

複数の順列に共通して現れるパターンの列挙法

井上 祐馬[†]

北海道大学 大学院情報科学研究科[†]

湊 真一[‡]

北海道大学 大学院情報科学研究科[‡]

1. はじめに

順列は1から n までの整数がちょうど1度ずつ出現する長さ n の数列である。順列にはパターン照合問題が存在し、テキスト順列 τ にパターン順列 π がマッチするとは、「 π と相対順序が一致するような τ の部分列が存在すること」と定義される。厳密には次節で定義する。順列パターンは、解が順列となる様々な問題の特徴付けに用いられ、例えば一時記憶としてスタックのみを使用するソート(スタックソート)で整列可能な入力、順列(2,3,1)にマッチしない順列に対応する[3]。

本稿では、複数のテキストすべてに共通してマッチするパターンを全列挙するアルゴリズムを提案する。提案手法は、順列集合を圧縮して効率的に処理するために、順列決定グラフ[2]というデータ構造を用いることで高速化を実現している。計算機実験により、ナイーブな手法に比べ10倍以上高速に動作することを確認した。

2. 順列パターン

テキスト順列 $\tau = (\tau_1, \dots, \tau_n)$ がパターン順列 $\pi = (\pi_1, \dots, \pi_m)$ を含むとは、あるインデックス $1 \leq i_1 < \dots < i_m \leq n$ が存在して、 $\tau_{i_x} < \tau_{i_y}$ のとき、かつそのときに限り $\pi_x < \pi_y$ であることをいう。例えばテキスト $\tau = (4, 2, 1, 3)$ 、パターン $\pi = (3, 1, 2)$ のとき、インデックス1,3,4に対する τ の部分列(4,1,3)がパターン(3,1,2)に対応するため、 τ は π を含む。

本稿では、複数のテキストに共通するパターンをすべて列挙する問題について考える。2つのテキスト順列の最長共通パターンを見つける問題がNP完全であることが知られている[1]など、1つの共通パターンを求めることでさえ難しい問題であり、全共通パターンの列挙に取り組んだ先行研究は筆者の知る限りない。

3. 順列決定グラフ

順列決定グラフは順列の集合を圧縮表現するデータ構造である。集合に含まれる各順列をある分解法で分解する手順を二分決定木で表し、それを圧縮することで順列

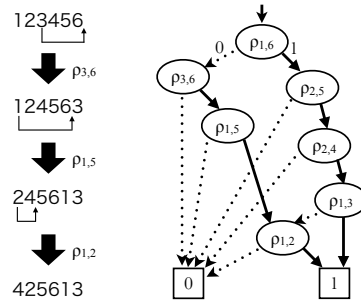


図1 Rot- π DD の例

決定グラフが得られる。本稿では、分解法として左ローテーションを採用した Rot- π DD [2] を用いる。ここで、区間 $[x, y]$ に対する左ローテーションは、 x 番目の要素を y 番目に移動し、 $x+1$ 番目から y 番目までの要素を左に1つ移動する操作であり、以降 $\rho_{x,y}$ と書く。

Rot- π DD の例を図1に示す。Rot- π DD の元となる決定木は、左ローテーション $\rho_{x,y}$ を表す内部節点と、 $\rho_{x,y}$ を行うか否かに対応する1-枝と0-枝を持つ。まず整列された順列 $(1, \dots, n)$ を持って、木を降り始める。木を降る際、1-枝を通るなら節点に書かれた $\rho_{x,y}$ を行い、0-枝を通るなら行わない。順列に対し左ローテーションを繰り返して終端節点に辿り着いたとき、対応する順列が集合に含まれるなら \perp 節点、含まれないなら \uparrow 節点を終端とする。この決定木に対して2つの圧縮ルール(1)1-枝が \uparrow 節点を指す節点を削除、(2)等価な節点($\rho_{x,y}$, 0-枝の先, 1-枝の先がそれぞれ同じ)をすべて1つに統合、を適用すると、圧縮された Rot- π DD が得られる。

Rot- π DD は2つの Rot- π DD 同士の演算が可能で、2つの順列集合の共通集合を求める演算や、片方の集合の各順列の並びに従ってもう片方の集合の各順列を並び替える演算などが存在する。演算の効率性は順列の個数ではなく Rot- π DD の節点数に依存している。

4. 提案法

提案法では1つの順列に含まれる全パターンを格納する Rot- π DD を構築する。それぞれの順列に含まれるパターンを表す Rot- π DD を構築すれば、それらの共通集合を Rot- π DD の演算で計算することで共通パターンが列挙できるためである。

Generation of Common Patterns in Permutations

[†] Yuma Inoue, Hokkaido University

[‡] Shin-ichi Minato, Hokkaido University

1つの順列に含まれるパターンを表す $\text{Rot-}\pi\text{DD}$ は、長さ k の全パターンの $\text{Rot-}\pi\text{DD}$ から長さ $k-1$ の全パターンの $\text{Rot-}\pi\text{DD}$ を構築する、という手順の繰り返しで構築する。入力順列 π の長さを n とすると、長さ n の唯一のパターン π のみを表す $\text{Rot-}\pi\text{DD}$ から始める。

長さ k のあるパターン π から長さ $k-1$ のパターンを作るには、 π の要素を1つ削除し、残った要素の順序を整理すればよい。これを長さ k の全パターンの全要素に対して行えば、重複を許して長さ $k-1$ の全パターンを列挙できる。長さ k のパターンから要素 i を削除するには、(1) i を k に、すべての i 以上の要素 j を $j-1$ に変更する、(2) k を一番右の k 番目に寄せて削除する、という手順を踏むことで可能になる。一見冗長に思えるが、この手順は $\text{Rot-}\pi\text{DD}$ の構造をうまく利用できる。

手順 (1) は、長さ $k-1$ の昇順の順列の i 番目に k を挿入した順列 $\sigma_{i,k}$ を考えることで簡単に実現できる。例えばパターン $(4, 2, 1, 3)$ の 2 を 4 に、3 と 4 をそれぞれ 2 と 3 に変更し、 $(3, 4, 1, 2)$ を作る場合を考える。これは、 $\sigma_{2,4} = (1, 4, 2, 3)$ を、元のパターン $\pi = (4, 2, 1, 3)$ に従って並べ替える、すなわち $\sigma_{2,4}$ の π_i 番目を i 番目に移動すればよい。例えば $\sigma_{2,4}$ の $\pi_1 = 4$ 番目は 3 なので、移動後の 1 番目が 3 になる。このようにすると、目的の $(3, 4, 1, 2)$ が得られる。 $\sigma_{i,k}$ は右ローテーションに対応するため、隣接の左ローテーションの繰り返し $\rho_{k-1,k} \dots \rho_{i,i+1}$ で表現できる。よって、長さ k のパターン集合を表す $\text{Rot-}\pi\text{DD}$ と $\{\sigma_{1,k}, \dots, \sigma_{k,k}\}$ を表す $\text{Rot-}\pi\text{DD}$ から、並べ替えの演算 1 回で手順 (1) が完了する。

手順 (2) では k を k 番目に寄せたいが、ここで $\text{Rot-}\pi\text{DD}$ の構造では最初 $(1, \dots, k)$ から降り始めていたことを思い出すと、最初から k は k 番目にあったのに、後の左ローテーションによって移動していることがわかる。よって、 $\rho_{x,k}$ を $\rho_{x,k-1}$ に変更することで、 k が k 番目から動かなくなり、結果的に k を k 番目に移動させたことになる。ただし、 $\rho_{x,k}$ の次の左ローテーション (=1-枝の先) が $\rho_{x',k-1}$ の場合、こちらも $\rho_{x',k-2}$ に変更する、というように、区間が k から連続する限り変更し続ける。

まとめると、順列 π に含まれる長さ k の全パターンを表す $\text{Rot-}\pi\text{DD}$ P_k の構築は、 $\{\pi\}$ のみを表す P_n から始め、(1) $P_k \leftarrow P_{k+1}$ と $\{\sigma_{1,k}, \dots, \sigma_{k,k}\}$ での並び替え演算結果の πDD 、(2) P_k の根から区間が k から連続する節点を修正、でできる。各手順はそのままでは効率的とは言えないが、 $\text{Rot-}\pi\text{DD}$ の演算の利用で高速化が期待できる。

5. 実験

提案法を実装し、ナイーブな手法 (テキストの全部分列を抜き出し、順序整理したものをパターンとして抽出) と比較した。C++ で実装し、3.20GHz CPU、64 GB メモリのマシンで実行した。入力は同じ長さ n のランダムな順列 2 つとし、2 つの順列の共通パターンを全列挙した。

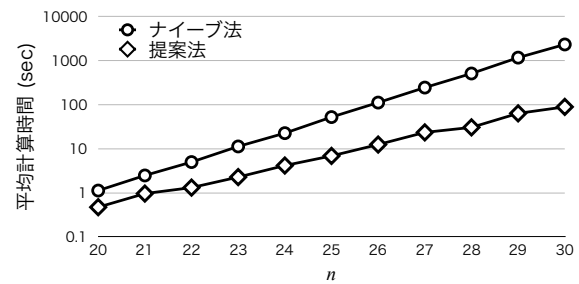


図2 各手法の実行時間の比較

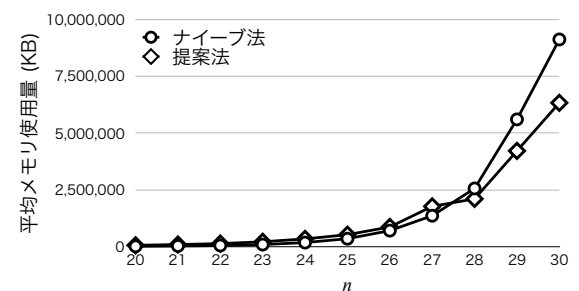


図3 各手法の使用メモリの比較

図2に実行時間の比較結果を示す。2つの順列の長さ n が大きくなるにつれ実行時間の差が広がることから、提案法が長い順列に対して相対的に効率が良いことがわかる。実験中最も長い $n=30$ については、ナイーブ法の約 20 倍高速に動作したことが確認できる。

図3に使用メモリの比較結果を示す。実行時間同様、順列が長いほど消費メモリに差がついた。 $n=30$ のケースでナイーブ法より 3GB ほどメモリ効率がよかった。

6. おわりに

本稿では、順列決定グラフ $\text{Rot-}\pi\text{DD}$ を用いた複数の順列に共通するパターンの全列挙手法を提案した。提案手法により、ランダムな順列に対しナイーブ法の 10 倍以上高速な共通パターン列挙が可能となった。

謝辞

本研究は JSPS 科研費 15J01665 および 15H05711 の助成を受けたものです。

参考文献

- [1] Mathilde Bouvel and Dominique Rossin: "The Longest Common Pattern Problem for two Permutations." *Pure Mathematics and Applications*, 17 (1-2), pp. 55-69, 2006.
- [2] Yuma Inoue and Shin-ichi Minato: "An Efficient Method for Indexing All Topological Orders of a Directed Graph." *In Proc. of 25th International Symposium on Algorithms and Computation (ISAAC)*, pp. 103-114, 2014.
- [3] Donald E. Knuth: "The Art of Computer Programming, Volume 1: Fundamental Algorithms." Addison-Wesley, 1968.