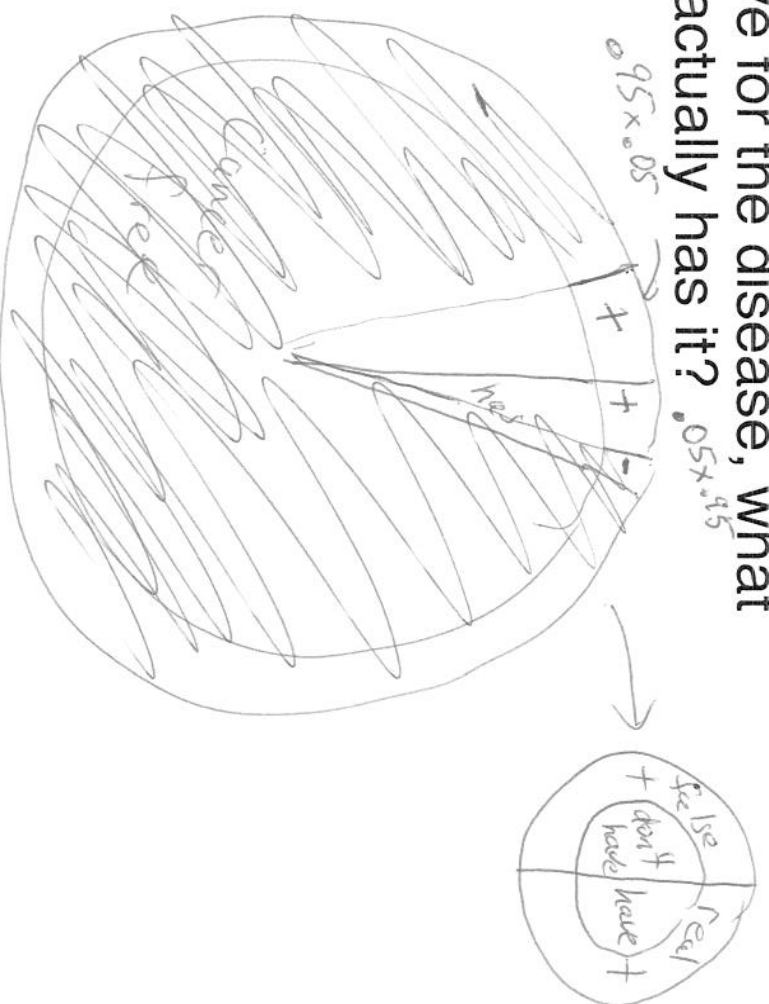Say that there is a disease which only about 5% of the population has and that there is a test for this disease which is roughly 95% accurate (that is someone who has the disease will test positive 95% of the time and negative 5% of the time, while someone who does not have the disease will test negative 95% of the time and positive 5% of the time).

Given that a patient tests positive for the disease, what is the probability that he or she actually has it?

Is the answer?

6  A) 95%
3  B) 90%
0  C) 75%
7  D) 50%
2  E) Don't know / care

$.95 \times .05$

$.05 \times .95$

has

Cancer free

false
+ don't
have

real
have +

real +

don't
have

# A simple test for monoalphabetic substitution

In English or monoalphabetic encrypted text we observe:

|  | English: | MISSISSIPPI |
|--|----------|-------------|
| Monoalphabetic: | | RDFFDFFDOOD |
| Vigenere : | | PQJLLAJBSXZ |

$\uparrow N-1$ pairs $\longleftarrow N = 11$ letters in text

While in polyalphabetic cyphertext we should observe:

$$P_{AA} + P_{BB} + P_{CC} + \cdots + P_{ZZ} \approx .027$$

$$P(\alpha\alpha \text{ occurrs in random cyphertext}) = \frac{1}{26} \approx .038$$

We should note (of course) that this only works for reasonably large amounts of text.

Look at cyphertext # double letters next = M
to each other
if there are N letters in text

$$\frac{M}{N-1} = .027 \text{ or } .038$$

If the cyphertext was obtained from a polyalphabetic cipher then the index of coincidence can also be used to estimate the period of the cipher.

Let $p$ be the period of the cyphertext and place the letters of the cyphertext into groups of $p$ so that the letters in the $i^{th}$ position of the groups are all encrypted with the same key.

- Let $M_\alpha^{(i)}$ equal the number of occurrences of the letter $\alpha$ that appears in the $i^{th}$ positions in the groups.
- If there are $M$ groups of $p$, then $\sum_{\alpha=A}^{Z} M_\alpha^{(i)} = M = \#$ of letters in the $i^{th}$ column
- We also have $N = Mp$
- Also we can estimate that $M_\alpha^{(i)} \approx M p_{\sigma(\alpha)}$ (again for some permutation for the alphabet $\sigma$)

$$
i = \overbrace{1 \ \ 2 \ \ 3 \cdots p}^{p} \ \ \leftarrow
$$

$$
M \left\{ \begin{array}{l} T\,H\,I\,S\,I \\ S\,M\,Y\,P\,L \\ A\,I\,N\,T\,E \\ X\,T\,O\,M\,O \\ A\,R\,E\,\xleftarrow{} \end{array} \right.
$$

Now, we calculate that

pairs of equal letters in same column

pairs of equal letters in different columns

$$
2D_c = \sum_{i=1}^{p} \sum_{\alpha=A}^{Z} M_\alpha^{(i)}(M_\alpha^{(i)} - 1) + 2\sum_{i=1}^{p} \sum_{j=i+1}^{p} \sum_{\alpha=A}^{Z} M_\alpha^{(i)} M_\alpha^{(j)}
$$
$$
\approx M^2 p(.065) - pM + M^2(.038)p(p-1)
$$
$$
= \frac{N^2}{p}(.027) - N + N^2(.038)
$$

The index of coincidence is defined as

$$I_c = \frac{\text{number of pairs of equal letters in ciphertext}}{\text{the total number of pairs of letters}}$$

That is if we set

- $N_\alpha =$ the number of occurrences of the letter $\alpha$ in the cyphertext

- 

$$D_c = \sum_{\alpha=A}^{Z} \binom{N_\alpha}{2}$$

ABCDE

$N = \#$ of letters

$\binom{N}{2} = \frac{N(N-1)}{2}$

$= \#$ of pairs
of letters

$D_c$ represents the number of pairs of equal letters in the cyphertext.

- then $I_c = \frac{D_c}{\binom{N}{2}}$

- where $N =$ the number of letters in the cyphertext

The index of coincidence is invariant under monoalphabetic cyphers and we estimate under this condition that $N_\alpha = N * p_{\sigma(\alpha)}$ for some permutation of the alphabet $\sigma$ and so

$$I_c = \frac{\sum_{\alpha=A}^{Z}(N_\alpha^2 - N_\alpha)/2}{N(N-1)/2}$$

$$\approx \frac{N^2(\sum_{\alpha=A}^{Z} p_\alpha^2) - N}{N(N-1)}$$

$$= \frac{N(.065) - 1}{N - 1}$$

$$\approx .065$$

$P_{\sigma(\alpha)} = $ prob that
$\alpha$ appears in
cyphertext

Index of
coincidence
of English
OR monoalphabetic
substitution.

Note that because $I_c = \frac{D_c}{\binom{N}{2}} = \frac{D_c}{N(N-1)/2}$, we have that

$$2D_c = N(N-1)I_c.$$

And we just derived that

$$2D_c \approx \frac{N^2}{p}(.027) - N + N^2(.038)$$

Therefore,

$$N(N-1)I_c \approx \frac{N^2}{p}(.027) - N + N^2(.038)$$

$$(N-1)I_c \approx \frac{N}{p}(.027) - 1 + N(.038)$$

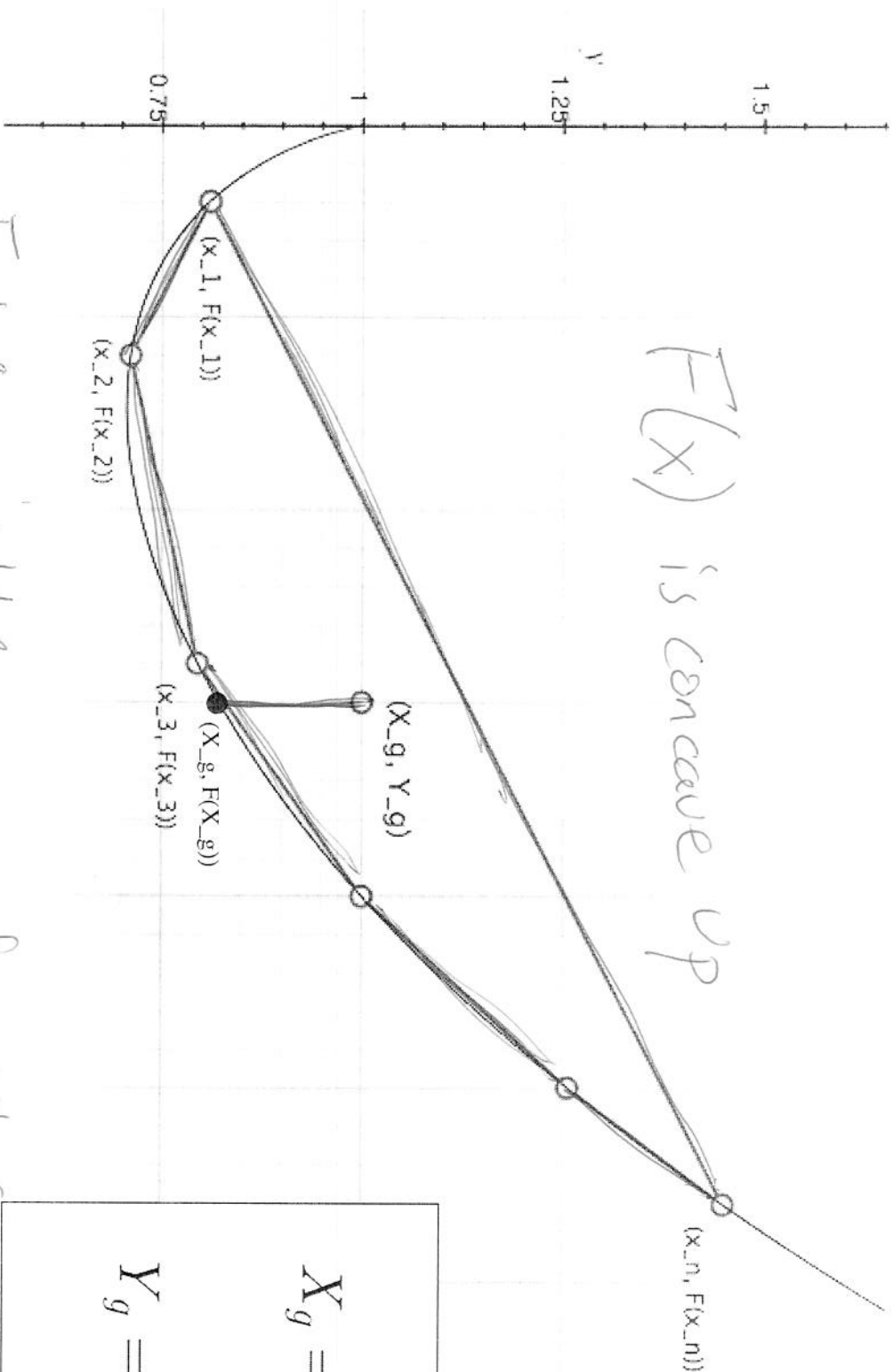$$(N-1)I_c + 1 \approx \frac{N}{p}(.027) + N(.038)$$

$$(N-1)I_c + 1 - N(.038) \approx \frac{N}{p}(.027)$$

$$p((N-1)I_c + 1 - N(.038)) \approx N(.027)$$

$$p \approx \frac{N(.027)}{(N-1)I_c + 1 - N(.038)}$$

formula for $p$ in terms of $N$ and the index of coincidence.

$F(x)$ is concave up



Graph point labels: (x_1, F(x_1)), (x_2, F(x_2)), (x_3, F(x_3)), (X_g, F(X_g)), (X_g, Y_g), (x_n, F(x_n))

Axis labels: y, 0.75, 1, 1.25, 1.5

$$X_g = \sum_{i=1}^{n} m_i x_i$$

$$Y_g = \sum_{i=1}^{n} m_i F(x_i)$$

$$\sum_{i=1}^{n} m_i = 1$$

$$m_i \geq 0$$

Fact 1: weighted average of points $(x_i, F(x_i))$ lies inside of convex polygon.

Fact 2: the point which is directly below the weighted average is outside the polygon (OR on edge)

$$F\left(\sum_{i=1}^{n} m_i x_i\right) = F(X_g) \leq Y_g = \sum_{i=1}^{n} m_i F(x_i) = \sum_{i=1}^{n} m_i y_i$$

$m_i$ = a set of probabilities $q_i$

where $\displaystyle\sum_{i=1}^{n} q_i = 1$

take another set of probabilities

$p_1, p_2, \cdots, p_n$

Set $\quad X_i = p_i / q_i \qquad m_i = q_i$

$$F(x) = x \log_b(x)$$

$$F'(x) = \log_b(x) + x \left(\frac{\ln(x)}{\ln(b)}\right)'$$

$$= \log_b(x) + x \cdot \frac{1}{x \ln(b)}$$

$$F''(x) = \frac{1}{x \ln(b)} \quad \text{concave up!}$$

$$F\left(\sum_{i=1}^{n} q_i \left(\frac{p_i}{q_i}\right)\right) \leq \sum_{i=1}^{n} q_i F\left(p_i / q_i\right)$$

$$\underset{0}{\shortparallel}$$

$$\shortparallel$$

$$= \sum_{i=1}^{n} q_i \left(\frac{p_i}{q_i}\right) \log\left(p_i / q_i\right)$$

$$= \sum_{i=1}^{n} p_i \left(\log p_i - \log q_i\right)$$

# Conclusion

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