' is defined to be one "1" if the $(n-i)$ -prefix (the first $n-i$ letters) of B' coincides with the $(n-i)$ -suffix of A' , and "0" otherwise.

e.g.

$$A = \underline{100100}$$

$$A = A'$$

$$B = \underline{001001}$$

$$B = B'$$

$$AB = (0, 1, 0, 0, 1, 1)$$

$$AA = (1, 0, 0, 1, 0, 0)$$

$$A = A' = B'$$

$$BB = (1, 0, 0, 1, 0, 0)$$

$$B = B' = A'$$

以下同 A, B 代表
from

$$BA = (0, 0, 1, 0, 0, 1)$$

$$A' = B, B' = A$$

(*) Let $AB = (c_0, c_1, \dots, c_{l-1})$ and $c_{m_1}, c_{m_2}, \dots, c_{m_k}$ be the bits of AB equal to 1. ($\text{Supp}(AB) = \{m_1, m_2, \dots, m_k\}$.)

In AB , $c_1 = 1, c_4 = 1, c_5 = 1$ and thus $m_1 = 1, m_2 = 4, m_3 = 5$.

$$(\text{Supp}(AB) = \{1, 4, 5\}.)$$

(5)

(*) Denote as \mathcal{A}_{AB} the set of k prefixes of $A = a_1 a_2 \dots a_k$ of length m_1, m_2, \dots, m_k :

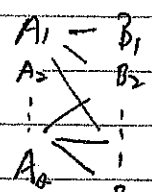
$$(a_1, \dots, a_{m_1}), (a_1, a_2, \dots, a_{m_1}, \dots, a_{m_2}), \dots, (a_1, a_2, \dots, a_{m_{k-1}}, \dots, a_{m_k})$$

e.g. If $A = \textcircled{X}YXYXY$, then

$B = \textcircled{Y}YXYXYX$
(5)

$$\mathcal{A}_{AB} = \{x, xYXYX, xYXYXY\}$$

1 4 5



Definition Let $A = a_1 a_2 \dots a_k$ and $B = b_1 b_2 \dots b_k$, then

the concatenation of A and B , $A * B = a_1 a_2 \dots a_k b_1 b_2 \dots b_k$. If

$\mathcal{A} = \{A\}$ and $\mathcal{B} = \{B\}$, then $\mathcal{A} * \mathcal{B} = \{A * B\}$ which has

possibly $|\mathcal{A}| |\mathcal{B}|$ words (perhaps with repeats).

Definition If W is an l -letter word, then let $P(W) = \frac{1}{2^l}$. For

a set $\mathcal{W} = \{W\}$, let $P(\mathcal{W}) = \sum_{W \in \mathcal{W}} P(W)$.

e.g. $P(\mathcal{A}_{AB}) = \frac{1}{2} + \frac{1}{2^4} + \frac{1}{2^5}$.

Lemma $K_{AB}(\frac{1}{2}) = P(\mathcal{A}_{AB})$.

Proof (Explain with an example)

In the above example $K_{AB}(t) = t + t^4 + t^5 = \frac{1}{2} + (\frac{1}{2})^4 + (\frac{1}{2})^5$.

(*) 在 \mathcal{A}_{AB} 中 长度分别为 1, 4, 5 的 words.

A = $\begin{matrix} XYYXY \\ XYYXY \\ XYYXY \\ XYYXY \\ XYYXY \\ XYYXY \end{matrix}$

AA
 $\begin{matrix} 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{matrix}$
 $H_{AA} = \{XYY\}$
 poly. $1+t^3$

B = $\begin{matrix} XYXYXYX \\ XYXYXYX \\ XYXYXYX \\ XYXYXYX \\ XYXYXYX \\ XYXYXYX \end{matrix}$

BB
 $\begin{matrix} 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{matrix}$
 $H_{BB} = \{YXY\}$
 $1+t^3$

$\begin{matrix} XYXYXYX \\ XYXYXYX \\ XYXYXYX \\ XYXYXYX \\ XYXYXYX \\ XYXYXYX \\ XYXYXYX \end{matrix}$

BA
 $H_{BA} = \{XY, XYXYXY\}$
 $t+t^5$

A $\begin{matrix} XYXYXYX \\ XYXYXYX \\ XYXYXYX \\ XYXYXYX \\ XYXYXYX \end{matrix}$
 B $\begin{matrix} XYXYXYX \\ XYXYXYX \\ XYXYXYX \\ XYXYXYX \end{matrix}$

$H_{AB} = \{X, XYXYXY, XYXYXY\}$
 $t+t^4+t^5$
 $\frac{K_{AB}-K_{BA}}{K_{AB}-K_{BA}} = \frac{\frac{9}{8} - \frac{19}{32}}{\frac{9}{8} - \frac{9}{32}} = \frac{17}{27}$

A word W is an A -victory if it contains A in the end and does not contain B .

A word W is an A -previctory if $W * A$ is an A -victory.

(*) $W * A$ 不一定是 A 赢 (有可能 A 没接完, B 已经赢了)。 (一般而言)

Let $S_A = \{A\text{-previctories}\}$, $S_B = \{B\text{-previctories}\}$ and

$\mathcal{J} = \{T : T \text{ is neither } A\text{-victory nor } B\text{-victory}\}$.

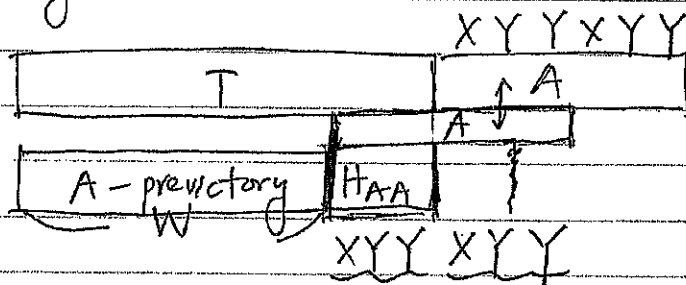
↑ the set of no-victor words.

$\forall T \in \mathcal{J}$,

Fact 1 $T * A$ corresponds to either an A -victory or a B -victory (接完的话是 A 赢)

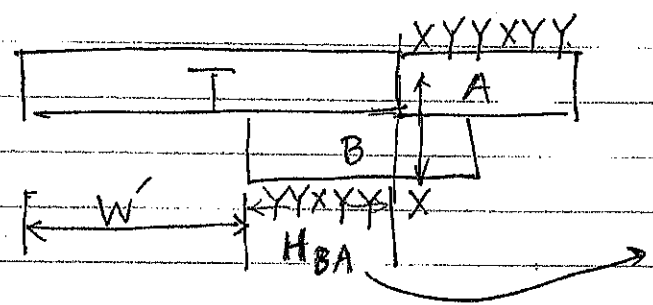
Fact 2 If $T * A$ corresponds to an A -victory, then

T can be represented as $W * H_{AA}$ where W is an A -previctory.



Fact 3 If $T * A$ corresponds to a B-victory, then

$T = W' * H_{BA}$ where W' is a B-previctory.



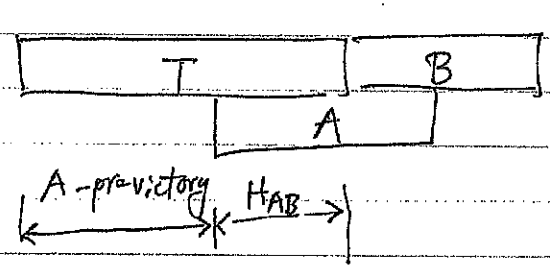
B的前段与A的后段相同

Fact 4 (Combine Fact 2 and Fact 3.)

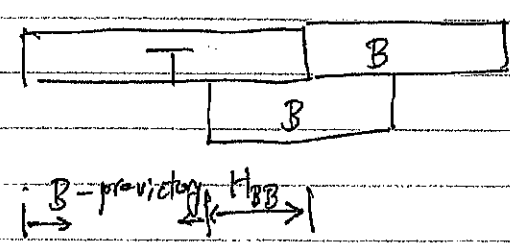
$$\mathcal{G} = \mathcal{G}_1 = (S_A * H_{AA}) \cup (S_B * H_{BA})$$

Fact 5 Similar to Fact 4,

$$\mathcal{G} = \mathcal{G}_2 = (S_A * H_{AB}) \cup (S_B * H_{BB})$$



$T * B \rightarrow A$ -victory



$T * B \rightarrow B$ -victory

Theorem The odds that B wins over A is

$$\frac{K_{AA} - K_{AB}}{K_{BB} - K_{BA}} \quad (K_{AB} = K_{AB}(\frac{1}{2}))$$

Proof

$$P(\mathcal{G}_2) = P(S_A * \mathcal{H}_{AB}) + P(S_B * \mathcal{H}_{BB})$$

$$= P(S_A)P(\mathcal{H}_{AB}) + P(S_B) \cdot P(\mathcal{H}_{BB})$$

$$= P(S_A)K_{AB} + P(S_B) \cdot K_{BB}$$

$$P(\mathcal{G}_1) = P(S_B)K_{BA} + P(S_A)K_{AA}$$

Since $P(\mathcal{G}_1) = P(\mathcal{G}_2)$,

$$P(S_A)K_{AB} + P(S_B)K_{BB} = P(S_B)K_{BA} + P(S_A)K_{AA}$$

$$\frac{P(S_B)}{P(S_A)} = \frac{K_{AA} - K_{AB}}{K_{BB} - K_{BA}}$$

(Note that the odds that B wins over A is equal to the ratio of $P(S_B)$ over $P(S_A)$, i.e., the probability of B-victories over the probability of A-victories)

Signals are everywhere!

~~(-)~~

A 0101

$$(t = \frac{1}{2})$$

B 選什麼?

Sol. Find a B such that $\frac{K_{AA} - K_{AB}}{K_{BB} - K_{BA}} \geq 1$. (" > 1 " is preferable.) K_{AA}

0101

0101

0101

0101

0101

1

0

1

0

$$\text{poly.} \approx 1 + t^2$$

$$K_{AA} = \frac{5}{4}$$

(c) Find B with smaller K_{BB}

$$\left(\begin{array}{l} \text{策略: } K_{AA} \geq K_{BB} \\ K_{AB} \leq K_{BA} \end{array} \right)$$

$$B = 1101$$

1101

1101

1101

1101

1101

1

0

0

1

$$\text{poly.} \approx 1 + t^3$$

$$K_{BB} = \frac{9}{8}$$

0101

1101

1101

1101

1101

0

0

0

1

 t^3 K_{AB}

1101

0101

0101

0101

0101

0

0

1

0

 t^2

$$K_{AB} = \frac{1}{8}$$

$$K_{BA} = \frac{1}{4}$$

$$\text{Odds} = \frac{\frac{5}{4} - \frac{1}{8}}{\frac{9}{8} - \frac{1}{4}} = \frac{\frac{9}{8}}{\frac{7}{8}} = \frac{9}{7}$$

Problem: Can we use the patterns (codewords) to determine the odds (B over A)?

What if $A = 1101$? $K_{AA} = \frac{9}{8}$

(* Smaller K_{BB} . $B = 0001$

$$K_{BB} = 1 \quad \begin{array}{cccc} 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \end{array} \quad \begin{array}{cc} & 1 \\ & 0 \\ & 0 \\ & 0 \end{array} \quad \begin{array}{cc} & 1 \\ & 1 \\ & \\ & \end{array}$$

$$K_{AB} \quad \begin{array}{cccc} 1 & 1 & 0 & 1 \\ 0 & 0 & 0 & 1 \end{array}$$

$$\text{Odds (B over A)} = \frac{\frac{9}{8}}{\frac{7}{8}} = \frac{9}{7}$$

$$K_{BA} \quad \begin{array}{cccc} 0 & 0 & 0 & 1 \\ 1 & 1 & 0 & 1 \end{array}$$

$$\text{rank}(0101) < \text{rank}(1101) < \text{rank}(0001)$$

$$K_{AA} = \frac{5}{4}, \quad K_{CC} = 1$$

$$\begin{array}{cccc} 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 1 \end{array}$$

$$K_{AC} = 0, \quad K_{CA} = \frac{1}{4}$$

$$\begin{array}{cccc} 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 1 \end{array} \quad \begin{array}{c} \\ \\ \\ \end{array}$$

$$\frac{\frac{5}{4} - 0}{1 - \frac{1}{4}} = \frac{5}{3} \quad (\text{even larger!})$$

So, can you find an example such that the transitivity of o.e.d.s does not hold?

Problems

$$x \in \mathbb{Z}_2^n$$

1. For which A and B , $\frac{K_{AA} - K_{AB}}{K_{BB} - K_{BA}} =_{\text{def}} K(A, B)$ is maximum?

2. For which $x \in \mathbb{Z}_2^n$, $K_{xx} = 1$?

Gene Prediction (Comments)

Charles Yanofsky, + Sydney Brenner et al. (1960's)

Fact 2

showed that a gene and its protein product are colinear structures with direct correlation between triplets of nucleotides in the gene and amino acids in the protein.

Later (1960's)

Overlapping genes and genes-within-genes were discovered.

(*) since, 註定要預測基因是一件極困難的工作。

尤其是在 1977, Split human genes 的發現, 終於造成 "computational gene prediction" puzzle.

1999

(Phillip Sharp and independently Richard Roberts)

Fact 1

① Most human genes are interrupted by junk DNA

and are broken into pieces call exons.