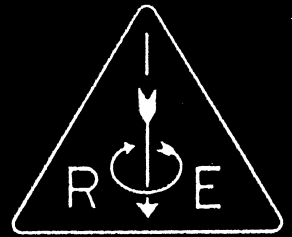


IRE Transactions



on INFORMATION THEORY

Vol. IT-2, No. 3

September, 1956

1956 SYMPOSIUM ON INFORMATION THEORY

held at

**Massachusetts Institute of Technology
Cambridge, Massachusetts**

September 10-12, 1956

**PUBLISHED BY THE
Professional Group on Information Theory**

ON an APPLICATION of SEMI GROUPS METHODS
TO SOME PROBLEMS in CODING

By M.P. Schützenberger
(C.N.R.S. Paris)

0. Introduction.

The current paper deals with a chapter in what could be called communication theory in extensive form : it starts with extremely restricted structures and it stops where begins the canonical problem of optimisation. It even ends sooner for no full use of the definitions is made and the main ergodic theorem is stated without proof.

Actually the nature itself of the question under study has commanded these restrictions together with the architecture of the paper : we give a abstract model of some sort of language and we try to show how semi group concepts apply fruitfully to it with the hope that some of them may be at least of stimulating interest to specialists working on natural languages.

As frequent in the field of cybernetic, the mathematics involved even if quite simple are far away from classical analysis and, indeed, many of the necessary tools had to be sharpened especially for the purpose.

Thus the paper is twofold : in a first part the model and its main properties are discussed at a concrete level on the simplest cases : the coding and decoding with length bounded codes. In a second part a selection of theorems are proved whenever the necessary semi group theoretic preliminaries are not exacting. The link along this tail of appendices is the theory developped verbally in the first part. Finally a special chapter provides a bridge toward probabilistic applications.

It is proper at this place to acknowledge the contributions of three authors who influenced deeply the building of the theory :

Sardinas and Patterson⁽¹⁾ who discussed first on a logical basis the general coding process.

B. Mandelbrot⁽²⁾ who recognised and studied extensively the role of "word units" in communication theory and related the problem to Feller's recurrent events.

P. Dubreuil⁽³⁾ and his school whose pionnering work on discrete semi groups has provided many basic concepts and arguments as it will be seen below.

Part I

1. Preliminary definition of a discrete semi group language :

We shall be concerned with the two basic sets of communication theory :

The set of all messages which may possibly be sent.

The set of all signals available for transmission along the line.

The main feature of the theory is the postulational requirement that the signals as well as the messages pertain both to some common class of structures so that coding and decoding not only be inverse operations but far more generally, be special instances of a quite broad new process, that of translation.

This identity of structure itself between two sets is a result from the basic restriction that they develop homogeneously in time - or more accurately that both admit a common partial order and composition operation.

That such requirements are rather stringent is clearly seen by the exemple of photography (two exposures give rarely a result which is, in any sense, equivalent to a third one) or even by harmonic modulation where Fourier transform exchanges so well time and frequency that finite signals cannot be fully adequate.

On the other hand, languages either spoken, written or gesticulated are somewhat akin with our consideration, and we shall use the name of "discrete semi group languages" (d.s.g.l.) for naming the elemental concepts of our study.

The definitions below are quite general and as said before, no full use of them will be made here - very little gain in simplicity would be achieved by using more restrictive ones.

DEFINITIONS :

I. A discrete semi group language will be a set Λ of object called "messages" satisfying the following conditions :

I.1. If λ_i and λ_j pertain to Λ so does their "product" $\lambda_k = \lambda_i \lambda_j$ made up of " λ_i " followed by " λ_j " (λ_i will be said a left divisor and λ_j a right divisor of λ_k).

I.2. If λ_i, λ_j and λ_k pertain to Λ and if $\lambda_i \lambda_j = \lambda_i \lambda_k$ and $\lambda_i \lambda_m = \lambda_j \lambda_k$ then $\lambda_i \lambda_m$ is identical with $\lambda_j \lambda_k$.

I.3. The "vacuous message" ρ pertains to Λ and satisfies $\rho \lambda_i = \lambda_i \rho = \lambda_i$ for all $\lambda_i \in \Lambda$.

I.4. There is a sub set Λ_c from Λ called "dictionary" or "basis" whose elements are called "words". Λ_c is such as :

I.4.1. ρ does not pertain to Λ_c .

I.4.2 for all $\lambda_i \in \Lambda - \Lambda_c$

either $\lambda_i \in \Lambda_c$

either these exist a unique finite set of words $\lambda_1, \lambda_2, \dots, \lambda_{im} \in \Lambda_0$

with
 $\lambda_i = \lambda_1 \lambda_2 \dots \lambda_{im}$

II. Given two d.s.g.l. Λ and M a correspondence θ between the elements of two subsets $\Lambda' \subset \Lambda$ and $M' \subset M$ will be said a translation if it satisfies :

II.1. The correspondence is one to one where ever it is defined.

II.2. If $\lambda_i, \lambda_j \in \Lambda', \theta \lambda_i = \mu_i, \theta \lambda_j = \mu_j$, then $\lambda_i \lambda_j \in \Lambda'$ and $\theta \lambda_i \lambda_j = \mu_i \mu_j$

II.3. The translation will be said :

Total from Λ to M , if $\Lambda' = \Lambda$.
Subtotal from Λ to M , if for all $\lambda_i \in \Lambda$ there is at least a $\lambda_j \in \Lambda'$ such as $\lambda_i \lambda_j \in \Lambda'$.

III. A neat coding of Λ into M will be a translation total from Λ to M and subtotal from M to Λ .

In algebraic form we could reduce our axiomatic to :

- I' : Λ is the free discrete semi group generated by Λ_0
- II' : A translation is an isomorphism between the sub semi groups $\Lambda' \subset \Lambda$ and $M' \subset M$
- III' : A translation is a neat coding if $\Lambda' = \Lambda$ and M' is a subsemigroup of M neat on the right. (Note that "subsemigroup" entails I.1, I.2 and I.3 ; "free" corresponds to unique in I.4.2 , "discrete" to finite at the same place).

2. Practical significance of the axiomatic :

Let us take a simple example in coding :

$\Lambda_0 = \{\lambda_1, \lambda_2, \lambda_3, \lambda_4\}$; $M_0 = \{+, -\}$
 (M_0 is the usual binary alphabet; Λ is the set of all strings of a finite number of the "elementary messages" λ_i ($i = 1, 2, 3, 4$) and M is built in the same way with the "letters" + and - .

When coding, we want to establish a correspondence between Λ and some subset M' of M satisfying two conditions :

1) to every $\lambda \in \Lambda$ corresponds at least one $\mu \in M'$ ("total" character of the coding)

2) to any distinct $\lambda, \lambda' \in \Lambda$ must correspond distinct $\mu, \mu' \in M'$ in order that the deciphering be free from ambiguity.

A priori any one to one correspondance between Λ and a subset M' from M would do - but usually this could imply that we cannot proceed to the sending of the message before we know it in its totality. So a further practical condition - which is not too easy to formulate rigourously - could be :

3) For a reasonably large number of messages λ the coding is such that for any right multiple λ' of λ (i.e. any $\lambda' = \lambda \lambda''$) the signals μ and μ' have a reasonably long common left divisor μ_1 (i.e. are of the form : $\mu = \mu_1 \mu_2$ and $\mu' = \mu_1 \mu_2'$).

The simplest way of fulfilling these desiderata is to assign to each $\lambda_{i0} \in \Lambda_0$ a string of binary letters μ_i (which very conveniently we may too call a word) and for any sequence $\lambda_1, \lambda_2, \dots, \lambda_{im}$ to send the corresponding sequence: $\mu_1, \mu_2, \dots, \mu_{im}$.

For example, with the correspondance : \mathcal{C}_1 :
 $\lambda_1 \rightarrow + = \mu_1$; $\lambda_2 \rightarrow +- = \mu_2$; $\lambda_3 \rightarrow -+- = \mu_3$;
 $\lambda_4 \rightarrow -- = \mu_4$

we would have :

$\lambda_2 \lambda_3 \lambda_1 \lambda_2 \lambda_4 \rightarrow + + - - + - + + - -$

It is not obvious however how the set M'_0 of the words μ_i has to be selected so that decoding be free from ambiguity :

At my knowledge, the question has been raised first and practically solved by Sardinas and Patterson in a pioneering paper(1).

With the help of semi group concepts we may however obtain a deeper insight into their whole procedure which was purely logical :

We are looking for a total translation from Λ to M and it is quite axiomatic that the decoding is unambiguous if and only if the sub semi group M'_0 generated by M'_0 is isomorphic to the free semigroup Λ - or - for short - that M' is a free subsemigroup of M .

Algebraic consequences of this simple remark are to be found in appendix 1.

Now would come a fourth requirement : (admissibility)

4) The length of the words μ_i must be as small as possible in respect of some a priori probability distribution on Λ .

As a matter of fact (4) will be met incidentally, so to say, in view of another condition we put in definition III :

That the translation from \mathcal{M} back to Λ be sub total :

What this means exactly is that any sequence μ of binary digit be a left divisor of at least one message $\mu' \in M'$ which can be completely and exactly retranslated into Λ .

This condition together with the possibility of one-to-one deciphering implies automatically that the code be unitary (as defined below)(see appendix O), and admissible in that sense that it meets the optimality requirement (4) in respect of at least one a priori probability distribution of the words. (*)

3. Discussion of the decoding methods : scansion

This being settled we have to look more closely at the decoding.

For avoiding repetition let us observe that Λ does not play any role by itself since the $\lambda_i \in \Lambda_0$ are in a one-to-one correspondance with the words $\mu_i \in M_0$. So we may perfectly well dispense from mentioning it altogether.

But in order to stress when a given string μ of binary symbols is really a set made up of a sequence of words and not any odd sequence of + and - we shall say that μ is a complete message (for instance : " +- - - + - " = $\mu_2 \mu_3$ is a complete message, but " + - - - " is not) and indicate it by enclosing it into two / signs, which shall denote too, end and beginning of the words.

Let us try to decode the following complete message in code \mathcal{C}_1 :

| + + - - + - - + - - + /
The only way open is trial and error : the first + may be :

- either μ_1 itself
- either the first letter from $\mu_2 = + + - - /$

so that we have the choice between :

| + | + - - + - - + - - + | and | + + - - / - - - + - - + /

In the first case no further doubt comes in and we are lead to :

| + | + - - | - - | + - - | + - - | = $\mu_1 \mu_1 \mu_4 \mu_1 \mu_4 \mu_1 \mu_4 \mu_1$

(*)If M is the free semi group of all phonemic sequences in English and M' the sub set of all "semantically correct sentences", M' is neat in M .

(For instance :

" /pri wat law cut chur coco feet .."(obtained from King Lear, Act III, scene I, with Tippet's help) is fitted into a complete message in M' by adding : "... and this, Gentlemen, was, may-be, my best example of a semantically void utterance /")

In the second we obtain :

| + + - - | - - | - - | - - | = $\mu_2 \mu_3 \mu_3 - +$

Since here - + is left at loose end (strictly speaking) the first translation was the good one, being known that the transmission is over. Observe that if, on the contrary, the signal was the same as before except for an added terminal - digit, the conclusion would be exactly opposite :

| + + - - | - - | - - | - - |

is the only fitting "scansion" as we could say by borrowing from prosody this term for its classical flavour.

So the inverse translation from M back to Λ does not look like satisfying very reasonably the above condition 3.

An obvious remedy to it would be to limit still more the set M'_0 . B. Mandelbrot, who has first discussed these problems has distinguished several possibilities :

- 1) Uniform codes : in which every word has the same length (i.e. number of letters), this criterium giving a direct scansion (examples: all the noise reducing codes introduced so far except for a proposal of "sequential coding" by Peter Elias(4) and some examples by Lemnael(5).)
- 2) More generally : what we shall call : Unitary codes : i.e. codes in which no word is a left divisor of another word (examples : Fano's, Huffman's, Shannon's codes)
- 3) Natural codes : (introduced by B. Mandelbrot) in which a special letter points out the end of the word (example : most of the spoken or written languages).

Further, Mandelbrot has shown that any unitary code is, at least asymptotically, as good from the point of view of economy of length as any other one. It could seem futile then to care for more extensive classes were we not prompted by other circumstances - and especially by the threat of a noise.

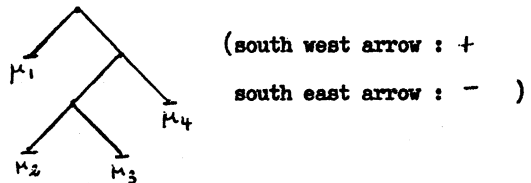
4. Noise absorption and eryodism.

Consider indeed the following code : \mathcal{C}_2

$\mu_1 = + ; \mu_2 = - + + ; \mu_3 = - - + - ; \mu_4 = - -$

(which is, parenthetically, just the previous one with the time arrow inverted)

It is unitary all right so that we may represent it by a "tree" in the familiar fashion :



The "neat" condition (subtotality of the translation from M back to Λ) is reflecting itself in the fact that any branch of the tree

Now the key point is that for any four finite fragments, μ_1, μ_2, μ_3 and μ_4
 $\mu_1 \equiv \mu_2$ and $\mu_3 \equiv \mu_4$ implies $\mu_1 \mu_3 \equiv \mu_2 \mu_4$.

The syntactic equivalence is thus fully compatible with the semi group structure of M and if we consider classes for \equiv (i.e. the subsets of elements from M which are syntactically equivalent between themselves), these classes make a new semi group \bar{M} which is an homomorphic image of M .

$\bar{M} \supset \bar{M}_0$, the fundamental semi group of the coding (f.s.g.) is most usually finite and is easily represented by matrices, but before we explain how, we need still a new concept: that of prefix:

Consider again two fragments μ and μ' but assume, now, that both are beginning at a / mark:

Even if μ and μ' are not syntactically equivalent, it could happen that under this supplement any restriction any further fragment which completes μ into a full message would do the same to μ' :

One could say that " μ and μ' as beginning of messages are syntactically equivalent on the right" (in symbols: $\mu \sim \mu'$)

For example:

$\mu = /--- \dots$ and $\mu' = /+--- \dots$ are not in the relation \equiv (since $/+---$ is a complete message although $/---$ is not complete), but $\mu \sim \mu'$ all the same for $\mu\mu''$ is a complete message if and only if $\mu'' = -/ \dots$ or $+/ \dots$ or $+ - / \dots$ just as well as for μ' .

We call prefixes the classes π_i of fragments for this new relation \sim .

For the code \mathcal{C}_2 , there are three prefixes:

$\pi_1 \ni / \emptyset$ (words and words only are bringing a $\mu \in \pi_1$ into a complete message.

π_2 contains all the words and its existence is typical of unitary coding).

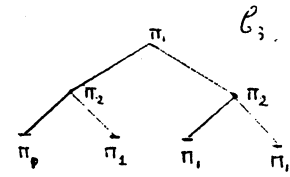
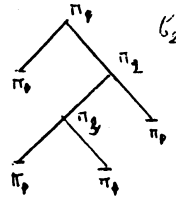
$\pi_2 \ni /- \dots$ (the corresponding right divisors are $-/ \dots$, $+/ \dots$ and $+ - / \dots$)

$\pi_3 \ni /-+ \dots$ (the corresponding right divisors are $+ / \dots$ or $- / \dots$).

Now if $\mu_1 \sim \mu_2$, one proves that

$\mu_1 \mu_3 \sim \mu_2 \mu_3$, too, whatever be μ_3

With unitary codes prefixes correspond to nodes of the tree in a one to many fashion: Two nodes being in relation \sim ("pertain to the same prefix") if the subtrees below them are identical. Such things does not occur in our \mathcal{C}_2 code (see below), but are quite typical of uniform codes.



In the code \mathcal{C}_3 of length 2 ($\mu_1 = ++$; $\mu_2 = +- ; \mu_3 = -+ ; \mu_4 = --$) there is only two prefixes: one, π_1 , corresponding to complete messages - i.e. to sequences with an even number of letters - and another one, π_2 , corresponding to odd length sequences.

6. Matrix representation of the fundamental semi group.

If we have started reading just at the beginning of the transmission, we may consider at any time t the prefix $\pi(t)$ to which pertain the initial fragment till the t -th letter as a "state" which changes at any new letter received.

For instance - apart from any meaning again - the sequence $/+---+---$ corresponds to the following sequence of prefixes:

$\pi_1, \pi_1, \pi_2, \pi_1, \pi_2, \pi_3, \pi_1, \pi_1, \pi_2, \dots$

It is easy to visualise "+" and "-" respectively as the transition matrices:

	π_1	π_2	π_3
π_1	1	0	0
π_2	0	0	1
π_3	1	0	0

(+)

	π_1	π_2	π_3
π_1	0	1	0
π_2	1	0	0
π_3	1	0	0

(-)

(+ lets π_1 invariant since it is a word. It sends π_2 into π_3 and makes a word from π_3 etc..)

These matrices correspond in a one to one fashion to the elements of the fundament semi group, for instance:

1	0	0
1	0	0
1	0	0

(++)

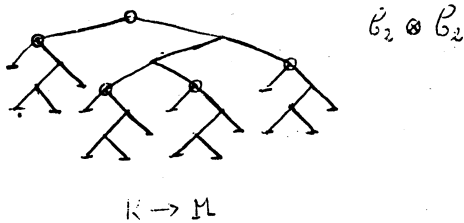
0	0	1
1	0	0
1	0	0

(-+)

filtering. If starting from the assumption that the Λ_i are provided independently with fixed probabilities by the source, we discover later on that, actually, they were just building blocks in some higher degree semantic units (sent again independently of each other as a second approximation) we can preserve at least some of the features of our initial approximation.

But the main point for us here lies in another aspect.

Suppose that the $K \rightarrow \Lambda$ coding be uniform. in general the $K \rightarrow M$ one will not be so, but it will fail to be ergodic just the same, giving us the second of the three exceptional families mentioned above. We shall call such codes "uniformly composed codes". An example is given below :



$K \rightarrow \Lambda$: uniform of length two

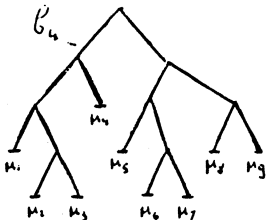
$\Lambda \rightarrow M$: our usual B_2

The nodes indicated with a \circ are the ones corresponding to nodes in the $K \rightarrow \Lambda$ coding.

8. Anagrammatic codes.

Let us come now into the last family. For this we produce the following horrible example : B_4

$M_1 = + + +$; $M_2 = + + -$; $M_3 = + + - -$; $M_4 = + - -$
 $M_5 = - + +$; $M_6 = - + -$; $M_7 = - + - -$; $M_8 = - - -$; $M_9 = - - - -$.



B_4 is not uniform - nor composed uniformly of a smaller code. But it has the property that by inverting its words we found again a unitary code and, indeed, its symmetric image (symmetric in respect of the N.S. line !)

Since ergodicity is somewhat synonymous of irreversibility of time, we are put on the alert by this oddity.

Indeed, absorption is linked very closely with the problem of reading "backward" messages with an inverted code, but, without entering this amusing theory, we can see at once that B_4 and all its family are not ergodic.

If a code is unitary the only sequences, which let Π_1 invariant are the complete messages, whose set is M' . In symbols, this means :

$$\mu_1, \mu_2 \in M' \text{ and } \mu_1 \in M' \text{ implies } \mu_2 \in M'$$

Suppose now that the same property be true on the other direction, i.e. that we had :

$$\mu_1, \mu_2 \in M' \text{ and } \mu_2 \in M' \text{ implies } \mu_1 \in M'$$

Let μ_1 be a complete message which is the unperturbed beginning of the transmission; μ'_1 , its noise corrupted form and μ_2 any other complete message. By the above condition $\mu'_1 \mu_2$ may have a final scansion like that of $\mu_1 \mu_2$ if and only if μ'_1 is a complete message, too.

As this is usually not the case the error will go on till the end.

Codes which are unitary for both directions of time (anagrammatic codes) are not yet fully explored but a construction for various infinite families of them is known. With binary alphabet, there is just the one given above and its symmetric for less than 16 words. It is conjectured that there is still no more than 38 other one below 32 words (on about 10^{16} distinct usual unitary neat codes of this size or less!).

So the family is really exceptionally interesting and deserves further studies since with the uniform and the uniformly composed codes, anagrammatic codes are the only length-bounded codes escaping ergodicity.

References.

1. A.Sardinas and Patterson (1953) Convention records of the I.R.E.
2. B.Mandelbrot .
 1953. Proc.Symp.Comm.Theory. London.
 1954. Proc. Symp. Inf. Network.
 1955. Proc. Symp. Comm. Theory . London.
3. P. Dubreil.
 1941. Mem. Acad. Sci. p. 1-52.
 1951. Rendiconti di Math. 81 p. 289 - 306.
 1953. Bull. Soc. Math. (10) p. 183 - 200.
4. P.Elias.
 1955. Proc. Symp. Comm. Theory. London.
5. A.E.Laemmel.
 1953. ibid.