

CODAGE DÉTERMINISTE IMITANT LE CODAGE ALÉATOIRE
DETERMINISTIC CODING WHICH MIMICS RANDOM CODING

Gérard Battail

Télécom Paris et URA 820 du CNRS, 46, Rue Barrault, 75634 PARIS CEDEX 13

RÉSUMÉ

On montre que la combinaison d'un codage q -aire de Reed-Solomon (n,k) et d'une application des symboles de son alphabet sur une constellation plane symétrique de q points a pour résultat une distribution des distances euclidiennes voisine de celle que l'on obtient en moyenne par codage aléatoire, pourvu que q , n et k soient assez grands. Il est donc possible d'obtenir par des moyens déterministes les mêmes performances sur le canal additif gaussien que par codage aléatoire.

1. Introduction

Un théorème fondamental de la théorie de l'information énonce qu'un canal bruyant étant donné, un procédé de codage approprié permet de rendre la probabilité d'erreur due à son utilisation aussi petite que l'on veut, à condition que l'entropie de la source connectée à son entrée soit inférieure à la capacité de ce canal (il faut, bien entendu, faire des hypothèses de régularité de la source et du canal, en particulier les supposer stationnaires, pour que les grandeurs mentionnées dans cet énoncé existent). En d'autres termes, le bruit du canal limite le débit d'information possible mais non la qualité (spécifiée par la probabilité d'erreur) avec laquelle le message issu de la source peut être restitué au destinataire. La probabilité d'erreur décroît en tendant vers 0 quand la taille des mots du code, et donc la complexité du codage, augmente indéfiniment.

Dans sa démonstration initiale de ce théorème, Shannon a utilisé le concept de *codage aléatoire*. La connaissance de moyens explicites pour obtenir le résultat énoncé faisant défaut, il considérait la famille probabilisée de tous les codes possibles d'un certain type (par exemple, l'ensemble des codes en blocs ayant un nombre donné M de mots, chacun formé d'un nombre donné n de symboles d'un alphabet donné) et en calculait la probabilité d'erreur *moyenne*, ou plutôt une borne supérieure de celle-ci. L'énoncé ci-dessus étant démontré pour cette probabilité moyenne, il existait au moins un code de la famille qui le vérifiait, ayant une probabilité d'erreur inférieure ou égale à la moyenne. Toutes les démonstrations ultérieures du théorème ont repris ce concept.

En outre, il a été établi qu'un code pris au hasard était presque sûrement bon, en ce sens que la probabilité qu'il lui corresponde une probabilité d'erreur très différente de la moyenne tend vers 0 quand la taille du code augmente. Cependant, l'emploi du codage aléatoire est pratiquement

ABSTRACT

We show that combining a q -ary (n,k) Reed-Solomon code with a mapping of its symbols into a 2D symmetric constellation having q points results in a Euclidean distance distribution close to that which results in the average from random coding, provided that q , n and k are large enough. It is therefore possible to obtain by deterministic means the same performance as by random coding over the additive Gaussian channel.

impossible à cause de la complexité de son décodage, qui ne peut être qu'exhaustif. Cette impossibilité de fait a conduit à chercher des procédés de codage explicites. Les résultats ainsi obtenus ont été souvent très inférieurs aux limites établies par la théorie. L'existence d'une incompatibilité fondamentale entre le caractère déterministe d'un codage et son aptitude à s'approcher de ces limites a même été envisagée (*all codes are good, except those we can think of*).

Nous nous proposons de montrer que les codes "séparables à distance maximale" (*maximum distance separable codes*), que nous désignerons en abrégé par "codes MDS", ont des propriétés voisines de celles du codage aléatoire. Les codes binaires de répétition et de parité exceptés, il s'agit de codes non binaires dont la longueur n est limitée supérieurement à une valeur voisine de la taille q de l'alphabet. La similitude avec le codage aléatoire ne concerne que les codes construits sur un alphabet de grande taille, mais des codes longs sont de toute façon nécessaires pour obtenir une probabilité d'erreur petite.

2. Les codes MDS et la distribution de leur poids de Hamming

La distance minimale d d'un code de longueur n construit sur un alphabet de taille q et comportant $M = q^k$ mots vérifie nécessairement la borne de Singleton :

$$d \leq n - k + 1, \quad (1)$$

qu'il soit linéaire ou non [1]. La démonstration de cette inégalité est particulièrement simple si le code est supposé linéaire : on ne restreint pas la généralité en le prenant sous forme systématique et, puisqu'il existe des mots où un seul des k symboles d'information est différent de 0, leur poids est le plus grand possible si aucun des symboles de contrôle n' est nul. Il vaut alors $n - k + 1$.



La distance minimale des codes MDS est la plus grande possible d'après (1), soit $d = n - k + 1$. Cette définition implique l'importante propriété que la donnée de k symboles arbitraires en des positions quelconques spécifie toujours un mot du code [1].

Les seuls codes binaires de ce type sont les codes de parité, pour lesquels k est arbitraire avec $n = k + 1$ et $d = 2$, et les "codes" de répétition, avec $k = 1$ et donc $d = n$, arbitraire. Sur un alphabet de taille $q > 2$, il existe des codes MDS pour lequel on a à la fois $d > 2$ et $k > 1$. Les codes de Reed-Solomon en sont la famille la plus connue. La longueur des mots y est $n = q - 1$, ou $n = q$ pour les codes dits étendus (*extended*) qui s'en déduisent par l'adjonction à chaque mot d'un symbole de contrôle supplémentaire égal à l'opposé de la somme de tous ses symboles.

La distribution des poids de Hamming des codes MDS est assez aisément calculée [2,3]. Le nombre A_j des mots de poids j a pour expression :

$$\begin{aligned} A_0 &= 1, \\ A_i &= 0, \quad 0 < i < n-k+1, \\ A_j &= C_n^j \sum_{i=0}^{j-1-(n-k)} (-1)^i C_j^i (q^{j-i-(n-k)} - 1) \end{aligned} \quad (2)$$

pour $n - k + 1 \leq j \leq n$. Avec j dans cette plage de valeurs, on peut encore écrire (2) sous la forme

$$\begin{aligned} A_j &= q^{-(n-k)} C_n^j (q - 1)^j \\ &- C_n^j \left[\sum_{i=j-(n-k)}^j (-1)^i C_j^i q^{j-1-(n-k)} + \sum_{i=0}^{j-1-(n-k)} (-1)^i C_j^i \right] \end{aligned} \quad (3)$$

où

- le premier terme est le nombre moyen de mots de poids j résultant d'un tirage aléatoire de q^k mots (chacun des n symboles d'un mot étant tiré avec une égale probabilité, indépendamment des précédents, tous les mots étant déterminés de la sorte indépendamment les uns des autres) ;
- le second terme est négligeable devant le premier si q et k sont suffisamment grands. En outre, l'approximation qui consiste à ne conserver que le premier terme de (3) est d'autant meilleure que j est plus grand.

Pour un code MDS, il suffit donc de choisir q et k suffisamment grands pour obtenir une distribution des distances de Hamming voisine de celle que l'on obtient en moyenne par codage aléatoire.

3. Description d'un système combinant codage et modulation

Nous considérons maintenant un système de communication pour le canal gaussien combinant un codage linéaire à une modulation décrite par une application bijective de l'alphabet sur une constellation plane de q points, avec $q \gg 2$, commune aux n symboles du code.

La distribution des distances entre les mots d'un code linéaire se réduit à celle de chacun de ses mots par rapport au mot nul, c'est-à-dire à la distribution de ses poids, beaucoup plus facile à calculer. Afin que cette propriété subsiste pour les distances euclidiennes, après modulation, nous supposons que la constellation utilisée est *symétrique*, en ce sens que l'ensemble des distances entre l'un de ses points et tous les autres y est indépendant du point considéré. La modulation

de phase satisfait évidemment à cette condition.

Le polynôme énumérateur de poids d'un code C (dont les coefficients sont donnés par (2) s'il est de type MDS) n'en constitue une description satisfaisante que si le poids de Hamming est le seul paramètre pertinent. Une description plus complète, particulièrement pertinente quand on envisage la modulation, est offerte par le *polynôme descripteur* introduit et étudié dans [4]. C'est un polynôme à n indéterminées x_1, x_2, \dots, x_n , dont nous désignons symboliquement l'ensemble par X . Un mot du code, soit $c = [c_1 \ c_2 \ \dots \ c_n]$, y est représenté par le monôme $x_1^{c_1} x_2^{c_2} \dots x_n^{c_n}$, noté symboliquement X^c . Le polynôme descripteur est la somme des monômes ainsi associés à tous les mots du code C :

$$L_C(X) = \sum_{c \in C} X^c. \quad (4)$$

Pour un code linéaire et une application symétrique, la distribution des distances euclidiennes entre les signaux associés aux différents mots se réduit à celle de leurs distances par rapport au signal qui représente le mot nul. Toutes les propriétés de distance du système ainsi défini peuvent être directement déduites du polynôme descripteur du code : le polynôme énumérateur des distances euclidiennes s'en déduit en y remplaçant chacune des indéterminées par une indéterminée unique x et chacun des symboles en exposant par le poids euclidien que l'application utilisée lui associe. On définit ce poids comme le carré de la distance euclidienne entre les points qui représentent le symbole considéré et le symbole 0.

Il est clair que le polynôme descripteur d'un code est équivalent à la liste de ses mots. Il est donc aussi complexe d'écrire le premier que la seconde : impossible en fait pour un code de taille utile. Le polynôme descripteur prend tout son intérêt quand le code est linéaire, car il devient possible d'en donner des expressions littérales condensées. Certaines d'entre elles permettent d'explicitier les relations entre le code considéré, son dual et l'ensemble des n -uplets, ce qui sous diverses variantes constitue une forme généralisée de l'identité de McWilliams [4].

Soit une matrice génératrice du code C de forme systématique

$$G = [I_k \ | \ P], \quad (5)$$

où I_k est la matrice unité d'ordre k et P une matrice à k lignes et $n-k$ colonnes. Alors, l'une des expressions du polynôme descripteur de C , convenablement généralisée à partir de [4], s'écrit :

$$q^{n-k} L_C(X) = B_C(X_u \ X'_s), \quad (6)$$

où le polynôme $B_C(X)$ est défini par :

$$B_C(X) = \sum_{u \in F_q^k} \sum_{s \in F_q^{n-k}} X_u^u X_s^s \chi(uPs^T), \quad (7)$$

où :

- la notation X_u désigne collectivement les k premières indéterminées associées aux symboles d'information et X_s les $n-k$ dernières, associées aux symboles de contrôle ;
- u et s sont les vecteurs, représentés par des matrices lignes, constitués des k premiers et $n-k$ derniers symboles d'un mot, respectivement ; l'indice supérieur T note la transposition ;



- la matrice \mathbf{P} a été introduite en (5) ;
- $\chi(\cdot)$ est un caractère non banal sur le groupe additif du corps \mathbb{F}_q ;
- X_s^i est un polynôme déduit du monôme X_s en substituant à ses facteurs le résultat de la transformation que spécifie \mathbf{W} , matrice carrée d'élément général $w_i^j = \chi(-\beta_i \beta_j)$, $0 \leq i, j < q$, i étant l'indice d'une ligne et j celui d'une colonne, les exposants formels des indéterminées étant interprétés ici comme des indices repérant les composantes d'un vecteur. Si q est premier, \mathbf{W} est la matrice de transformation de Fourier discrète d'élément général $w_i^j = [\exp(-2\pi i j^{q-1} / q)]$; si le corps est de caractéristique 2, \mathbf{W} est une matrice d'Hadamard, à éléments réels +1 et -1.

On notera que le polynôme $B_C(X)$ est commun en fait au code C et à son dual C' , pour lequel il suffit d'échanger les rôles des vecteurs \mathbf{u} et \mathbf{s} , le premier désignant alors les symboles de contrôle et le second les symboles d'information. Ce polynôme constitue une représentation du classique tableau de décodage (*standard array*).

En effet, chacun des q^n termes du polynôme $B_C(X)$ dépend, d'après (7), de l'un des q^k vecteurs \mathbf{u} et de l'un des q^{n-k} vecteurs \mathbf{s} . On peut donc le représenter à l'aide d'un tableau comportant q^k lignes et q^{n-k} colonnes indexées respectivement par les vecteurs d'information et les vecteurs de contrôle, contenant à l'intersection de la ligne associée à \mathbf{u} et de la colonne associée à \mathbf{s} le coefficient dans (7) qui dépend de ces vecteurs, soit $\chi(\mathbf{u}\mathbf{P}\mathbf{s}^T)$. Bien entendu, ce tableau est immense pour des ordres de grandeur raisonnables et son écriture relève alors de "l'expérience de pensée" ; la figure 1 en donne un exemple, dans un cas très simple.

	00	01	02	03	10	11	12	13	20	21	22	23	30	31	32	33
00	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
01	0	0	1	1	1	1	0	0	1	1	0	0	0	0	1	1
02	0	1	1	0	1	0	0	1	0	1	1	0	1	0	0	1
03	0	1	0	1	0	1	0	1	1	0	1	0	1	0	1	0
10	0	1	1	0	0	1	1	0	1	0	0	1	1	0	0	1
11	0	1	0	1	1	0	1	0	0	1	0	1	1	0	1	0
12	0	0	0	0	1	1	1	1	1	1	1	1	0	0	0	0
13	0	0	1	1	0	0	1	1	0	0	1	1	0	0	1	1
20	0	1	0	1	1	0	1	0	1	0	0	1	0	0	1	0
21	0	1	1	0	0	1	1	0	0	1	1	0	0	1	1	0
22	0	0	1	1	0	0	1	1	1	1	0	0	1	1	0	0
23	0	0	9	0	1	1	1	0	0	0	0	0	1	1	1	1
30	0	0	1	1	1	1	0	0	0	0	1	1	1	1	0	0
31	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1
32	0	1	0	1	0	1	0	1	0	1	0	1	0	1	0	1
33	0	1	1	0	1	0	0	1	1	0	0	1	0	1	1	0

Figure 1. Tableau représentant le polynôme $B_C(X)$, défini dans le texte par (7), pour le code MDS (4,2) de génératrice

$$\mathbf{G} = \begin{pmatrix} 1 & 0 & \alpha & \alpha^2 \\ 0 & 1 & \alpha^2 & \alpha \end{pmatrix}$$

construit sur le corps \mathbb{F}_4 , dont α désigne un élément primitif. Dans la figure, il a été noté 2 et son carré α^2 a été noté 3, ce qui est conforme à la représentation binaire des éléments. Les coefficients du tableau prennent les deux valeurs +1 et -1 et sont représentés sur la figure par l'exposant 0 ou 1 de (-1). Les lignes et les colonnes y ont été repérées par les différents k -uplets et $(n-k)$ -uplets (et non par les monômes correspondants). Les éléments de ce tableau qui correspondent aux mots du code ont été imprimés en italique gras.

On peut interpréter l'expression (6) comme signifiant que le codage opère une *décimation* sur l'ensemble des n -uplets, réduisant leur nombre de q^n à q^k , car la transformation des monômes X_s^i en produits X_s^j de polynômes entraîne que leur somme sur une ligne du tableau, pondérée par les coefficients de (7), se réduit à un unique monôme dont le coefficient est q^{n-k} . Ce monôme survivant représente le vecteur de contrôle \mathbf{r} tel que

$$\mathbf{r} = \mathbf{u}\mathbf{P}, \quad (8)$$

ce qui exprime simplement l'appartenance au code du vecteur $\mathbf{u}\mathbf{G}$, où \mathbf{G} est la génératrice (5).

4. Cas des codes MDS

Nous examinons maintenant les propriétés de ce tableau lorsque que le code utilisé est de type MDS. La matrice \mathbf{P} de la relation (8) est alors telle qu'aucune sous-matrice carrée qui en est extraite n'est singulière [1] (\mathbf{P} est dite de ce fait "superrégulière" [5]). Cela implique en particulier qu'aucun des éléments de \mathbf{P} n'est nul et qu'il existe une correspondance bijective entre chacune des lignes distinctes du tableau et le monôme X_s^i survivant. La superrégularité de la matrice \mathbf{P} entraîne que $\mathbf{u}\mathbf{P}$, et donc \mathbf{r} vérifiant (8), n'est jamais nul pour $\mathbf{u} \neq \mathbf{0}$.

Nous avons déjà rappelé la propriété des codes MDS que tout ensemble de k symboles en des positions arbitrairement choisies détermine un unique mot du code. Il s'ensuit que, par combinaison linéaire de ses lignes, on peut transformer une matrice génératrice \mathbf{G} du code donnée par (5) en une autre, qui se déduit d'une matrice de même forme (5), avec \mathbf{P} superrégulière, par permutation de ses colonnes.

On vérifie immédiatement que la superrégularité de la matrice \mathbf{P} entraîne pour le tableau défini ci-dessus que, si $2k \geq n$ (l'énoncé qui est vrai si $2k < n$ se déduit du suivant en échangeant les mots "lignes" et "colonnes" d'une part, les entiers k et $n-k$ d'autre part) :

- q^{n-k} lignes distinctes y sont présentes car les matrices-lignes $\mathbf{u}\mathbf{P}$ dans (7) correspondent à tous les $(n-k)$ -uplets distincts ;
- chacune des lignes distinctes y apparaît le même nombre de fois, q^{2k-n} ;
- les colonnes y sont toutes distinctes ;
- une seule des lignes (colonnes) distinctes du tableau a tous ses éléments égaux à 1 ; toutes les autres lignes (colonnes) distinctes sont équilibrées en ce sens que la somme de leurs éléments est nulle.

5. Les codes MDS imitent le codage aléatoire

Sous sa forme la plus directe, le codage aléatoire consiste à choisir chacun des symboles d'un mot avec une même probabilité indépendamment les uns des autres, tous les M mots du code étant déterminés ainsi indépendamment les uns des autres. Des procédés plus raffinés combinent un tirage de ce type avec une sélection (*expurgation*) qui améliore la distribution de distances obtenue.

Une variante possible du codage aléatoire consiste en un codage *systématique*, en ce sens que les k premiers symboles de chaque mot en sont les symboles d'information. Soit une table ayant q^k lignes indexées par tous les vecteurs d'information possibles et q^{n-k} colonnes indexées par tous



les vecteurs q -aires possibles ayant $n - k$ composantes. Cette table comporte q^n entrées. Nous faisons correspondre à chacune d'entre elles le n -uple résultant de la concaténation du vecteur u et du vecteur s qui repèrent respectivement la ligne et la colonne auxquelles elle appartient. Un code systématique ayant $M = q^k$ mots peut en être déduit par un procédé de décimation consistant à ne conserver qu'une entrée par ligne pour désigner les n -uples appartenant au code, faisant ainsi correspondre un vecteur de contrôle s à chacun des vecteurs d'information u . Cette décimation est *aléatoire* si le choix de l'entrée conservée dans chaque ligne est fait au hasard, avec une égale probabilité pour toutes, indépendamment d'une ligne à la suivante. Des procédures de sélection peuvent être envisagées pour mieux répartir les distances ainsi obtenues.

Nous avons décrit ci-dessus le codage d'une façon très semblable. Cependant, la règle de décimation, qui désigne l'entrée survivante de chaque ligne du tableau, n'y est pas aléatoire, puisqu'elle s'exprime par la relation déterministe (8), mais d'une complexité qui peut être évaluée en remarquant que la superrégularité de la matrice P est invariante par toute permutation de ses lignes ou de ses colonnes. L'existence d'une seule matrice superrégulière entraîne donc son appartenance à un ensemble d'environ $k!(n-k)!$ matrices équivalentes (en fait un peu moins, car elles ne sont pas toutes distinctes). Or, ce nombre dépasse largement celui des $(n-k)$ -uples qu'il s'agit d'associer aux k -uples d'information si n et k sont suffisamment grands. Par exemple, si $n = 2k = q$, l'approximation par la formule de Stirling du nombre de matrices superrégulières équivalentes conduit à $2\pi k(k/e)^{2k}$, nombre qui croît plus rapidement en fonction de k que celui des $(n-k)$ -uples de contrôle, qui est ici de $(2k)^k$.

Une dépendance complexe, comme celle que spécifie (8), imite de ce fait un choix aléatoire. Certains auteurs comme Chaitin assimilent d'ailleurs le hasard à la complexité et fondent une théorie de l'information algorithmique sur une mesure de l'information associée à un objet par la taille minimale du programme qui le décrit [6].

En outre, les propriétés énoncées subsistent, comme nous l'avons vu, quel que soit l'ensemble de k symboles dans le mot qui est choisi pour indexer les lignes du tableau. Toutes les positions de symbole jouent donc exactement le même rôle, propriété de symétrie qui est évidemment vérifiée en moyenne par le codage aléatoire.

Nous sommes donc en mesure de confirmer l'assertion énoncée dans [7], où elle n'était que le constat de résultats obtenus par simulation : si les paramètres q , n et k sont suffisamment grands, la combinaison d'un codage MDS et d'une modulation symétrique conduit à une distribution des distances euclidiennes très voisine de celle qui résulterait en moyenne du codage aléatoire. Les codes MDS, très structurés, imitent donc paradoxalement le codage aléatoire.

6. Perspectives et conclusion

Contrairement à ce que l'on a cru, il existe ainsi des moyens de codage explicites dont les propriétés intrinsèques sont très proches de celles du codage aléatoire ; ils sont même connus depuis plus de 30 ans. Leur utilisation doit en principe permettre d'approcher la capacité du canal par un codage déterministe, mais leur exploitation effective dépend surtout

de la possibilité d'un décodage presque optimal dont la complexité reste, cependant, raisonnablement petite.

Ces résultats suggèrent notamment la possibilité d'un système de modulation-codage destiné à un canal additif gaussien combinant un codage MDS et une représentation des symboles de son alphabet par une application bijective sur une constellation symétrique. Sous peine d'une dégradation inacceptable, le décodage doit alors être pondéré. On peut envisager pour cela l'algorithme de type séquentiel décrit dans [8-10].

Les performances accessibles ainsi peuvent d'autre part être estimées en assimilant la distribution des distances euclidiennes à celle qui serait obtenue en moyenne par codage aléatoire sans contrainte sur la position des points représentatifs des signaux. Des calculs de Shannon supposant le décodage exactement optimal [11] ont été repris à cet effet. Leur validité implique que l'on néglige la quantification de la position des points due à l'emploi d'un alphabet fini, hypothèse légitime quand la norme du vecteur bruit est supérieure en moyenne à la distance minimale entre ces points, donc quand le codage est utile.

Références

- [1] R.C. SINGLETON, Maximum distance q -nary codes, IEEE Trans., vol IT-10, avril 1964, pp. 116-118
- [2] F.J. McWILLIAMS et N.J.A. SLOANE, The theory of error-correcting codes, North-Holland, Amsterdam, 1977
- [3] K.M. CHEUNG, More on the decoder error probability for Reed-Solomon codes, IEEE Trans. on Inf. Th., vol. 35, n° 4, juil. 1989, pp. 895-900
- [4] G. BATAIL, Description polynomiale des codes en blocs linéaires, Annales Téléc., vol. 38, n° 1-2, jan.-fév. 1983, pp. 3-15
- [5] R.M. ROTH et A. LEMPEL, On MDS codes via Cauchy matrices, IEEE Trans. on Inf. Th., vol. 35, n° 6, nov. 1989, pp. 1314-1319
- [6] G.J. CHAITIN, Algorithmic information theory, Cambridge Univ. Press, 1987
- [7] G. BATAIL, H. MAGALHÃES de OLIVEIRA et ZHANG W., Coding and modulation for the Gaussian channel, in the absence or in the presence of fluctuations, EUROCODE 90, Udine, Italie, 5-9 nov. 1990
- [8] G. BATAIL, Décodage pondéré optimal des codes linéaires en blocs I.- Emploi simplifié du diagramme du treillis, Annales Télécommunic., vol. 38, n° 11-12, nov.-déc. 1983, pp. 443-459
- [9] G. BATAIL et J. FANG, Décodage pondéré optimal des codes linéaires en blocs II. - Analyse et résultats de simulation, Annales Télécommunic., vol. 41, n° 11-12, nov.-déc. 1986, pp. 580-604
- [10] J. FANG, Décodage pondéré des codes en blocs et quelques sujets sur la complexité du décodage, Thèse de Docteur de l'ENST, soutenue le 18 mars 1987
- [11] C.E. SHANNON, Probability of error for optimal codes in a Gaussian channel, BSTJ, vol 38, n° 3, mai 1959, pp. 611-656