

Substituição Homofônica Ótima com Restrição

Valdemar C. da Rocha Jr., Cecilio Pimentel e Marcos Müller Vasconcelos

Resumo—A substituição homofônica é uma técnica através da qual cada símbolo de uma fonte é representado por um ou mais símbolos denominados homofonemas. Este artigo apresenta um esquema ótimo de substituição homofônica, para tratar o caso no qual cada palavra de homofonema tem como símbolos variáveis aleatórias independentes e identicamente distribuídas, obedecendo uma distribuição de probabilidade arbitrária. Um algoritmo é apresentado para decompor as probabilidades dos símbolos de uma fonte, o qual produz a cada passo o conjunto de homofonemas de mínima entropia. É apresentada uma prova da otimalidade deste algoritmo no sentido de que minimiza a redundância resultante.

Palavras-Chave—Substituição homofônica, criptografia, teoria da informação.

Abstract—Homophonic substitution is a technique by which each source symbol is represented by one or more symbols called homophones. This paper presents an optimum homophonic substitution scheme to handle the case in which the symbols of each homophonic codeword are independent and identically distributed random variables, distributed according to an arbitrary probability distribution. An algorithm is presented for decomposing the source symbol probabilities which produces at each iteration the set of homophones of least entropy. A proof is presented of the optimality of this algorithm in the sense that it minimizes overall redundancy.

Keywords—Homophonic substitution, cryptography, information theory.

I. INTRODUÇÃO

Substituição homofônica [1], [2] é uma técnica criptográfica para reduzir a redundância de uma mensagem a ser cifrada, ao custo de uma certa expansão do texto claro. Esta técnica consiste na substituição (*one-to-many*) de cada letra da mensagem original por um substituto ou *homofonema*, em um alfabeto maior, para formar a mensagem de texto claro que é então cifrada [3]. Cada homofonema é representado (*one-to-one*) por uma palavra código, cujos símbolos geralmente são uniformemente distribuídos e estatisticamente independentes, e que conseqüentemente tornam a cifra mais segura por aumentar a sua distância de unicidade [4].

A fim de simplificar nosso tratamento do assunto, consideraremos apenas a substituição homofônica aplicada na seqüência de saída U_1, U_2, U_3, \dots de uma fonte discreta sem

memória (DMS) K -ária, porém a teoria apresentada aplica-se também a fontes com memória, bastando apenas substituir a distribuição de probabilidade de U_i pela distribuição de probabilidade condicional de U_i dados os valores observados de U_1, U_2, \dots, U_{i-1} . Para uma DMS K -ária o problema da substituição homofônica reduz-se ao caso de uma única variável aleatória U . Seja U uma variável aleatória assumindo valores no conjunto finito $\{u_1, u_2, \dots, u_K\}$. Vamos supor, sem perda de generalidade essencial, que todos os K valores de U possuem probabilidade não-nula de ocorrência e que $K \geq 2$. O homofonema V para U assume valores no conjunto $\{v_1, v_2, \dots\}$, o qual pode ser finito ou infinito contável, e é caracterizado pelo fato de que para cada j existe exatamente um i tal que $P(V = v_j | U = u_i) \neq 0$. Para a substituição homofônica D -ária de comprimento variável, ao homofonema V é associada uma *palavra de homofonema* (X_1, X_2, \dots, X_W) , na qual X_i é uma variável aleatória D -ária e o seu comprimento W é também em geral uma variável aleatória. Em geral requer-se que as palavras de homofonema sejam alocados de tal modo que (X_1, X_2, \dots, X_W) seja uma codificação de V livre de prefixo, i.e., tal que as palavras de homofonema (x_1, x_2, \dots, x_w) sejam todas distintas e que nenhuma delas seja prefixo de uma outra. Se os componentes X_1, X_2, \dots, X_W da palavra associada ao homofonema V forem variáveis aleatórias D -árias estatisticamente independentes e uniformemente distribuídas, a substituição homofônica é dita ser *perfeita* [2]. A substituição homofônica foi definida em [2] como *ótima* se ela for perfeita e se minimizar o comprimento médio $E(W)$ das palavras de homofonema, sobre substituições homofônicas perfeitas [5].

Em 2001 foi introduzido um algoritmo [6] (*vide Apêndice*) para realizar substituição homofônica com alfabeto binário, no qual os símbolos em cada palavra de homofonema são variáveis aleatórias independentes e identicamente distribuídas, obedecendo uma distribuição de probabilidade arbitrária. A importância deste algoritmo decorre do fato dele satisfazer um limitante superior para a codificação homofônica binária ótima livre de prefixo, o qual estipula que a diferença entre a entropia dos homofonemas e a entropia da fonte não é maior que $h(p)/p$ bits [6], sendo $h(\cdot)$ a função entropia binária e $p \geq 1/2$ uma das probabilidades dos símbolos das palavras de homofonema. No que vem a seguir iremos nos referir a este algoritmo como o *algoritmo de máxima entropia*, ou algoritmo MAX-ENT, por razões que ficarão claras mais adiante. O algoritmo MAX-ENT representa uma solução para o problema formulado por Knuth [7, p.427], sobre a geração de distribuições de probabilidade usando uma moeda viciada, e nos foi comunicado por Julia Abrahams [8]. Nossa motivação para desenvolver o presente trabalho decorreu da observação de que há casos, como é ilustrado na *Seção 2, Exemplo 1*, nos

Grupo de Pesquisa em Comunicações - CODEC, Departamento de Eletrônica e Sistemas, Caixa Postal 7800 Universidade Federal de Pernambuco, Recife 50711-970 PE, BRASIL. E-mails: vcr@ufpe.br, cecilio@ufpe.br, fuzz@terra.com.br. Os autores agradecem ao CNPq pelo apoio parcial recebido através dos projetos 304214/77-9, 300987/96-0 e PIBIC 30161, respectivamente.

quais pudemos realizar a substituição homofônica de modo mais eficiente que pelo uso do algoritmo MAX-ENT. Na *Seção 2* nós também provamos dois lemas os quais indicam um modo para expandir o alfabeto da fonte, com incremento mínimo na entropia do alfabeto expandido. Na *Seção 3* nós introduzimos um algoritmo para realizar a substituição homofônica D -ária com restrição, ou seja, na qual os símbolos das palavras de homofonema são variáveis aleatórias D -árias ($D \geq 2$), e cuja distribuição de probabilidade não é necessariamente uniforme. Este procedimento é então demonstrado ser ótimo no sentido de que minimiza a redundância total. Finalmente, na *Seção 4* nós apresentamos algumas conclusões sobre este trabalho.

II. MOTIVAÇÃO E RESULTADOS BÁSICOS

Apresentaremos a seguir um exemplo ilustrando que em algumas situações é possível realizar a substituição homofônica de modo mais eficiente que usando o algoritmo MAX-ENT. Esta foi a nossa maior motivação para tentar melhorar o desempenho do algoritmo MAX-ENT.

Exemplo 1: Seja U uma DMS binária com $P_U(u_1) = 5/9$ e $P_U(u_2) = 4/9$. Consideremos a substituição homofônica binária perfeita aplicada a U quando $\Pi_2 = \{2/3, 1/3\}$ for a distribuição de probabilidade dos símbolos das palavras de homofonema. Aplicando o algoritmo MAX-ENT [6] (vide Apêndice) obtemos

$$P_U(u_1) = 4/9 + \sum_{i=0}^{\infty} 8/3^{4+2i}$$

$$P_U(u_2) = 1/3 + 2/27 + \sum_{i=0}^{\infty} 8/3^{5+2i},$$

que produz um comprimento médio $E(W) = 19/9$ para as palavras de homofonema e uma redundância de $E(W) - H(U) = 2,111 - 0,991 = 1,12$ bits. Por outro lado, por tentativa e erro obtivemos

$$P_U(u_1) = 2/9 + 3/9$$

$$P_U(u_2) = 4/9,$$

que produz um comprimento médio de $5/3$ para as palavras de homofonema e uma redundância de $E(W) - H(U) = 1,667 - 0,991 = 0,676$ bits, i.e., uma redundância representando 60% daquela obtida com o algoritmo MAX-ENT.

Os dois lemas a seguir provêm a base para a construção do algoritmo que apresentaremos na *Seção 3*, para a realização da substituição homofônica D -ária ótima com restrição.

Lema 1: Seja U uma DMS K -ária com distribuição de probabilidade

$$P_U = \{P_U(u_1), P_U(u_2), \dots, P_U(u_K)\}$$

e entropia $H(U)$. A entropia $H(V)$ associada à distribuição de probabilidade

$$P_V = \{P_U(u_1), P_U(u_2), \dots, \beta P_U(u_i), (1 - \beta)P_U(u_i), \dots, P_U(u_K)\},$$

contendo $K + 1$ termos, obtida expandindo-se o símbolo u_i de U em dois novos símbolos com probabilidades $\beta P_U(u_i)$ e

$(1 - \beta)P_U(u_i)$, respectivamente, onde $0 < \beta < 1$, é dada por

$$H(V) = H(U) + P_U(u_i)h(\beta),$$

na qual $h(\beta)$ denota a função entropia binária.

Demonstração: Seja $\bar{\beta} = 1 - \beta$. Nós provaremos este lema expandindo e simplificando a expressão de $H(V)$ do seguinte modo

$$\begin{aligned} H(V) &= - \sum_{j=1, j \neq i}^K P_U(u_j) \log P_U(u_j) \\ &\quad - \beta P_U(u_i) \log [\beta P_U(u_i)] \\ &\quad - \bar{\beta} P_U(u_i) \log [\bar{\beta} P_U(u_i)] \\ &= - \sum_{j=1, j \neq i}^K P_U(u_j) \log P_U(u_j) \\ &\quad - \beta P_U(u_i) \log \beta \\ &\quad - \beta P_U(u_i) \log P_U(u_i) \\ &\quad - \bar{\beta} P_U(u_i) \log \bar{\beta} \\ &\quad - \bar{\beta} P_U(u_i) \log P_U(u_i) \\ &= - \sum_{j=1}^K P_U(u_j) \log P_U(u_j) \\ &\quad + P_U(u_i) [-\beta \log \beta - \bar{\beta} \log \bar{\beta}] \\ &= H(U) + P_U(u_i)h(\beta). \end{aligned}$$

Lema 2: Seja U uma DMS K -ária com distribuição de probabilidade

$$P_U = \{P_U(u_1), P_U(u_2), \dots, P_U(u_K)\}$$

e entropia $H(U)$. Seja

$$P_{V_1} = \{P_U(u_1), P_U(u_2), \dots, \alpha, P_U(u_i) - \alpha, \dots, P_U(u_K)\} \quad (1)$$

a distribuição de probabilidade obtida pela expansão do símbolo u_i de U em dois novos símbolos com probabilidades α e $P_U(u_i) - \alpha$, respectivamente, e seja

$$P_{V_2} = \{P_U(u_1), P_U(u_2), \dots, \alpha, P_U(u_j) - \alpha, \dots, P_U(u_K)\} \quad (2)$$

a distribuição de probabilidade obtida pela expansão do símbolo u_j de U , $j \neq i$, em dois novos símbolos com probabilidades α e $P_U(u_j) - \alpha$, respectivamente, i.e., ambos P_{V_1} e P_{V_2} contêm $K + 1$ termos cada e $0 < \alpha < \min\{P_U(u_i), P_U(u_j)\}$. A entropia $H(V_1)$, associada com P_{V_1} , é maior que a entropia $H(V_2)$, associada com P_{V_2} , se e somente se $P_U(u_i) > P_U(u_j)$.

Demonstração: Seja $\Delta = P_U(u_i) + P_U(u_j) - \alpha > 0$. Se $P_U(u_i) > P_U(u_j)$ teremos

$$P_U(u_i) > P_U(u_i) - \alpha \geq P_U(u_j) > P_U(u_j) - \alpha$$

ou

$$P_U(u_i) > P_U(u_j) \geq P_U(u_i) - \alpha > P_U(u_j) - \alpha.$$

De modo equivalente,

$$\begin{aligned} P_U(u_i)/\Delta &> [P(u_i) - \alpha]/\Delta \\ &\geq P_U(u_j)/\Delta \\ &> [P(u_j) - \alpha]/\Delta \end{aligned} \quad (3)$$

ou

$$\begin{aligned} P_U(u_i)/\Delta &> P_U(u_j)/\Delta \\ &\geq [P_U(u_i) - \alpha]/\Delta \\ &> [P_U(u_j) - \alpha]/\Delta. \end{aligned} \quad (4)$$

Como $P_U(u_i)/\Delta > [P_U(u_j) - \alpha]/\Delta$ segue que $P_U(u_i)/\Delta > 1/2$, e em vista de (3) e (4) concluímos que $P_U(u_i)/\Delta > \max\{P_U(u_j)/\Delta, [P(u_i) - \alpha]/\Delta\} > 1/2$. Portanto, nós temos

$$h[P_U(u_i)/\Delta] < h[P_U(u_j)/\Delta],$$

i.e.,

$$h[P_U(u_j)/\Delta] - h[P_U(u_i)/\Delta] > 0.$$

Subtraindo $h[P_U(u_i)/\Delta]$ de $h[P_U(u_j)/\Delta]$ obtemos

$$\begin{aligned} h[P_U(u_j)/\Delta] - h[P_U(u_i)/\Delta] &= \\ (1/\Delta) \{ &-P_U(u_j) \log P_U(u_j) \\ &-[P_U(u_i) - \alpha] \log [P_U(u_i) - \alpha] \\ &+ P_U(u_i) \log P_U(u_i) \\ &+ [P_U(u_j) - \alpha] \log [P_U(u_j) - \alpha] \} > 0. \end{aligned}$$

Entretanto, subtraindo $H(V_2)$ de $H(V_1)$ obtemos

$$\begin{aligned} H(V_1) - H(V_2) &= \\ &-P_U(u_j) \log P_U(u_j) \\ &-[P_U(u_i) - \alpha] \log [P_U(u_i) - \alpha] \\ &+ P_U(u_i) \log P_U(u_i) \\ &+ [P_U(u_j) - \alpha] \log [P_U(u_j) - \alpha] \\ &= \Delta \{ h[P_U(u_j)/\Delta] - h[P_U(u_i)/\Delta] \} > 0. \end{aligned} \quad (5)$$

Tendo provado a condição de necessidade, observamos que a prova da condição de suficiência é imediata a partir da equação (5). ■

Fazendo $\alpha = \beta_i P_U(u_i)$ em (1) e fazendo $\alpha = \beta_j P_U(u_j)$ em (2) obtemos pelo Lema 1 $H(V_1) = H(U) + P_U(u_i)h(\beta_i)$ e $H(V_2) = H(U) + P_U(u_j)h(\beta_j)$, ou seja, $H(V_1) - H(V_2) = P_U(u_i)h(\beta_i) - P_U(u_j)h(\beta_j)$. Como consequência de (5) podemos escrever

$$H(V_1) - H(V_2) = P_U(u_i)h(\beta_i) - P_U(u_j)h(\beta_j) \geq 0. \quad (6)$$

Corolário 1: (Aos Lemas 1 e 2) Seja u_r um símbolo da fonte para o qual $P_U(u_r) - \alpha \geq 0$ é um mínimo, $r \in \{1, 2, \dots, K\}$. A fim de expandirmos de um símbolo o alfabeto da fonte, segundo o Lema 2, com o menor incremento possível na entropia do alfabeto expandido resultante, precisamos substituir o símbolo u_r por dois novos símbolos cujas probabilidades são α e $P_U(u_r) - \alpha$, respectivamente.

III. UM ALGORITMO ÓTIMO

Na substituição homofônica D -ária padrão, i.e., aquela na qual $\Pi_D = \{1/D, 1/D, \dots, 1/D\}$, o projetista beneficia-se do fato de que uma dada probabilidade de símbolo $P_U(u_i)$, $0 < P_U(u_i) < 1$, possui essencialmente uma única decomposição na base D . Esta afirmação decorre da observação de que $P_U(u_i)$ ou tem uma única decomposição que é uma soma de potências negativas de D com um número infinito de termos, ou então possui tanto uma decomposição que é uma soma contendo um número finito de potências negativas de D e uma decomposição que é uma soma contendo um número infinito de potências de D , na qual o termo menor de todos na decomposição finita é expandido numa soma contendo um número infinito de potências negativas de D . Por exemplo, para $D = 3$, $P_U(u_i) = 4/9$ pode ser decomposto tanto como $P_U(u_i) = 1/3 + 1/9$, ou como $P_U(u_i) = 1/3 + (1/27) \sum_{i=0}^{\infty} (2/3)^i$.

A substituição homofônica com restrição infelizmente não herda a propriedade de ser essencialmente única a decomposição da probabilidade descrita acima. Isto significa que, a fim de separarmos cada símbolo da fonte em homofonemas, precisamos considerar a cada passo o conjunto total de probabilidades que é produzido a partir das probabilidades dos símbolos da fonte. Lembramos que no caso D -ário, com distribuição uniforme dos símbolos dos homofonemas, podemos considerar a decomposição da probabilidade de cada símbolo isoladamente. Descreveremos a seguir uma maneira de realizar a substituição homofônica D -ária com restrição, na forma de um algoritmo que denominaremos de *algoritmo de mínima entropia* ou, de modo abreviado, algoritmo MIN-ENT.

A. Algoritmo de mínima entropia

Seja $\Pi_D = \{\pi_0, \pi_1, \dots, \pi_{D-1}\}$ a distribuição de probabilidade dos dígitos das palavras de homofonema. Os homofonemas são selecionados como nós terminais na árvore D -ária enraizada T com probabilidades, de tal modo que de cada nó emanam D ramos com probabilidades $\pi_0, \pi_1, \dots, \pi_{D-1}$, respectivamente. O rótulo de um caminho em T é representado por uma seqüência D -ária, formada pelos números inteiros $0, 1, 2, \dots, D-1$, associada aos ramos que constituem este caminho. Denotemos por $\pi_0^{\lambda_0}, \pi_1^{\lambda_1}, \dots, \pi_{D-1}^{\lambda_{D-1}}$ a probabilidade de um caminho em T , de comprimento $\sum_{i=0}^{D-1} \lambda_i$ dígitos, contendo λ_i vezes o dígito i , $0 \leq i \leq D-1$. Para uma dada fonte o algoritmo MIN-ENT simultaneamente encontra a decomposição da probabilidade de cada símbolo da fonte, como uma soma finita ou infinita de termos $\pi_0^{\lambda_0} \pi_1^{\lambda_1} \dots \pi_{D-1}^{\lambda_{D-1}}$, e a correspondente palavra livre de prefixo, na qual o dígito i , $0 \leq i \leq D-1$ ocorre λ_i vezes. Denotemos por $v(i, j)$ o j -ésimo homofonema alocado ao símbolo da fonte u_i , $1 \leq i \leq K$, $j = 1, 2, \dots$, e denotemos por $\alpha(i, j)$ a probabilidade de $v(i, j)$.

Definição 1: Definimos a soma corrente de símbolo $\gamma_m(i)$ para $U = u_i$ na m -ésima iteração do algoritmo MIN-ENT como

$$\gamma_m(i) = P_U(u_i) - \sum_{k=1}^{j-1} \alpha(i, k),$$

com $\gamma_m(i) = P_U(u_i)$ para $j = 0$, na qual j denota o número de homofonemas alocados a u_i até a m -ésima iteração.

Definição 2: Definimos o conjunto de soma corrente Γ_m na m -ésima iteração do algoritmo MIN-ENT como

$$\Gamma_m = \{\gamma_m(i) | \gamma_m(i) > 0, 1 \leq i \leq K\},$$

com $\Gamma_0 = \{P_U(u_1), P_U(u_2), \dots, P_U(u_K)\}$.

Seja $\gamma_{\max} = \max \gamma_m(i) \in \Gamma_m, 1 \leq i \leq K$. Quando $m = 0$ no algoritmo MIN-ENT nós construímos T a partir da raiz, começando com apenas D folhas. Daí então nós expandiremos cada nó terminal em T , cuja probabilidade exceda γ_{\max} , em um número mínimo de ramos suficiente para fazer com que as probabilidades dos nós terminais estendidos resultantes sejam iguais ou menores que γ_{\max} . Chamaremos a árvore resultante de *árvore D -ária enraizada e processada com probabilidades*, T_p . Na m -ésima iteração, $m \geq 1$, um homofonema é alocado a um nó terminal da correspondente T_p , de modo que o nó terminal não utilizado que possua a maior probabilidade, denotada por P_M , seja alocado como um homofonema para o símbolo u_r que apresente o mínimo valor não negativo para a diferença entre sua soma corrente de símbolo $\gamma_m(r)$ e P_M , i.e., tal que $\min_i \{\gamma_m(i) - P_M | (\gamma_m(i) - P_M) \geq 0\} = \gamma_m(r) - P_M \geq 0, 1 \leq i \leq K$. O algoritmo consiste dos seguintes passos.

- 1) Faça $m = 0$. Faça $\gamma_0(i) = P_U(u_i), 1 \leq i \leq K$. Faça $\Gamma_0 = \{P_U(u_1), P_U(u_2), \dots, P_U(u_K)\}$.
- 2) Determine γ_{\max} e construa a árvore T_p para a m -ésima iteração expandindo cada nó terminal não usado na árvore construída para a $(m - 1)$ -ésima iteração, $m \geq 1$, cuja probabilidade exceda γ_{\max} , em um número mínimo de ramos suficiente para fazer a probabilidade dos nós terminais resultantes menores ou iguais a γ_{\max} .
- 3) Encontre em T_p , o caminho não usado E_l cuja probabilidade $P(E_l)$ seja a maior dentre as dos caminhos não usados, i.e., $P(E_l) = P_M$. Denotemos por l o comprimento de E_l .
- 4) Se, para $1 \leq i \leq K, \min_i \{\gamma_m(i) - P_m | (\gamma_m(i) - P_m) \geq 0\} = \gamma_m(r) - P_m \geq 0$, nós então associamos a u_r o homofonema (nó terminal) $v(r, j)$ e a palavra de homofonema D -ária de comprimento l , cujos símbolos constituem o rótulo de E_l em T_p . Isto implica $\alpha(r, j) = P_M$. Compute a soma corrente de símbolo $\gamma'_m(r)$ após esta decomposição e faça $\Gamma'_m = \Gamma_m - \{\gamma_m(r)\}$. Se $\gamma'_m(r) = 0$ então faça $\Gamma_{m+1} = \Gamma'_m$. A decomposição de $P_U(u_r)$ estará agora concluída e conterà j homofonemas, e se $\Gamma_{m+1} = \emptyset$ então FIM. Em caso contrário, i.e., se $\gamma'_m(r) > 0$, então faça $\Gamma_{m+1} = \Gamma'_m \cup \{\gamma'_m(r)\}$.
- 5) Faça $m \leftarrow m + 1$.
- 6) Vá para o passo 2.

Exemplo 2: Seja U uma DMS com $K = 3$ e $P_U = \{53/81, 16/81, 4/27\}$. Consideraremos a seguir a substituição homofônica binária perfeita com restrição aplicada a U quando $\Pi_2 = \{2/3, 1/3\}$ é a distribuição de probabilidade do alfabeto das palavras de homofonema. Aplicando o algoritmo MAX-ENT [6] obtemos

$$53/81 = 4/9 + 4/27 + 4/81 + 2/243 + \sum_{i=0}^{\infty} 8/3^{7+2i}$$

$$16/81 = 4/27 + 4/81$$

$$4/27 = 1/9 + 2/81 + \sum_{i=0}^{\infty} 8/3^{6+2i},$$

que conduz a um comprimento médio das palavras de homofonema de $E(W) = 214/81$ e a uma redundância de $E(W) - H(U) = 2,642 - 1,27 = 1,372$ bits. Por outro lado, aplicando o algoritmo MIN-ENT à fonte dada obtemos

$$P_U(u_1) = 53/81 = 4/9 + 1/9 + 2/27 + 2/81$$

$$P_U(u_2) = 16/81 = 4/27 + 4/81$$

$$P_U(u_3) = 4/27,$$

que conduz a um comprimento médio das palavras de homofonema de $E(W) = 68/27$ e a uma redundância de $E(W) - H(U) = 2,52 - 1,27 = 1,25$ bits, i.e., uma redundância representando 91% daquela obtida com o algoritmo MAX-ENT.

Proposição 1: Seja U uma DMS K -ária com distribuição de probabilidade $P_U = \{P_U(u_1), P_U(u_2), \dots, P_U(u_K)\}$ e entropia $H[P_U(u_1), P_U(u_2), \dots, P_U(u_K)]$. Seja $\Pi_D = \{\pi_1, \pi_2, \dots, \pi_D\}$ a distribuição de probabilidade dos dígitos das palavras de homofonema. A substituição homofônica D -ária perfeita com restrição aplicada a U , da maneira descrita no algoritmo MIN-ENT, minimiza a redundância $E(W) - H(U)$ e portanto é ótima.

Demonstração: É um fato conhecido [10] que a entropia $H(V)$ dos homofonemas na substituição homofônica perfeita está relacionada com o comprimento médio $E(W)$ das palavras de homofonema pela expressão

$$H(V) = H(\pi_1, \pi_2, \dots, \pi_D)E(W).$$

Portanto, ao minimizarmos $H(V)$ estamos também minimizando $E(W)$. A cada passo do algoritmo MIN-ENT o alfabeto da fonte original é estendido de um símbolo e, pelo *Lema 2*, decorre que a entropia do alfabeto estendido sofre o incremento mínimo permitido. Este procedimento é repetido até que a decomposição de U em homofonemas V seja concluída. Como consequência $H(V)$ tem o valor mínimo possível para um dado P_U , o que prova esta proposição. ■

IV. CONCLUSÕES

Introduzimos neste trabalho um algoritmo para realizar a substituição homofônica D -ária ótima com restrição, o qual tem a característica interessante de realizar simultaneamente a seleção das probabilidades dos homofonemas e as respectivas palavras de homofonema. Além disso, deduzimos propriedades deste algoritmo que nos permitiram provar sua otimalidade na minimização da redundância da substituição homofônica D -ária perfeita com restrição. Finalmente, registramos que o algoritmo MIN-ENT produz os mesmos resultados obtidos por tentativa e erro no *Exemplo 1*.

APÊNDICE

Apresentamos a seguir os passos do algoritmo MAX-ENT [6], seguindo a mesma notação da *Seção 3*.

- 1) Faça $m = 0$. Seja Γ_0 o conjunto cujos elementos são as probabilidades dos símbolos $P_U(u_i)$, $1 \leq i \leq K$, ordenados em ordem decrescente.
- 2) Determine γ_{\max} e construa a árvore D -ária enraizada e processada T_p , com probabilidades, para a m -ésima iteração expandindo na árvore construída para a $(m - 1)$ -ésima iteração, $m \geq 1$, cada nó terminal não usado cuja probabilidade exceda γ_{\max} , em um número mínimo de ramos suficiente para fazer a probabilidade dos nós terminais resultantes menores ou iguais a γ_{\max} . Faça $(i, j) = (i, 1)$ e $\gamma_0(i) = P_U(u_i)$, $1 \leq i \leq K$.
- 3) Encontre em T_p o caminho não usado E_l cuja probabilidade $P(E_l)$ é a maior dentre aquelas dos caminhos não-usados, i.e., $P(E_l) = P_M$. Denotemos por l o comprimento de E_l .
- 4) Seja u_r o símbolo da fonte ao qual corresponde a máxima soma corrente de símbolo γ_{\max} na m -ésima iteração. Associe a u_r o homofonema (nó terminal) $v(r, j)$ e a palavra D -ária de homofonema de comprimento l , cujos símbolos constituem o rótulo de E_l in T_p . Isto implica $\alpha(r, j) = P(E_l)$. Faça $(r, j) \leftarrow (r, j + 1)$. Compute a soma corrente de símbolo $\gamma'_m(r)$ após esta decomposição e faça $\Gamma'_m = \Gamma_m - \{\gamma_{\max}\}$. Se $\gamma'_m(r) = 0$ então faça $\Gamma_{m+1} = \Gamma'_m$. A decomposição de $P_U(u_r)$ estará agora concluída e conterà j homofonemas, e se $\Gamma_{m+1} = \phi$ então FIM. Em caso contrário, i.e., se $\gamma'_m(r) > 0$ então faça $\Gamma_{m+1} = \Gamma'_m \cup \{\gamma'_m(r)\}$.
- 5) Faça $m \leftarrow m + 1$.
- 6) Vá para o passo 2.

Concluimos pelo *Lemma 2* que o passo 4, no algoritmo acima, causa o máximo incremento de entropia possível no alfabeto expandido, daí a denominação de algoritmo MAX-ENT.

REFERÊNCIAS

- [1] Ch. G. Günther, "A universal algorithm for homophonic coding", pp. 405-414 in *Advances in Cryptology-Eurocrypt'88*, Lecture Notes in Computer Science, No.330. Heidelberg and New York: Springer, 1988.
- [2] H. N. Jendal, Y. J. B. Kuhn and J. L. Massey, "An information-theoretic approach to homophonic substitution", pp. 382-394 in *Advances in Cryptology-Eurocrypt'89* (Eds. J.-J. Quisquater and J. Vandewalle), Lecture Notes in Computer Science, No.434. Heidelberg and New York: Springer, 1990.
- [3] V. C. da Rocha Jr. and J. L. Massey, "On the entropy bound for optimum homophonic substitution", *Proc. IEEE International Symposium on Information Theory*, Ulm, Germany, 29 June - 4 July, 1997, p.93.
- [4] C. E. Shannon, "Communication theory of secrecy systems", *Bell System Tech. J.*, vol. 28, pp. 656-715, Oct., 1949.
- [5] V. C. da Rocha Jr. and J. L. Massey, "Better than "optimum" homophonic substitution", *Proc. IEEE International Symposium on Information Theory*. Sorrento, Italy, 25 - 30 June, 2000, p. 241.
- [6] V.C. da Rocha Jr. and C. Pimentel, "Binary-constrained homophonic coding", *VI International Symposium on Communication Theory & Applications*, 15 - 20 July 2001, Ambleside, England, pp.263-268.
- [7] D.W. Knuth and A.C. Yao, "The complexity of random number generation", In J.F. Traub, editor, *Algorithms and Complexity: Recent Results and New Directions. Proceedings of the Symposium on New Directions and Recent Results in Algorithms and Complexity*, Carnegie Mellon University, 1976. Academic Press, New York, 1976.
- [8] Julia Abrahams, "Generation of Discrete Distributions form Biased Coins", *IEEE Trans. Inform. Theory*, vol. IT-42, pp.1541-1546, September 1996.
- [9] M. Hoshi and T.S. Han, "Interval algorithm for homophonic coding", *IEEE Trans. Inform. Theory*, vol. IT-47, pp.1021-1031, March 2001.
- [10] J.L. Massey, "Applied Digital Information Theory", *Fach Nr. 35-417 G, 7. Semester*, Class notes at the ETH Zurich, Chapter2, Wintersemester 1988-1989.
- [11] J.L. Massey, "The Entropy of a rooted tree with probabilities", in *Abstracts of Papers, IEEE Int. Symp. on Info. Th.*, 1983, p.127.