# 187 - Day 10

Wednesday, April 25, 2018　　9:34 AM

Day 10 – Breaking Rectangular Transposition

## Conditional Probability

**Definition 1.** *The **conditional probability** of an event B is the probability that this event will occur, given the knowledge that another event A has already occurred.*

$$\mathbb{P}(B|A) = \frac{\mathbb{P}(B \text{ and } A)}{\mathbb{P}(A)},$$

*assuming that* $\mathbb{P}(A) > 0.$

b/c A has occured, it will block some of the outcomes

**Example.** Two women state the following:

- A: "I have two children, the eldest is a girl."

- B: "I also have two children, and one of them is a girl."

Which of them is more likely to have two girls?

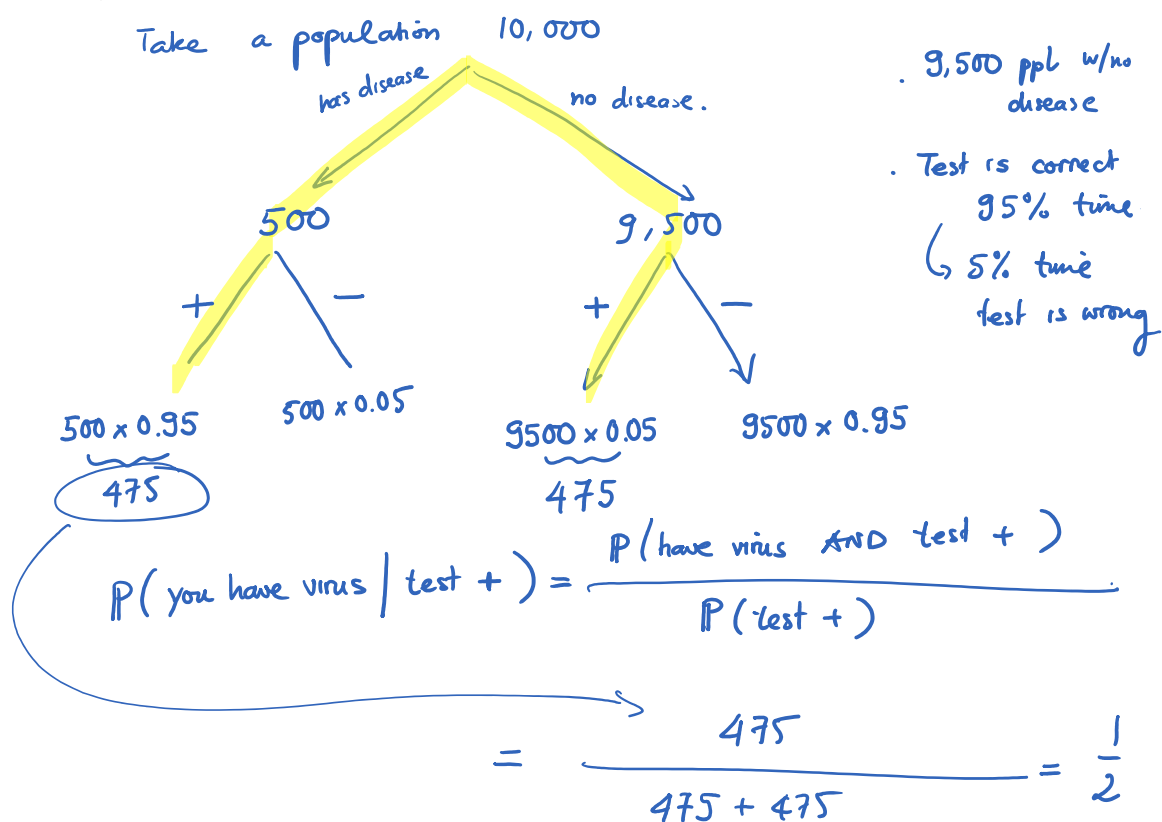2 children : Sample space = set of all outcomes
$$\{ BB, BG, GB, GG \}.$$

### A
New sample space $\{GB, \underline{GG}\}$

$$\mathbb{P}(2G's \mid 1st = G) = \frac{1}{2}.$$

$$= \frac{\mathbb{P}(2G's \text{ and } 1st = G}{\mathbb{P}(1st = G)} \quad \nearrow \mathbb{P}(GG)$$

$$= \frac{1/4}{1/2} = \frac{1}{2}$$

### B
New sample space = $\{B\widehat{G}, GB, \underline{GG}\}$

$$\mathbb{P}(2G's \mid \text{at least 1 is a } G) = \frac{1}{3}$$

**Example.**

- The Z-virus is a rare disease that hits about **5% of the population**.

- Meanwhile, Umbrella Corporation claim that they have a test that can detect Z-virus with **95% accuracy**. → 95% it will get the correct result.

Suppose that you are **tested positive** for having Z-virus (by Umbrella's test), what is the probability that you actually carry the deadly disease?

Take a population 10,000

has disease      no disease.

. 9,500 ppl w/no disease

. Test is correct 95% time
↳ 5% time test is wrong

500      9,500

$+$    $-$      $+$    $-$

$500 \times 0.05$      $9500 \times 0.95$

$500 \times 0.95$      $9500 \times 0.05$

475      475

$$P\left(\text{you have virus} \mid \text{test } +\right) = \frac{P\left(\text{have virus AND test } +\right)}{P\left(\text{test } +\right)}$$

$$= \frac{475}{475 + 475} = \frac{1}{2}$$

$\overline{\mu} \quad \overline{\lambda}$

**Example.** Given an English text. What is the probability that a randomly chosen letter $\lambda$ is $A$?

$$\mathbb{P}(\lambda = \text{"}A\text{"}) = P_A = 0.08399 \quad (\text{from table on Day 09})$$

Now suppose that we also know about the letter $\mu$ that is immediately to the left of $\lambda$. What is the probability that $\lambda = \text{"}A\text{"}$ given that

only letter after $Q$ is $U$

- $\mu = \text{"}Q\text{"}$?   $\sim Q \; \underset{?}{A} \sim$

$$\mathbb{P}(\lambda = A \mid \mu = Q) = \frac{\mathbb{P}(\lambda = A \text{ and } \mu = Q)}{\mathbb{P}(\mu = Q)} = \frac{0}{\mathbb{P}(\mu = Q)} = 0.$$

$\sim E \; \underset{?}{A} \sim$   • $\mu = \text{"}E\text{"}$?

#in cell EA

$$\mathbb{P}(\lambda = A \mid \mu = E) = \frac{\mathbb{P}(\mu\lambda = EA)}{\mathbb{P}(\mu = E)} = \frac{110/10{,}000}{1237/10{,}000} = \frac{110}{1237}$$

$\sum$ of E-row.

- $\mu = \text{"}L\text{"}$?

$$\mathbb{P}(\lambda = A \mid \mu = L) = \frac{\mathbb{P}(\mu\lambda = LA)}{\mathbb{P}(\mu = L)} = \frac{40}{391}$$

( Now suppose that $\mu$ and $\lambda$ are far apart. What is the probability that $\mu = \text{"}L\text{"}$ and $\lambda = \text{"}A\text{"}$?

$- \underset{=}{L} \sim\sim\sim \underline{A} \sim$

If the letters are far apart.
↳ they're independant.

$$\mathbb{P}(\lambda = A \text{ and } \mu = L) = \mathbb{P}(\mu = L) \cdot \mathbb{P}(\lambda = A)$$

3

$$= P_L \cdot P_A$$

from Day 09 Table.

## Breaking Rectangular Transposition

Some remarks:

1. Single letter frequencies are useless here

2. Not all pairs of adjacent English letters are equally probable

3. Table of bi-letter frequencies should reveal which pairs of letters of ciphertext were adjacent in the plaintext

**Definition 2.** *A function $y = f(x)$ is **convex** on an interval $[a, b]$ provided that $f''(x) \geq 0$ for all $a \leq x \leq b$. In particular, the first derivative $f'(x)$ is increasing on the interval $[a, b]$.*

weights

**Theorem 1.** *Let $x_1, x_2, \ldots, x_n \in [a, b]$ and let $p_1, p_2, \ldots, p_n$ be real numbers such that $p_1 + \cdots + p_n = 1$. If $f$ is **convex** on $[a, b]$ then*

$$f(p_1 x_1 + p_2 x_2 + \cdots + p_n x_n) \leq p_1 f(x_1) + p_2 f(x_2) + \cdots + p_n f(x_n).$$

*Here, equality occurs if and only if $x_1 = x_2 = \cdots = x_n$.*

average for $f(x_i) = \bar{y}$

$f(\bar{x}) \leqslant \bar{y}$

**Corollary 1.1.** *Let $f(x) = \log\left(\dfrac{1}{x}\right)$ in the above theorem to obtain*

$$\log\left(\frac{1}{p_1 x_1 + p_2 x_2 + \cdots + p_n x_n}\right) \leq p_1 \log\left(\frac{1}{x_1}\right) + \cdots + p_n \log\left(\frac{1}{x_n}\right)$$

**Theorem 2.** *Let $p_1, \ldots, p_n$ be probabilities with $p_1 + \cdots + p_n = 1$. Then for any set of probabilities $q_1, \ldots, q_n$ such that $q_1 + \cdots + q_n = 1$, we have*

$$\sum_{i=1}^{n} p_i \log(q_i) \leq \sum_{i=1}^{n} p_i \log(p_i)$$

key to break rec. trans.

To prove this, use Corollary 1.1 with $x_i = q_i / p_i$.

4

The steps for breaking rectangular transposition:

1. Guess a length for the decrypting permutation, says $k$.

2. Arrange the ciphertext into $k$ columns and let $N$ be the height (i.e. number of rows) of the resulting rectangle.

3. For each pair $1 \leq i \neq j \leq k$, extract the columns $i$ and $j$ and count the number of occurrence of the pair of letters $\alpha\beta$ and call this $n_{\alpha\beta}^{(ij)}$.

*do this for all pairs $\alpha\beta$ that appear in these columns*

4. For each pair $\alpha\beta$, let $p_{\alpha\beta}$ be the probability of the pair $\alpha\beta$ in the English language (obtain from the table of frequency for letter pairs). Compute

$$C_{ij} = \sum_{\alpha,\beta} p_{\alpha\beta} \log(n_{\alpha\beta}^{(ij)}).$$

*now repeat for all values $i, j$ in $1, 2, \ldots, k$. with $i \neq j$*

5

$k = 10,\ N = 23,\ i = 3,\ j = 7$

$k = 10$

$N = 23$

| E | C | T | I | H | N | O | H | G | I |
|---|---|---|---|---|---|---|---|---|---|
| O | K | R | O | B | C | A | O | H | F |
| E | I | N | S | G | N | N | S | A | A |
| E | T | C | N | I | I | E | C | N | H |
| O | A | S | R | E | E | H | C | T | L |
| H | S | A | A | T | E | I | B | N | E |
| S | F | N | E | U | C | N | O | E | R |
| R | E | T | I | U | S | S | S | A | A |
| R | E | O | C | U | W | S | O | I | F |
| M | N | D | A | O | D | I | D | V | A |
| T | E | C | H | E | X | O | T | T | E |
| H | O | F | E | T | C | E | R | L | A |
| I | I | A | T | S | O | E | S | M | S |
| M | S | T | E | I | O | N | K | W | N |
| N | I | C | S | O | S | F | S | O | T |
| X | Y | S | T | I | U | H | F | R | O |
| A | R | E | G | X | S | A | A | E | M |
| S | M | C | Y | H | L | Z | B | I | O |
| B | A | E | Y | D | R | I | P | T | A |
| L | R | C | A | U | R | N | A | A | R |
| M | N | G | E | E | F | I | T | S | O |
| T | A | X | R | S | H | A | I | T | G |
| B | O | N | R | D | N | I | K | L | E |

$i = 3 \qquad j = 7$

$\Longrightarrow$

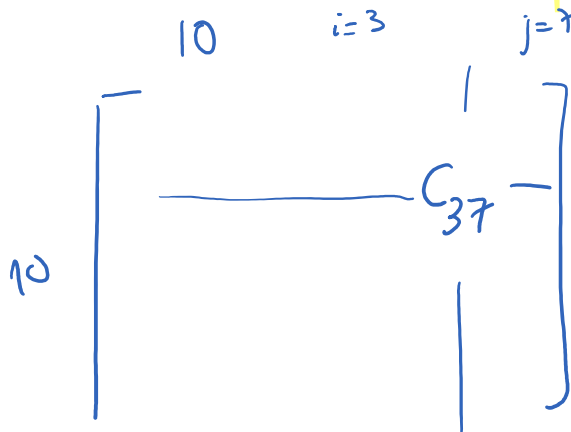| T | O |
|---|---|
| R | A |
| N | N |
| C | E |
| S | H |
| A | I |
| N | N |
| T | S |
| O | S |
| D | I |
| C | O |
| F | E |
| A | E |
| T | N |
| C | F |
| S | H |
| E | A |
| C | Z |
| E | I |
| C | N |
| G | I |
| X | A |
| N | I |

count the pairs

1 TO, 1 RA

2 NN', ...

$n_{TO}^{(3,7)} = 1$

$n_{RA}^{(3,7)} = 1$

$n_{NN}^{(3,7)} = 2$

$\vdots$

$$C_{37} = P_{TO} \cdot \log\left(n_{TO}^{(3,7)}\right)^{\frac{1}{}} {}^{0} + P_{RA} \cdot \log\left(n_{RA}^{(3,7)}\right)^{1} + P_{NN} \cdot \log\left(n_{NN}^{(3,7)}\right) + \dots$$

10

i = 3 \qquad j = 7

10 $\left[\ \underline{\qquad\qquad} C_{37} \underline{\quad}\ \right]$ 6

Compute all $C_{i;j}$

*If col j not follow col i then $C_{ij}$ small.*

Define $f_{\alpha\beta}^{(ij)} = \dfrac{n_{\alpha\beta}^{(ij)}}{N}$. When two columns $i$ and $j$ **were not consecutive in the plaintext**, then

$$C_{ij} = \sum_{\alpha,\beta} p_{\alpha\beta} \log(N \cdot f_{\alpha\beta}^{(ij)})$$

$$= \log(N) + \sum_{\alpha,\beta} p_{\alpha\beta} \log(f_{\alpha\beta}^{(ij)})$$

$$\leq \sum_{\alpha,\beta} p_{\alpha\beta} \log(p_{\alpha\beta})$$

$$\begin{bmatrix} s & s & 0 & & s & s & B & s \end{bmatrix}$$

so $C_{ij}$ is much smaller comparing to when two columns $i$ and $j$ were consecutive in the plaintext.

So **if we guessed the correct period** then the matrix $[C_{ij}]_{1 \leq i \neq j \leq k}$ will have a substantially bigger number in **each row, except one.**

- If $C_{ij}$ is the substantially big number on row $i$ then $j$ follows $i$ in the decryption permutation.

- If row $k$ is the only row with no substantially big entry, then $k$ is the first entry in the decryption permutation.

$$\begin{bmatrix} 0 & s & s & B & s \\ s & 0 & B & s & s \\ s & s & 0 & s & s \\ s & s & s & 0 & B \\ s & B & s & s & 0 \end{bmatrix}$$

*B's cannot be on the same row/col.*

7