

## 文字の出現確率に基づく第2類の文字列類似度

Similarity of strings -2nd kind- based on existence probability of character

飯箸泰宏\*

Yasuhiro IIHASHI

\*明治大学

Meiji University

あらまし：昨年、本学会で、ランダムテキストに基づく文字列類似度を発表し、「飯箸法」として一定の普及を見ることができた。著作権や工業所有権をめぐる争いはますます増加している。それらの争いの解決には、よりいっそう精緻な文字列類似度の計測法が必要である。今回は、ランダムに出現する文字ではなく、対象文字列空間に出現する文字の出現頻度に基づく正規化された新たな尺度を提案する。文字列類似度計測法の飯箸法第2類である。

キーワード：文字列、テキスト、ソースコード、類似度、文字の出現頻度、計算例、塩基配列

### 1 はじめに

昨年、本学会で、ランダムテキストに基づく文字列類似度の計測法を発表した[5]。今回は現実の文字列比較空間に存在する文字の出現確率を基にしたキーテキストを採用することによって、精緻化を図る。

昨年発表したランダムテキストに基づく文字列類似度の計測法[5]は、教育的効果が高いので、「(飯箸法)第1類」として残しておくことにする。

今回は、現実の文字列比較空間に存在する文字の出現確率を基にした、比較空間に共有されるキーテキストを用いて、オリジナリティと類似性に関する正規化された汎用的尺度(飯箸法第2類)を提案する。また、実際に2組のテキストを用いて、この新しい尺度(飯箸法第2類)を用いた計算結果と、昨年発表したランダムテキストに基づく文字列類似度の計測法[5](飯箸法第1類)とLevenshtein距離の3種類の値を比較する。

### 2 文字列類似度の汎用的尺度 テキスト間の相互類似度

$$rS_{a-b} = \frac{rS_{a/b} + rS_{b/a}}{2}$$

### テキスト間の相対類似度

$$rS_{a/b} = \frac{S_{a/b}}{O_a}, rS_{b/a} = \frac{S_{b/a}}{O_b}$$

対象文字列空間に出現する文字の出現頻度に基づく正規化された新たな尺度(飯箸法第2類)における類似度は次のとおりである。

$$\begin{aligned} \text{類似度 } S_{a/b} &= \ln(I_{a/b}) = \sum_{i=1}^{m_a} (n_{k_{a_i}} \times p_{k_{a_i}}) \\ &= \sum_{i=1}^{m_a} \left( n_{k_{a_i}} \times \sum_{j=1}^{k_{a_i}} m_{c_{a_{ij}}} \times p_{c_{a_{ij}}} \right) \end{aligned}$$

$$\begin{aligned} \text{類似度 } S_{b/a} &= \ln(I_{b/a}) = \sum_{i=1}^{m_b} (n_{k_{b_i}} \times p_{k_{b_i}}) \\ &= \sum_{i=1}^{m_b} \left( n_{k_{b_i}} \times \sum_{j=1}^{k_{b_i}} m_{c_{b_{ij}}} \times p_{c_{b_{ij}}} \right) \end{aligned}$$

$$\text{オリジナル度 } O_a = \ln \left( \prod_{i=1}^{t_a} p_i \right) = \sum_{i=1}^{t_a} \ln(p_i)$$

$$\text{オリジナル度 } O_b = \ln \left( \prod_{i=1}^{t_b} p_i \right) = \sum_{i=1}^{t_b} \ln(p_i)$$

諸元は次のとおりである。

使用される文字種の数	$N$
テキストAのテキスト長	$t_a$
テキストBのテキスト長	$t_b$
キーテキストの数	$m_a, m_b$
それぞれのキーテキストのテキスト長	

