

部分文字列増幅法による共通パターン発見アルゴリズム

池田, 大輔
九州大学情報基盤センター

山田, 泰寛
九州大学大学院システム情報科学府

廣川, 佐千男
九州大学情報基盤センター

<https://hdl.handle.net/2324/2969>

出版情報：情報処理学会研究報告：数理モデル化と問題解決. 2003 (122), pp. 45-48, 2003-12. 情報処理学会

バージョン：

権利関係：ここに掲載した著作物の利用に関する注意 本著作物の著作権は（社）情報処理学会に帰属します。本著作物は著作権者である情報処理学会の許可のもとに掲載するものです。ご利用に当たっては「著作権法」ならびに「情報処理学会倫理綱領」に従うことをお願いいたします。

部分文字列増幅法による 共通パターン発見アルゴリズム

池田 大輔、山田 泰寛*、廣川佐千男

九州大学情報基盤センター

*九州大学システム情報科学府

1本の文字列

例：HTML/XMLのファイル

2003.12.12

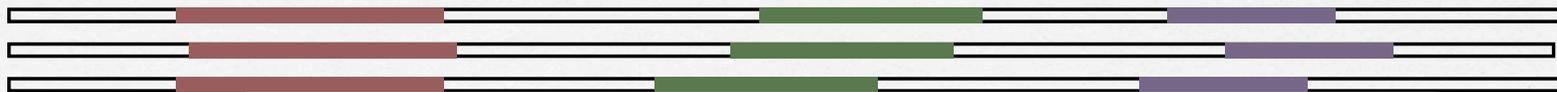
蛋白質配列



問題：共通な部分を探す

例：情報抽出のテンプレート発見

生物学的に意味のある部分配列



モデル化の例：最長共通部分列問題

2003.12.12

文字列の共通部分列

acbbacc**babaa**

ccbbcbabbabca

cbabbcbcab

cbbbab

同じ機能を持つ
部分配列?

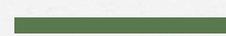
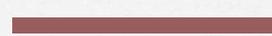
最長共通部分列問題
NP完全(Maier'78)

主な結果

- ☑ 入力：未知のパタンからベキ分布に従って生成された語の集合
- ☑ すべての語に共通な部分文字列の列を探す問題は平均的に $O(nm|t|^2 \log |t|)$ 時間で解ける

n は入力長で、ほぼ線形時間

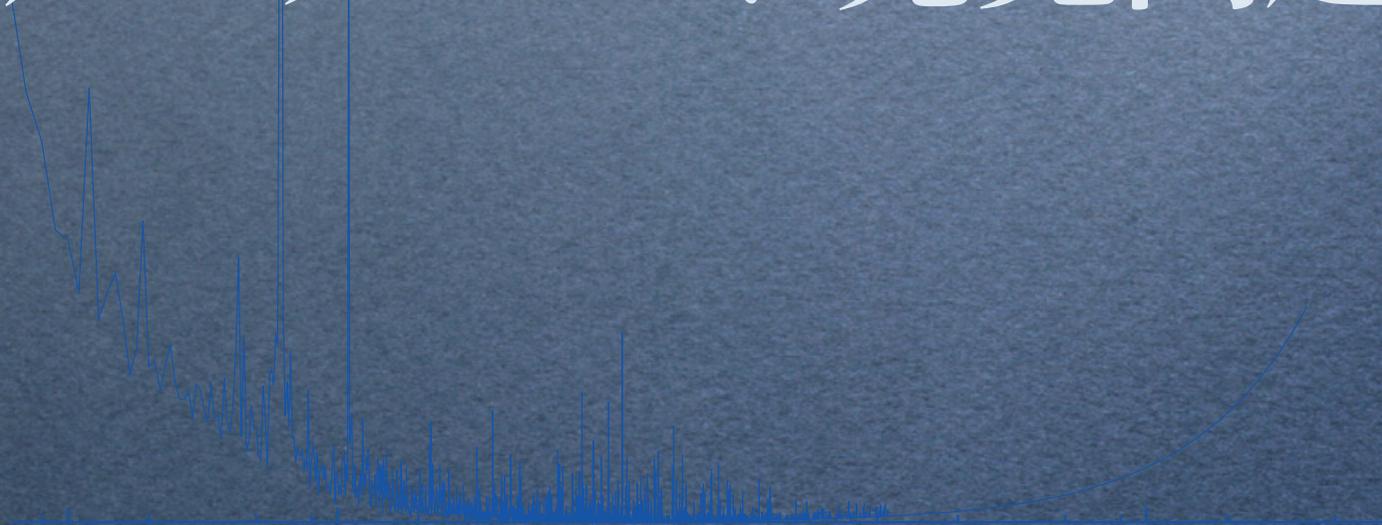
$m = 3$



目次

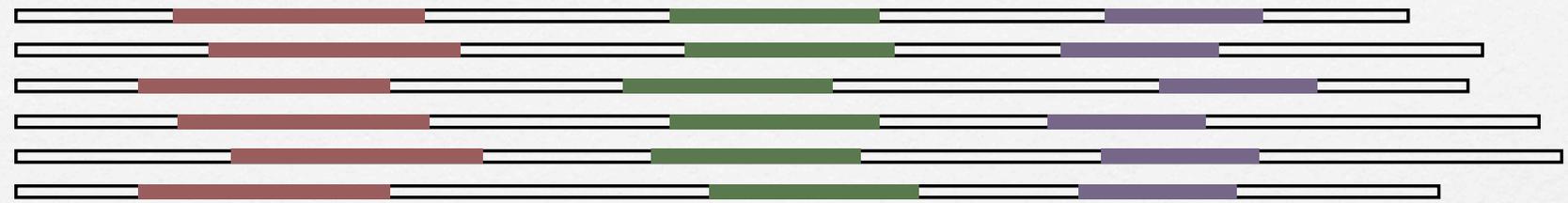
- テンプレート発見問題
 - 主に情報抽出の観点からの問題の見直し
 - 共通部分の表現方法
- 部分文字列増幅法
 - 共通部分の高速な特定方法
- 実験
 - Web上のデータからの情報抽出

テンプレート発見問題



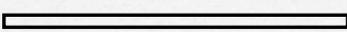


入力文字列の性質

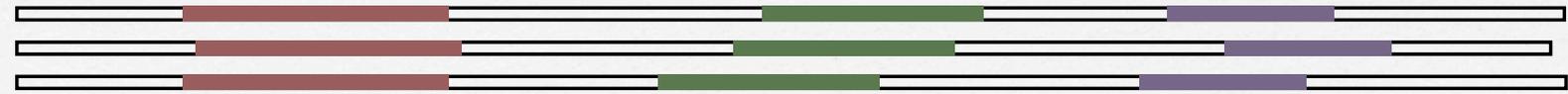


正例が与えられる

(, , )

 には自然な代入

テンプレート



パターン言語

Σ : 有限アルファベット

V : 変数の集合 $(\Sigma \cap V = \phi)$

パターン $\Leftrightarrow \Sigma \cup V$ 上の文字列

例 : $axbyax$ ($\Sigma = \{a, b\}, V = \{x, y\}$)

代入 $\Leftrightarrow \theta : V \rightarrow \Sigma^*$

パターン p の言語

$$L(p) = \{w \in \Sigma^* \mid \exists \theta; p\theta = w\}$$

テンプレート

正規パターン \Leftrightarrow 各変数の出現は高々1回

例： $axby$ $axbx$

正規パタンの言語族は正規言語族の真の部分クラス

$a.*b.*$

$w_0x_1w_1 \dots w_{m-1}x_mw_m$: 正規パターン
 $(w_i \in \Sigma^*, x_j \in V)$

テンプレート \Leftrightarrow 定数文字列の列 (w_0, \dots, w_m)

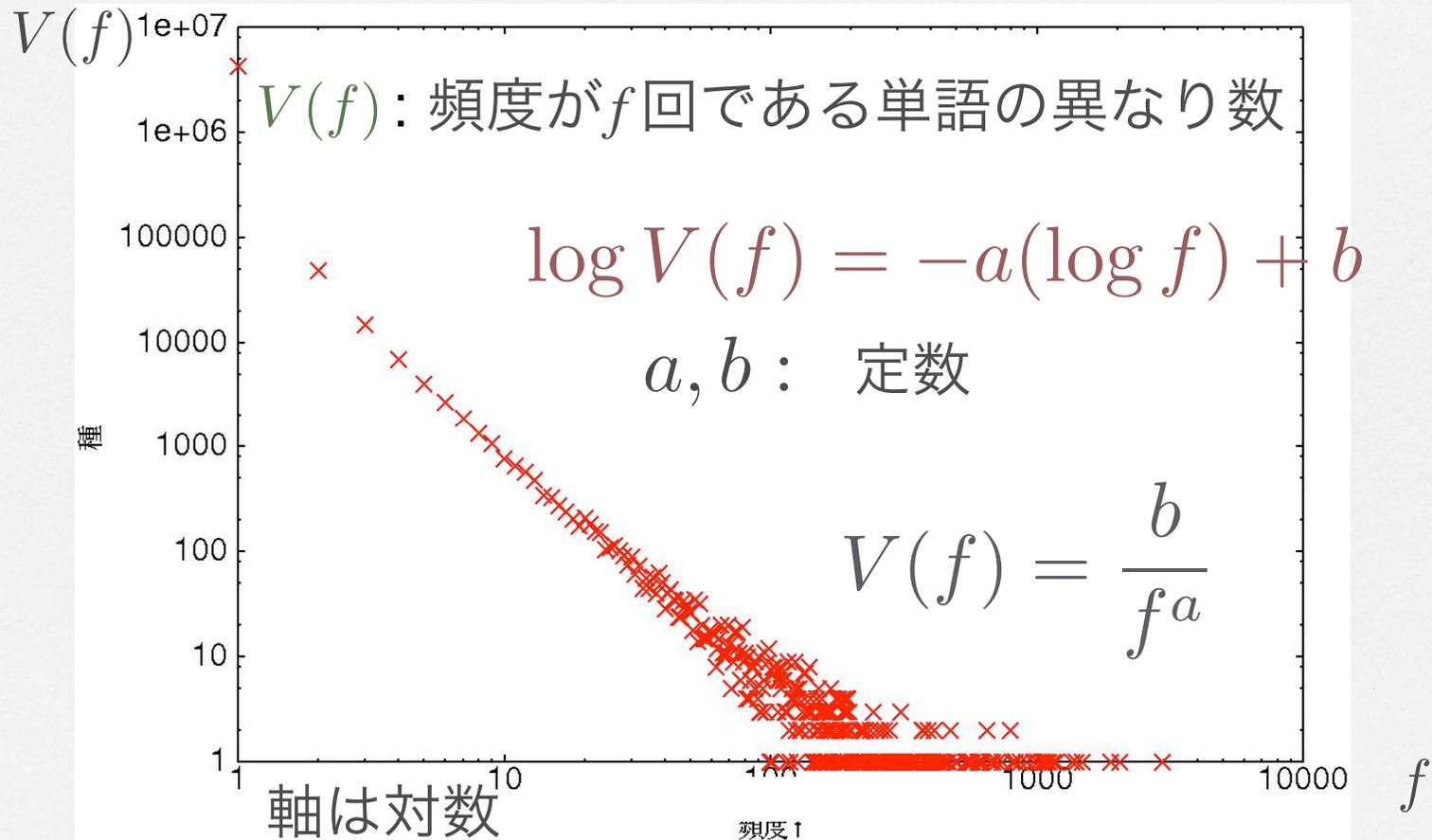
テンプレートと代入

仮定：代入される文字列の頻度は
ベキ分布に従う

$$p = x_1 A x_2 B x_3 C x_4$$

$$(A, B, C \in \Sigma^*, x_i \in V)$$

ベキ分布：Zipfの(第二)法則



テンプレート発見問題

2003.12.12

(池田'03)

入力：文字列の有限集合 S

S は未知のパタンからベキ分布に従って生成された

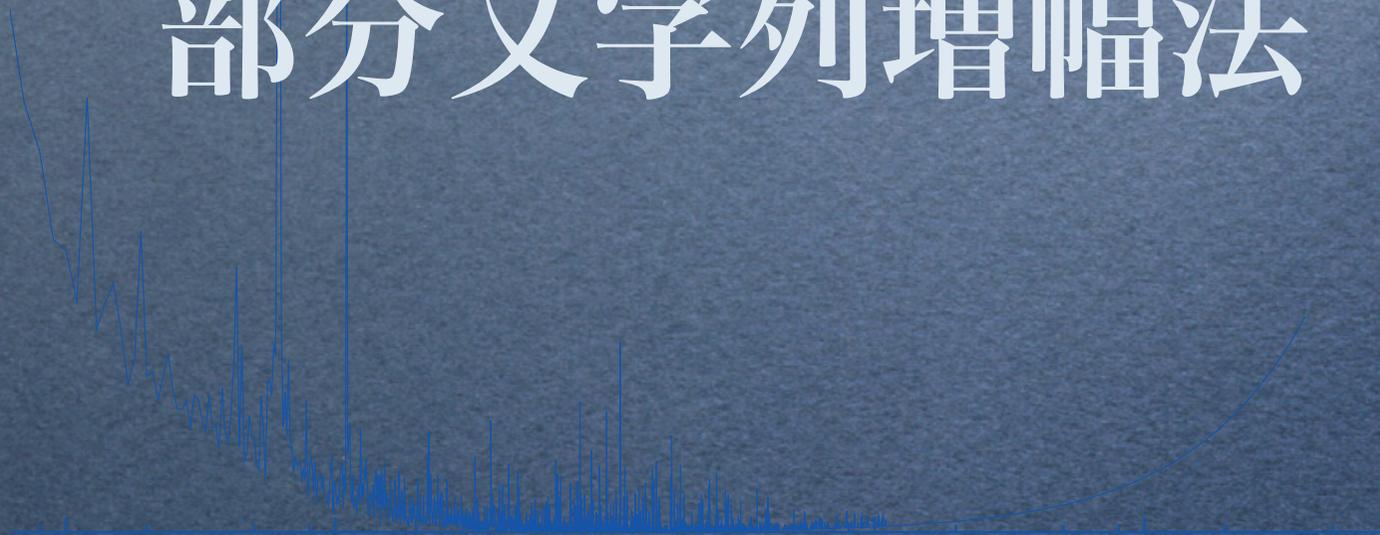
問題：以下を満すテンプレート t を見つけよ

t は正規パタン p のテンプレートである

$S \subseteq L(p)$ かつ $\nexists q; S \subseteq L(q) \subset L(p)$

p は S に対し極小

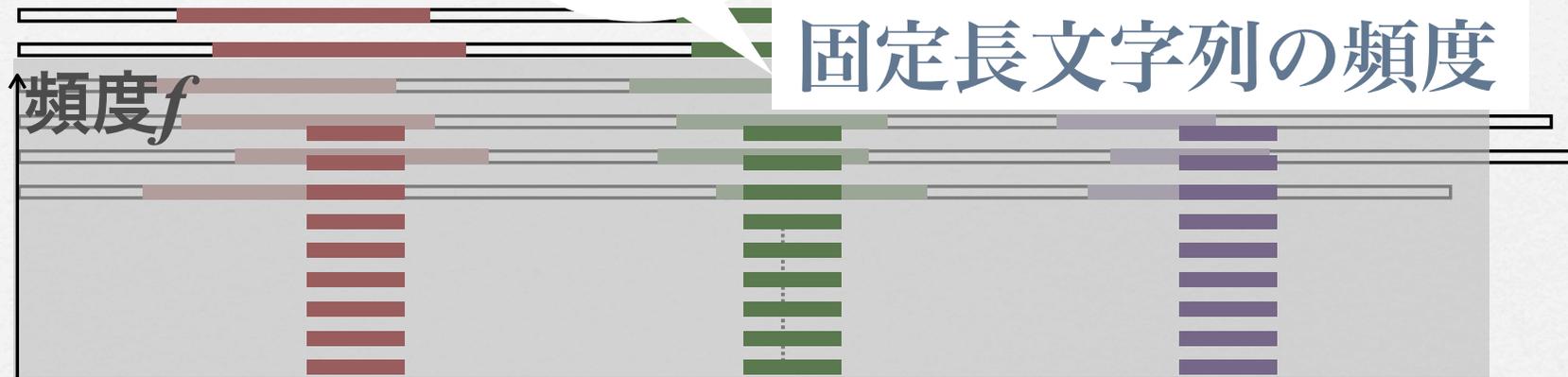
部分文字列增幅法



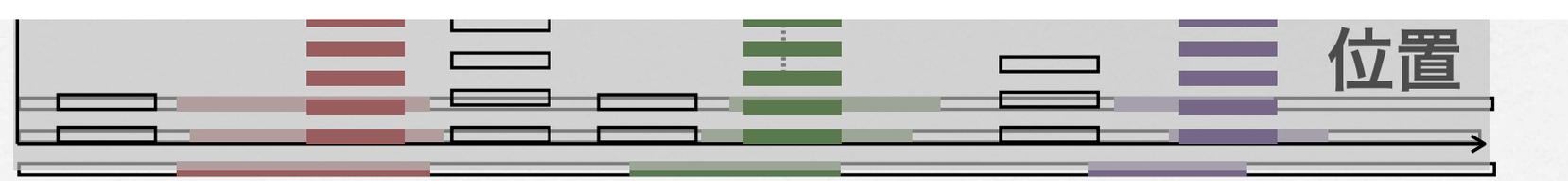
共通パターンと頻度

例えば長さ10

1ファイルでみる
固定長文字列の頻度



長さをどう決めるかが問題



文字列の総出現数

$F(f)$: 頻度(出現回数)が f 回である
ような文字列の総出現数

$$F(f) = fV(f)$$

$$\log V(f) = -a(\log f) + b$$

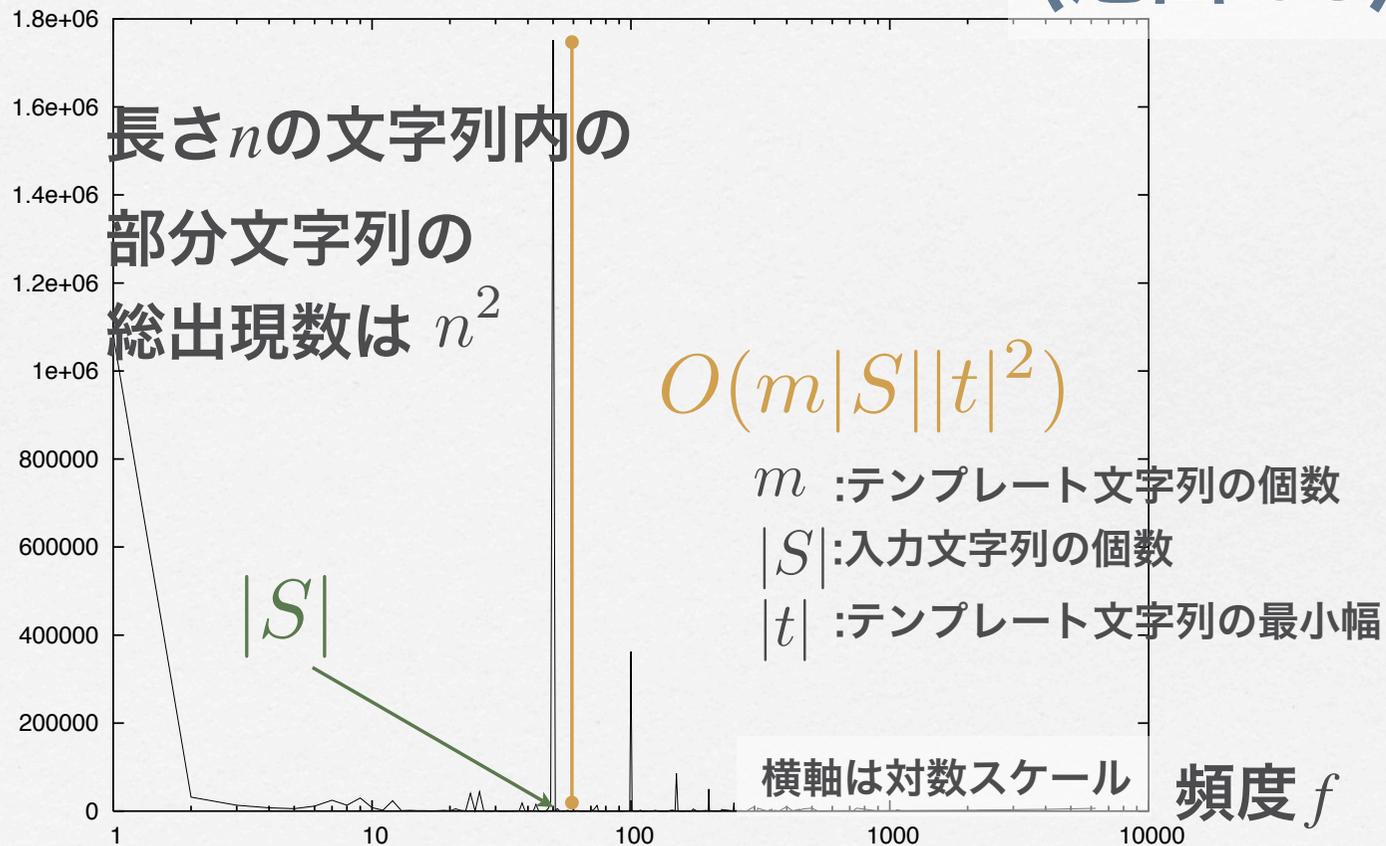
$$\log f + \log V(f) = (1 - a)(\log f) + b$$

$$\log fV(f) = (1 - a)\log f + b$$

$$\log F(f) = (1 - a)\log f + b$$



部分文字列増幅法 (池田'03)

 $F(f)$


abcabcabc の接尾辞木

abcabcabc\$

bcabcabc\$

cbcabc\$

bcabc\$

cabc\$

abc\$

bc\$

c\$

\$

\$

\$

\$

\$

\$

\$

\$

\$

\$

\$

\$

\$

\$

abcabcabc\$

abc\$

bcabc\$

bcabcabc\$

bc\$

cabc\$

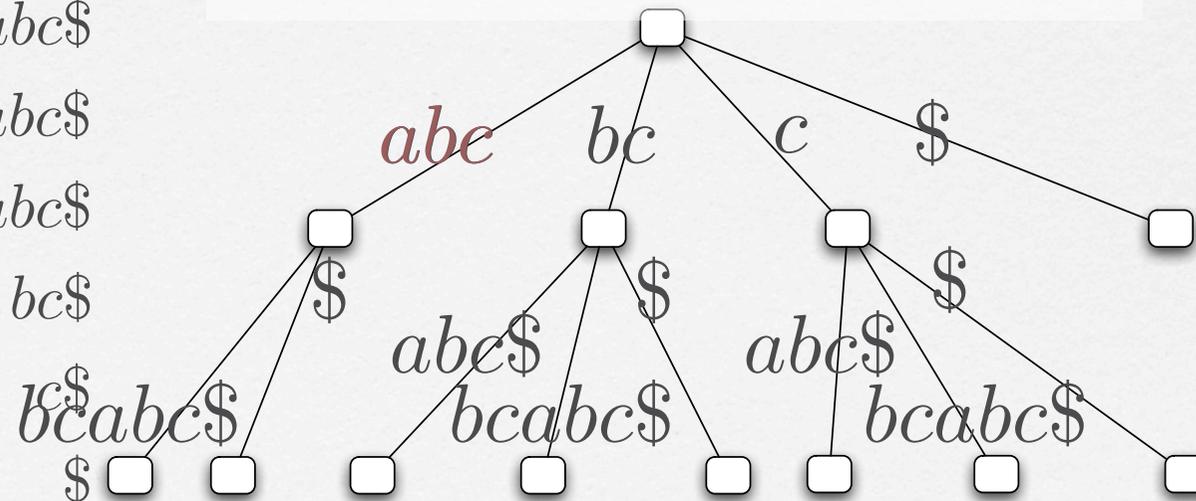
cbcabc\$

c\$

\$

branching word

⇔ 根からあるノードまでの
辺のラベルをつなげたもの



Lemma 1

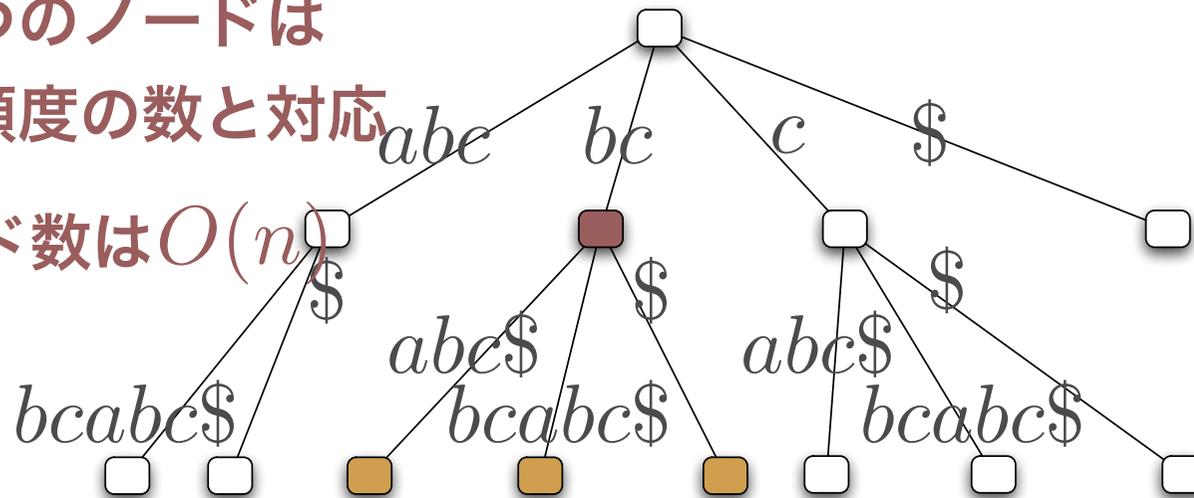
サンプル S 中に現われる部分文字列の異なる
頻度の数は高々 $O(n)$ 個である

ただし n は S 中の全て文字列の長さの和である

一つのノードは

ある頻度の数と対応

ノード数は $O(n)$



Lemma 2

テンプレート文字列内の全ての部分文字列に対し、
同じ文字列の総出現回数は高々 $O(m \log |t|)$
個である

m : テンプレート文字列の個数

$|t|$: テンプレート文字列の最小幅

```
function Template(S):template
begin
    V:=Count(S); #V(f)={f回出現するbranching wordの集合}
    for f in keys(V):
        F(f):=f | V(f) |;
    end
    P:=FindPeaks(F);
    t:=Reconstruct(V, P);
    return(t)
end
```

$$|P| \leq O(m \log t)$$

```
subroutine Reconstruct(S):template
```

```
begin
```

```
  for f in P:
```

```
    W := V(f); #V(f) = {f回出現するbranching wordの集合}
```

```
    for w in W:
```

sにおける出現回数と位置をチェック

```
    end
```

tに追加

```
  end
```

```
  return(t)
```

```
end
```

$$O(nm|t|^2 \log |t|)$$

Theorem

Templateは平均的に正しくテンプレート発見問題を解き、その計算量は $O(nm|t|^2 \log |t|)$ である

m : テンプレート文字列の個数

$|S|$: 入力文字列の個数

$|t|$: テンプレートの幅

○示すべきこと

$$S \subseteq L(p)$$

極小である

2003.12.12

極小性の証明

○示すべきこと

余計なものをテンプレートとしない

テンプレートの見落としがない

テンプレート文字列の
一部が**偶然**代入された



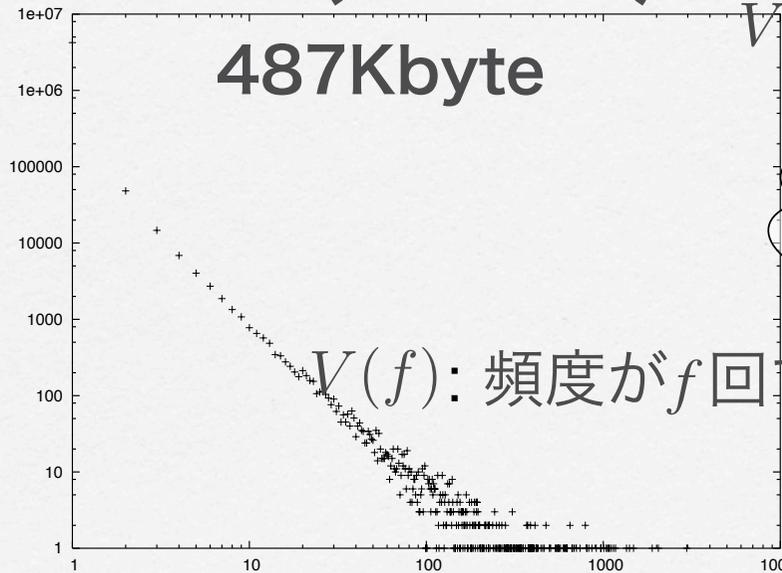
≡ の分だけピークが低くなる!

23

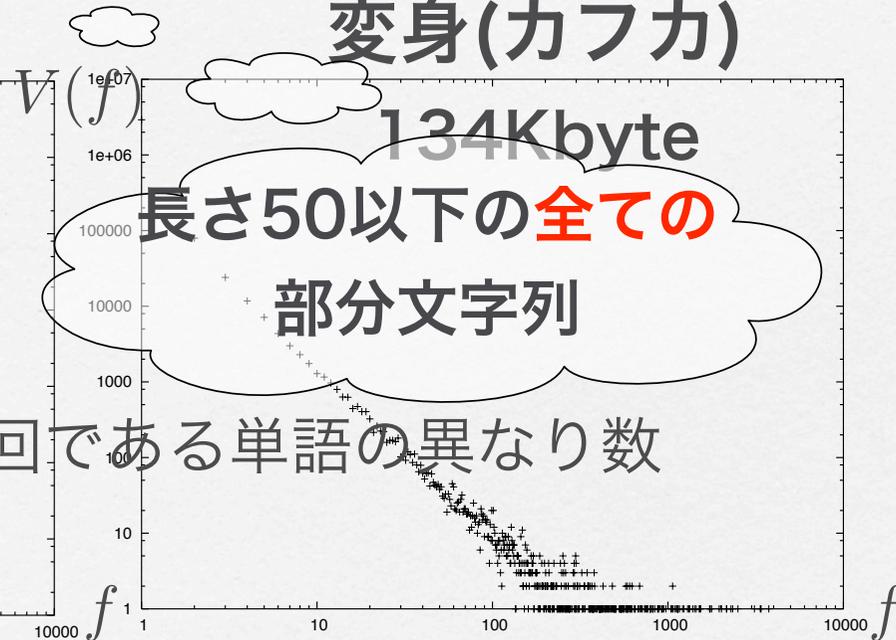
実験

小説中の部分文字列の分布

こころ(夏目漱石)



変身(カフカ)



青空文庫

<http://www.aozora.gr.jp/>

Project Gutenberg

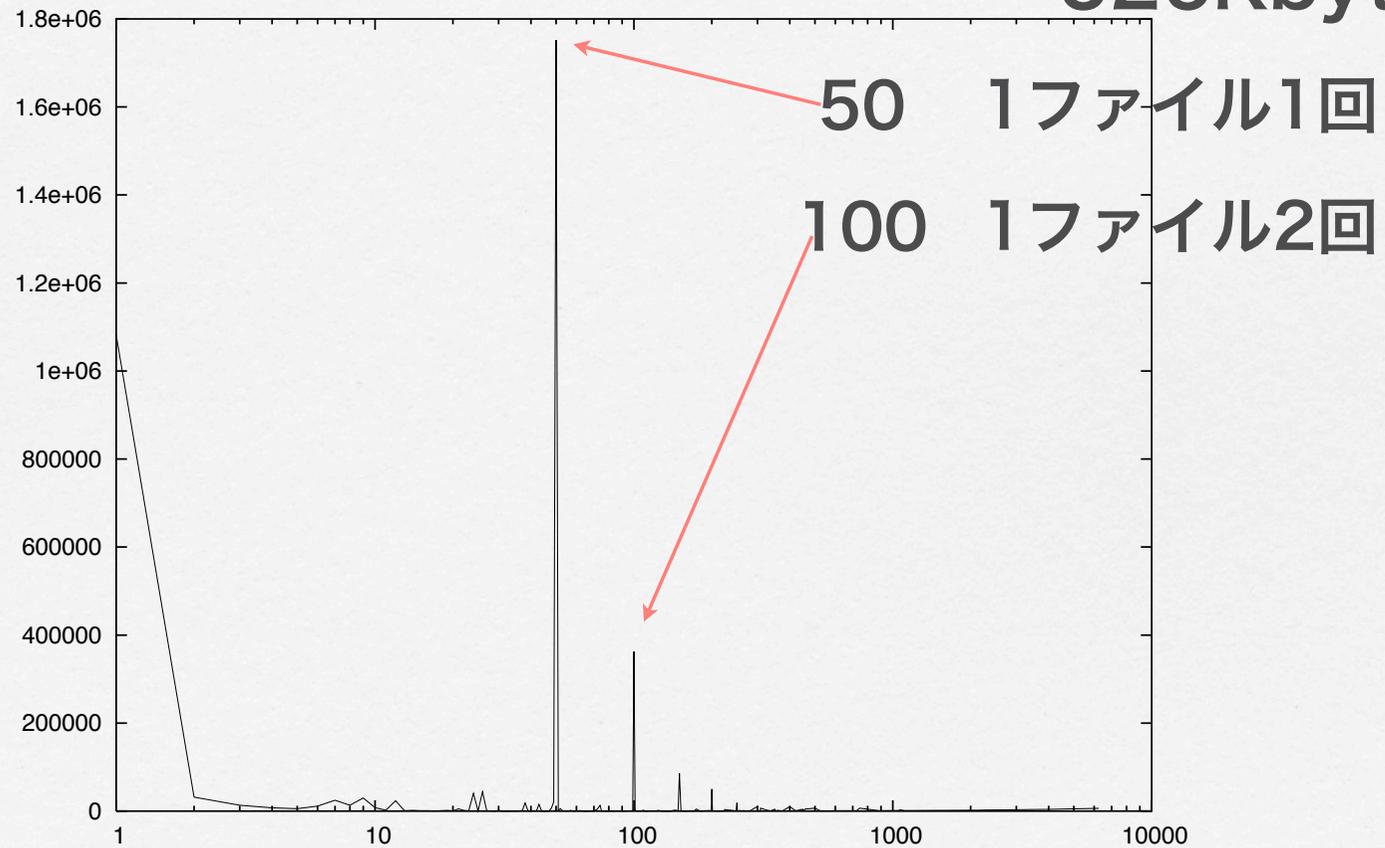
<http://www.promo.net/pg/>

実験データ

- データ：HTMLファイル
- 前処理はしない
 - 空白文字もタグもそのまま用いる
- $F(f)$ のグラフを作るところまで
 - `Reconstruct()`を実行しない

2003.12.12

産経新聞の記事ファイル50 526Kbyte



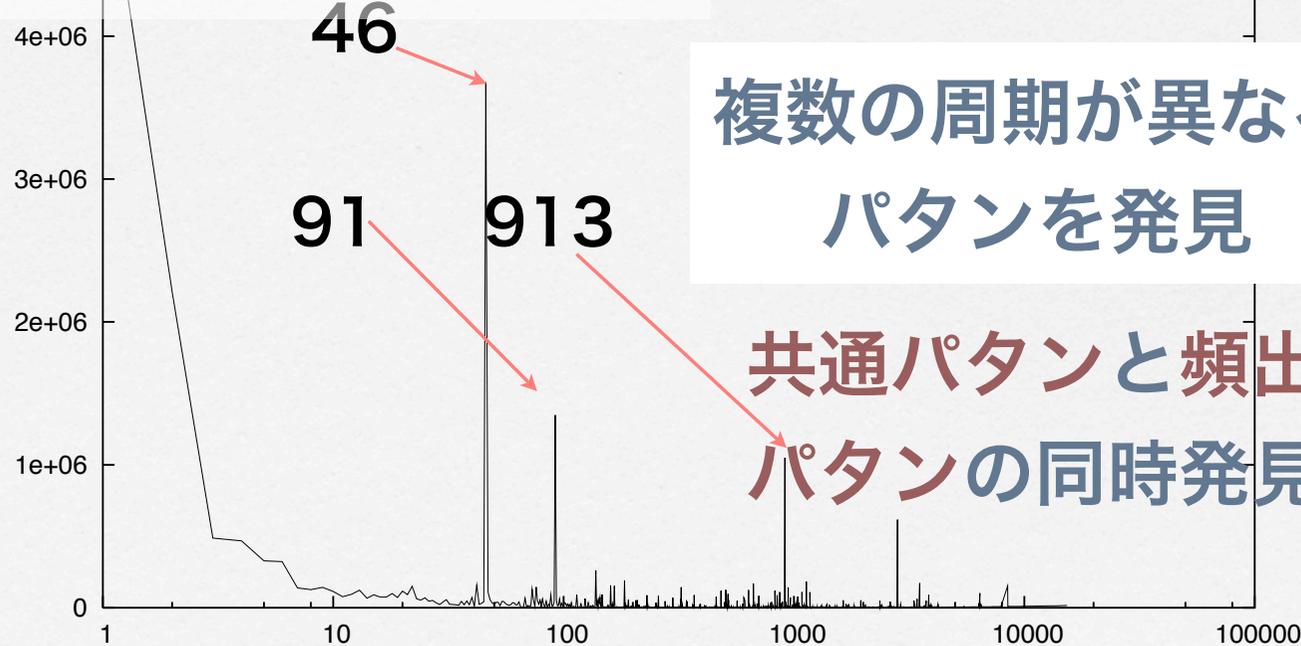
Yahoo!検索結果

- 46 : 各ファイルごと
- 91 : 一致したカテゴリ数
- 913 : 各検索結果ごと

46ファイル

(1212Kbyte)

913件

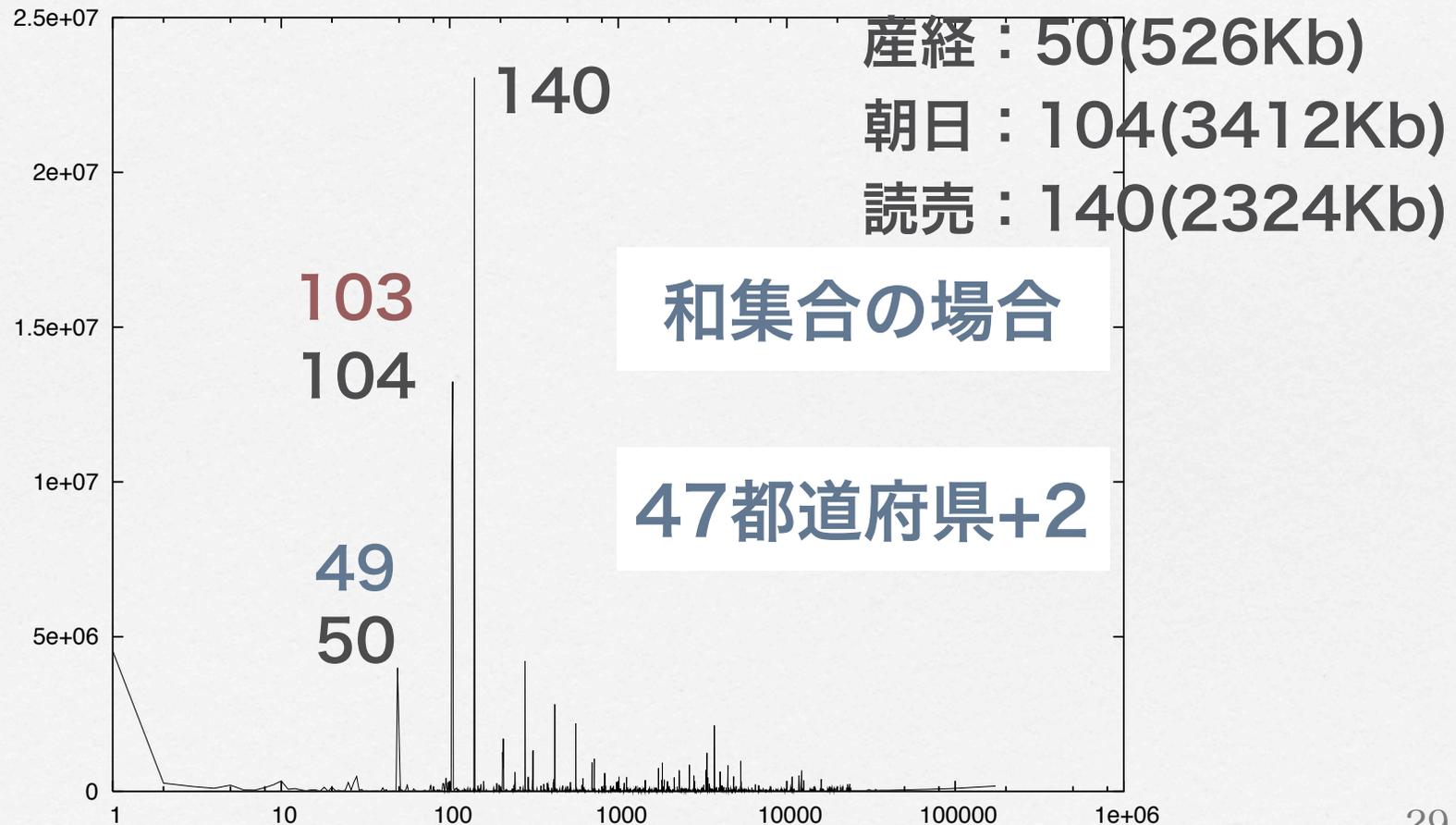


複数の周期が異なる
パターンを発見

共通パターンと頻出
パターンの同時発見

2003.12.12

3サイトのニュース記事



29

正規表現との関係

- 正則パターン言語族は正規言語族の真の部分クラス⇒表現力半構造化データに必須に劣る
- 1回以上の繰り返し $(abx)^+$
 - 同一ファイル内で出現回数が増える
 - ab は共通でかつ(ファイル内で)頻出
- 複数パタンの和 $(\alpha|\beta|\dots)$

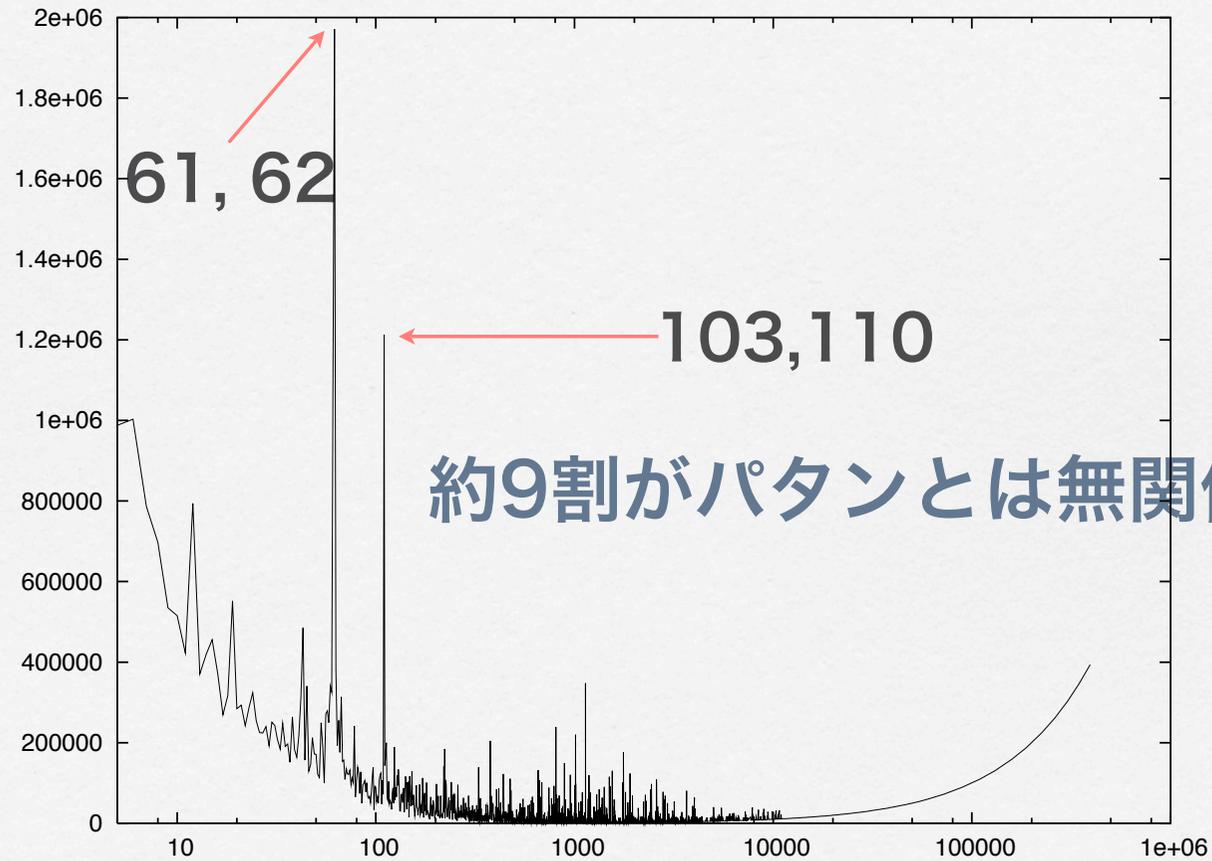
九大サイト内



九大トップページから深さ3まで
リンクをたどり598(5584Kb)ファイル

<http://www.kyushu-u.ac.jp/>

九大サイト内(Cont.)



九大サイト内(Cont.)

The screenshot shows a web browser window with the address bar displaying "九州大学__病院__". The page header includes the Kyushu University logo and name, along with navigation links for "Mail to Web master" and "サイトマップ". A search bar is present with the text "キーワードを入れてください" and a "検索" button. The main content area is titled "病院" (Hospital) and lists several affiliated hospitals: "九州大学の病院", "医学部附属病院", "歯学部附属病院", and "生体防御医学研究所附属病院". A sidebar on the left contains various navigation buttons such as "日本語", "English", "学部・大学院", "教育・学生生活", "研究", "社会との連携", and a list of services including "行事・イベント", "入試情報", "国際交流・留学", "教官情報", "図書館", "博物館", "病院", "同窓会・後援会", "改革・評価", "新キャンパス", "広報誌", "教職員向け情報", and "マップ". The footer contains the copyright notice "Copyright 2003 kyushu University. All rights reserved." and a page number "33".

まとめ

まとめ

- テンプレート発見問題の定式化
 - 正例とベキ分布
- アルゴリズムの提案と計算量
 - 頻度 f を持つ部分文字列による共通部分列の表現
 - $F(f)$ による部分文字列増幅法
 - 計算量 $O(nm|t|^2 \log |t|)$

まとめ(Cont.)

- 情報抽出の実験
 - HTMLファイルを利用
- 検証された結果
 - ベキ分布の検証
 - $F(f)$ による共通部分列の抽出
 - 繰り返しや和の演算にも対応可能
 - ノイズ耐性の高さ

2003.12.12

今後の課題

- 近似文字列への拡張
 - ○○%の違いを許容
- ゲノム、異常値検出等への応用
- 再現率と適合率による評価
- 誤り率の平均的解析による見積もり