
THE APPLICATION OF INFORMATION THEORY TO VOWEL-RECOGNITION EXPERIMENTS

J. G. BLOM*

This paper deals with the application of Information theory to the transmission of natural and artificial vowels. I want to start with a brief explanation of the main concepts of information theory for those not acquainted with them.

Information theory describes the phenomena of transmission as perceived by an outside observer who has full knowledge of both sides of the transmission channel.

The symbols to be coded by the transmitting part of the channel—in our case in sounds—will be referred to as input, the symbols decoded at the receiving end as output.

In this case we are only interested in the most simple situation in which the auto-correlation of the string of input symbols is zero, which means that the input symbols are in a random order. Incidentally this does not exclude the possibility that the decoding process is affected by the actual succession of two or more sounds. The number of different symbols will be finite.

Due to imperfections or instability of transmitter and receiver, distortion or interference, the string of output symbols will not be an exact replica of the string of input symbols. We speak therefore of a transmission channel with noise.

Let the number of different symbols be n .

The symbols can then be referred to as $S_1, S_2 \dots S_n$.

The performance of the channel can be depicted by a table of confusion probabilities (see fig. 1).

$$\sum_j p_{j0} = \sum_k p_{0k} = \sum_j p_{jk} = 1$$

For a noise-free channel

$$\begin{array}{ll} p_{j0} = p_{jk} = p_{0k} & \text{for } j = k \\ p_{jk} = 0 & \text{for } j \neq k \end{array}$$

For a channel with no correlation between input and output (that means no transmission at all, the receiver is only guessing)

$$\begin{array}{ll} p_{jk} = p_{j0} \cdot p_{0k} & \text{for } j = 1, \dots, n \\ & k = 1, \dots, n \end{array}$$

* University of Amsterdam, Institute of Phonetic Sciences.

Input		Output			
		S_1	S_2	S_k	S_n
S_1	p_{10}	p_{11}	p_{12}	p_{1k}	p_{1n}
S_2	p_{20}	p_{21}	p_{22}	p_{2k}	p_{2n}
S_j	p_{j0}	p_{j1}	p_{j2}	p_{jk}	p_{jn}
S_n	p_{n0}	p_{n1}	p_{n2}	p_{nk}	p_{nn}
Total	1	p_{01}	p_{02}	p_{0k}	p_{0n}

Fig. 1. Confusion Probability Matrix.

p_{j0} = probability that S_j is the input symbol,
 p_{0k} = probability that S_k is the output symbol,
 p_{jk} = probability of the combination of S_j as input symbol and S_k as output symbol.

In a forced-choice situation

A real channel will be somewhere between these extremes. Now we have to deal with different amounts of information. The information of the input H_x , that of the output H_y and that of the combination of input and output H_{xy} .

The unit of information is called a Bit.

One bit is the amount of information contained in the answer to a question to which there are two mutually exclusive answers with equal probability of occurrence. Take for example the information contained in the position of a coin. So the amount of information in Bits is the minimal number of questions of the type just mentioned necessary to obtain full knowledge. The amounts of information can be easily calculated using the formulae of fig. 2.

$$H_x = \sum_j -p_{j0} \log p_{j0}$$

$$H_y = \sum_k -p_{0k} \log p_{0k}$$

$$H_{xy} = \sum_{j,k} -p_{jk} \log p_{jk}$$

Fig. 2.

When information is transmitted by the channel we have the following inequality

$$H_x + H_y > H_{xy}$$

This means that given the output and our knowledge about the confusion matrix, we can make a good guess at the input.

The relations between H_x , H_y , and H_{xy} can be shown in simple Wenn-diagrams. (See fig. 3.)

$$T_{xy} = H_x + H_y - H_{xy}$$

The cross-section between H_x and H_y is called the transmission T_{xy} .
The physical meaning of the transmission is that part of the information of the input which we know when the output is known, in other words, the transmission is the information transmitted by the channel.

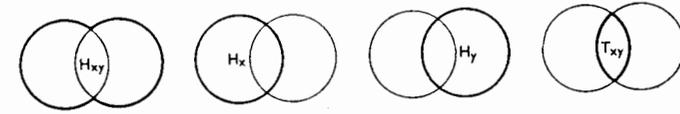


Fig. 3.

In order to calculate the transmission we have to make use of the confusion frequency matrix resulting from an experiment (fig. 4).

Input	Total	Output			
		S_1	S_2	S_k	S_n
S_1	m_{10}	m_{11}	m_{12}	m_{1k}	m_{1n}
S_2	m_{20}	m_{21}	m_{22}	m_{2k}	m_{2n}
S_j	m_{j0}	m_{j1}	m_{j2}	m_{jk}	m_{jn}
S_n	m_{n0}	m_{n1}	m_{n2}	m_{nk}	m_{nn}
Total	M	m_{01}	m_{02}	m_{0k}	m_{0n}

Fig. 4. Confusion Frequency Matrix.

In this table the m 's represent observed frequencies, the subscripts have the same meaning as in the probability matrix.

Taking the quotients m/M as best estimates for p 's we can calculate the transmission.

The necessary calculations can easily be programmed for evaluation by an electronic computer.

All our calculations were carried out with the IBM 1130 system of the Institute of Phonetic Sciences of the University of Amsterdam.

To get some insight into the process of vowel perception we applied information theory to some data published in the literature.

We started with the well-known experiment by Peterson and Barney on formant measurements on vowels of different speakers (JASA 1952) (fig. 5).

Suppose we have a vowel-recognition system that relates the sounds within a specific contour to one and only one vowel-class.

We determined the confusion frequency matrix for such a system shown in fig. 6 by a simple counting procedure, any sound falling in the cross-section of two areas being scored as 0.5 for each area. All frequencies are multiplied by 10 to avoid fractions.

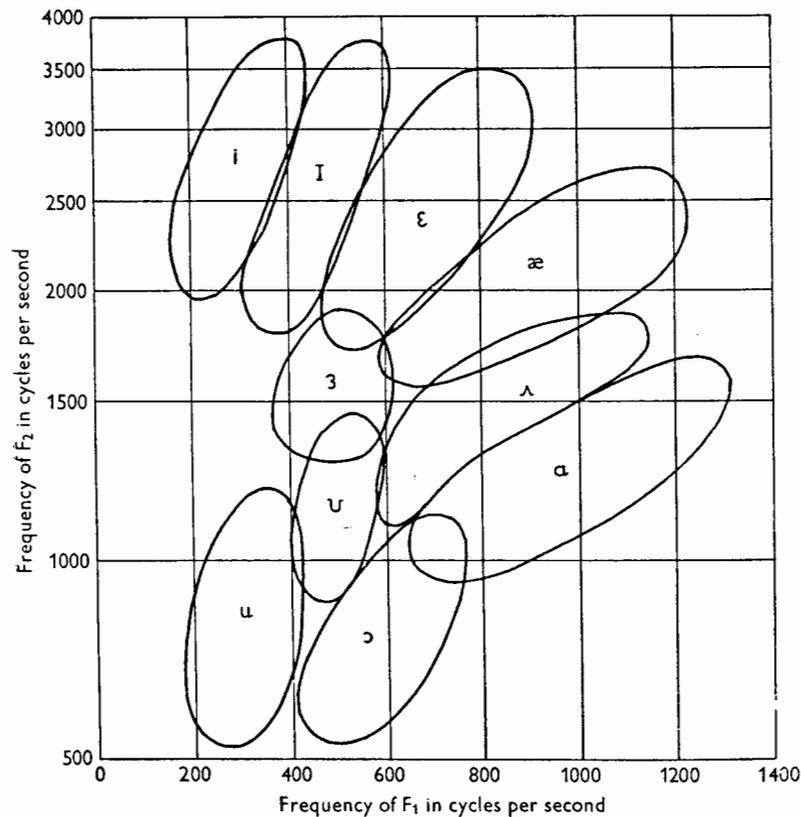


Fig. 5.

As we see, the information of the input is 3.32 Bits, the transmission 2.19 Bits. The same sounds were presented to a group of listeners. Peterson and Barney published the confusion matrix which is shown here as fig. 7.

When we apply our formulae to their matrix we find a transmission of 2.98 Bits.

It appears that human listeners do better than our hypothetical vowel recognition system. Our conclusion must be that man uses factors additional to the first two formants. These factors might be fundamental frequency, duration, the connection with surrounding consonants, and knowledge of the particular vowel system of an individual speaker. Although the speechsounds of different speakers were randomized, some knowledge of the position of the vowel system in the two-formant plane was available, due to the high correlation between fundamental freq. and the formant frequencies (Mol 1964).¹

As no confusion occurs when we listen to the sounds of a familiar voice we can list and add up our data as follows.

¹) Proceedings of the 5th Intern. Congress of Phonetic Sciences.

Peterson & Barney Formant Measurements											
Input	Total	Out-put 1	2	3	4	5	6	7	8	9	10
1	740	680	60	0	0	0	0	0	0	0	0
2	750	30	650	55	0	0	0	0	0	0	15
3	755	1	96	571	11	1	1	1	1	1	71
4	740	1	1	121	566	26	1	1	1	1	21
5	730	0	0	0	30	610	45	35	0	5	5
6	760	3	3	3	3	53	603	83	3	3	3
7	740	2	2	2	2	22	92	572	12	32	2
8	740	2	2	2	2	12	2	2	567	92	57
9	750	1	1	1	1	41	1	11	76	491	126
10	775	2	22	97	52	2	2	2	7	127	462
Total	7480	722	837	852	667	767	747	707	667	752	762

$H(X) = 3.32$ $H(Y) = 3.31$ $H(XY) = 4.45$ $T(XY) = 2.18$

Fig. 6.

Information of input		3.32 Bits
Contributed by formant positions alone	maximal	2.19 Bits
Contributed by other factors than system	at least	0.79 Bits
	Sum of these factors	2.98 Bits
Contributed by specific knowledge of an individual speakers vowel system		0.34 Bits
		3.32 Bits

The next data to be examined are published by Cohen, Slis & 't Hart (Phonetica 1967) in an article entitled "On Tolerance and Intolerance in vowel perception". They utterly failed to grasp the meaning of a paper by Blom & Uys, entitled "Some Notes on the Existence of a 'Universal Concept' of Vowels" (Phonetica 1966), but they presented a highly interesting confusion matrix for a system of 12 synthetic vowels. They used 12 fixed two-formant positions and introduced duration as an extra parameter. The spacing of the vowels in the F_1 , F_2 plane is somewhat exaggerated. The matrix is shown in fig. 8.

The information of the input is 3.62 Bits in formant positions and 1.55 Bits in duration which is redundant.

// XEQ
Peterson & Barney Listening Experiment

Input	Total	Out-put	1	2	3	4	5	6	7	8	9	10
1	10 280	10 267	4	6	0	0	3	0	0	0	0	0
2	10 279	5	9 549	694	2	1	1	0	0	0	0	26
3	10 277	0	257	9 014	949	1	3	0	0	0	2	51
4	10 278	0	1	300	9 919	2	2	0	0	0	15	39
5	10 273	0	1	0	19	8 936	1 013	69	0	228	7	
6	10 279	0	0	1	2	590	9 534	71	5	62	14	
7	10 279	0	0	1	1	16	51	9 924	96	171	19	
8	10 279	0	0	1	0	2	0	78	10 196	0	2	
9	10 277	0	1	1	8	540	127	103	0	9 476	21	
10	10 279	0	0	23	6	2	3	0	0	2	10 243	
Total	102 780	10 273	9 813	10 041	10 906	10 090	10 737	10 245	10 297	9 956	10 422	

$H(X) = 3.32$ $H(Y) = 3.32$ $H(Y) = 3.66$ $T(XY) = 2.98$

Fig. 7.

The transmission is 2.93 Bits. As the experimental conditions are comparable with the situation in which a person is listening to the sounds of one individual speaker, part of the information is lost. (Of course, some of the factors operating in experiments where monosyllabic words are used are absent in experiments with isolated sounds).

This low transmission is in agreement with our findings. It seems that a transmission channel operates less stably with artificial vowel-like sounds than with natural vowels.

From the results of the scaling experiment described by my colleague Meinsma an estimate can be made as to the confusion occurring between different areas of the perceptive vowel-triangle. We estimate the following data:

$$H_x = 3.6 \text{ Bits}$$

$$H_y = 3.6 \text{ Bits}$$

$$H_{xy} \approx 5.3 \text{ Bits}$$

$$T_{xy} \approx 1.9 \text{ Bits}$$

This means that the duration factor introduced by Cohen and collaborators must have contributed about 1 Bit of the 1.15 Bits of partly redundant transmitted information.

The present study is part of a larger programme which aims at the generation of vowel systems of optimal efficiency for the production of artificial speech.

Cohen, Slis.'t Hart

Input	Total	Out-put	1	2	3	4	5	6	7	8	9	10	11	12
1	1 670	1 628	42	0	0	0	0	0	0	0	0	0	0	0
2	1 669	82	1 584	0	0	1	0	0	0	0	2	0	0	0
3	1 664	0	25	1 570	0	0	12	0	0	0	57	0	0	0
4	1 670	1	0	47	1 546	68	1	4	1	0	0	2	0	0
5	1 669	4	7	7	3	1 628	0	3	10	6	0	0	1	0
6	1 669	0	0	5	0	1	1 475	128	6	53	0	0	1	0
7	1 670	0	0	0	0	0	126	1 536	2	4	0	2	0	0
8	1 670	0	0	0	0	1	63	10	1 592	0	2	0	2	0
9	1 665	0	30	318	0	0	2	3	0	1 309	1	2	0	0
10	1 663	0	1	221	13	12	19	4	0	28	1 358	7	0	0
11	1 670	0	0	34	274	3	1	5	0	45	322	981	5	0
12	1 670	0	0	0	0	269	1	0	3	0	109	4	1 284	0
Total	20 019	1 715	1 689	2 202	1 836	1 983	1 700	1 693	1 614	1 504	1 792	998	1 293	0

$H(X) = 3.58$ $H(Y) = 3.56$ $H(XY) = 4.21$ $T(XY) = 2.93$

Cohen, Slis.'t Hart Vowels in one Durationclass Added

Input	Total	Out-put	1	2	3	4	5	6	7	8	9	10	11	12
1	6 680	1 629	1 820	1 545	985	42	71	3	49	5	81	128	322	
2	8 343	86	3	16	6	1 621	1 899	1 605	1 317	1 287	325	66	112	
3	4 996	0	13	132	7	26	13	6	138	1	1 796	1 506	1 358	
Total	20 019	1 715	1 836	1 693	998	1 698	1 983	1 614	1 504	1 293	2 202	1 700	1 792	

$H(X) = 1.55$ $H(Y) = 3.56$ $H(XY) = 3.97$ $T(XY) = 1.13$

Fig. 8.

DISCUSSION

Newel:

The experiments I have performed on human perception of vowel sounds with and without prior knowledge of the speaker would indicate that the remark "as no confusion occurs when we listen to the sounds of a familiar voice" is incorrect with reference to data of the Peterson and Barney type. This will surely invalidate those conclusions made on the basis of this premise.