

# 講義「アルゴリズムとデータ構造」

## 第x回 その他の精選トピックス

大学院情報科学研究科 情報理工学専攻  
情報知識ネットワーク研究室  
喜田拓也

# 今日の内容

喜田'sセレクション:

文字列索引データ構造

接尾辞木

接尾辞配列

Burrows-Wheeler変換 (BWT)

※ 確率と情報の科学 高速文字列解析の世界  
ーデータ圧縮・全文検索・テキストマイニングー,  
岡之原大輔, 岩波書店, 2012年.

幾何的アルゴリズム

点・線・線分

符号付面積

線分の交差判定

凸多角形の判定

※ アルゴリズム・サイエンス シリーズ10  
計算幾何 ー理論の基礎から実装までー,  
浅野哲夫 著, 共立出版, 2007年.

# 文字列索引データ構造

パターン照合を行う際、テキスト側のデータをあらかじめ前処理することができるならば、「索引」を作っておけばいいじゃないか！

転置ファイル： キーワードの逆引き辞書をリストで持つ

接尾辞木： テキストのすべての接尾辞を木で持つ

接尾辞配列： 接尾辞の位置をソートした配列で持つ

FM-index： **BWT**文字列をウェーブレット木で持つ

ブロックソート圧縮法 (bzip2など)  
に用いられている

# 接尾辞木 (Suffix Tree)

テキストのすべての接尾辞を木で持つ索引データ構造

索引 = 元の文字列 + 接尾辞木  $ST(T)$

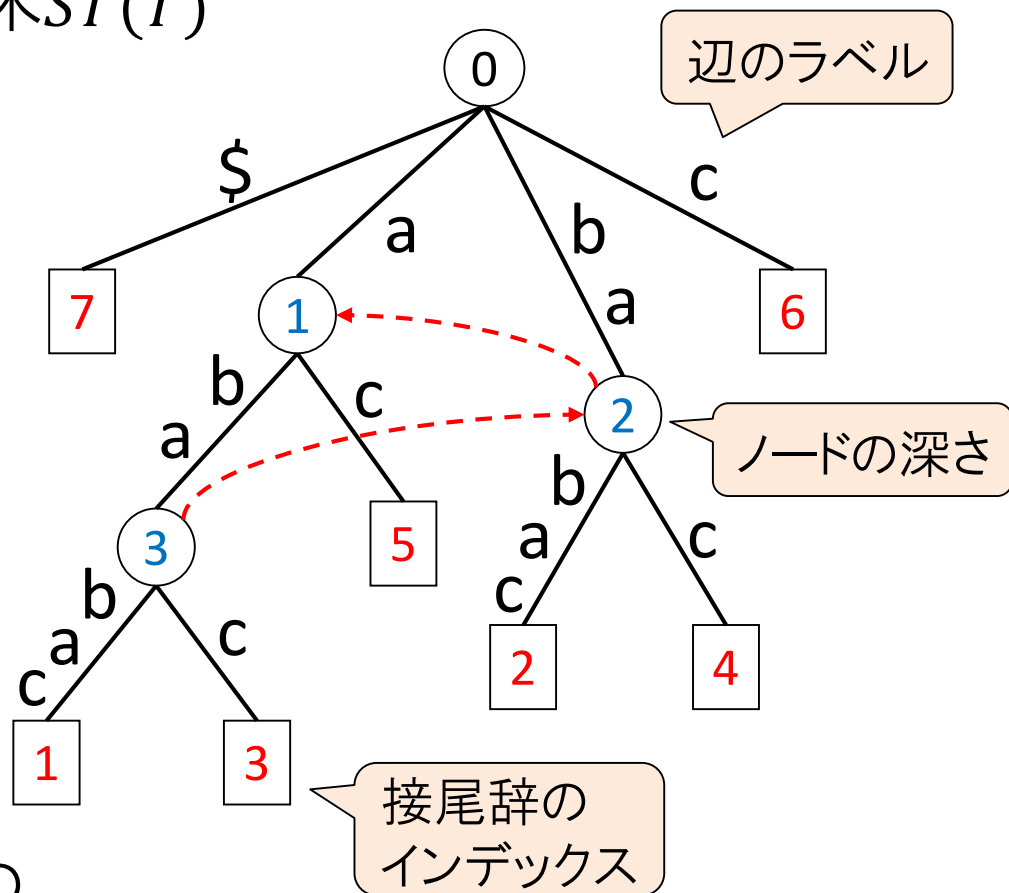
$O(n \log n)$  ビットの領域を使用

検索はパターン長  $m$  に対して

$O(m \log |\Sigma|)$  時間

1 2 3 4 5 6 7  
 $T = a b a b a c \$$

-----> Suffix link



辺のラベルは、 $T$ 上の文字列の  
(開始位置, 終了位置)として持つ

# 接尾辞配列 (Suffix Array)

テキストのすべての接尾辞を辞書順でソートし, その順番に接尾辞の開始位置(インデックス)を配列として保持したもの

索引 = 元の文字列 + 接尾辞配列  $SA_T[i]$

$n \log n + n \log |\Sigma|$  ビットの領域を使用

検索はパターン長  $m$  に対して  $O(m \log n)$  時間

1 2 3 4 5 6 7  
 $T = a b a b a c \$$

単純な構築法だと  $O(n \log n)$  時間かかるが,  $O(n)$  時間でできるうまい方法が知られている

$i$	$SA[i]$	
1	7	\$
2	1	a b a b a c \$
3	3	a b a c \$
4	5	a c \$
5	2	b a b a c \$
6	4	b a c \$
7	6	c \$

# Burrows-Wheeler変換 (BWT; BW変換)

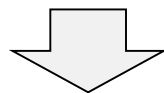
テキスト $T$ のソートした各接尾辞の  
1つ前の文字を並べたもの

$$BW_T[i] = T[SA[i] - 1]$$

テキスト $T$ の並び替えになっている  
 $BW_T$ から元の $T$ を復元できる

$BW_T$ は圧縮しやすい  
(同じ文字が連続しやすい)

1 2 3 4 5 6 7 8 9  
 $T = a c a g c a g g \$$



$BW_T = g \$ c c a g g a a$

$SA[i]$	$BW_T[i]$	
9	g	\$
1	\$	a c a g c a g g
3	c	a g c a g g \$
6	c	a g g \$
2	a	c a g c a g g \$
5	g	c a g g \$
8	g	g \$
4	a	g c a g g \$
7	a	g g \$

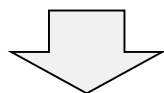
先頭の1つ前は末尾とする

# BW逆変換

$T$ は $BW_T$ から復元できる

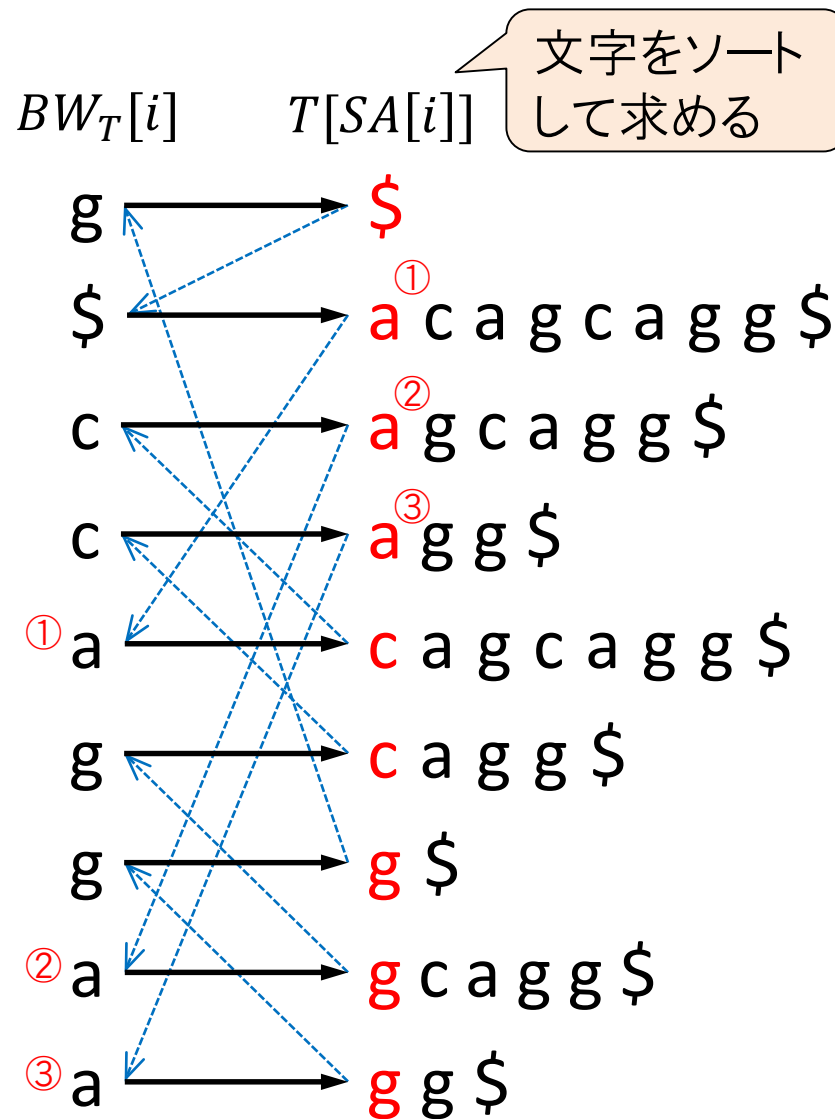
1.  $BW_T$ の文字の並びをソートして  $T[SA[i]]$ を得る
2.  $BW_T$ 上の\$のインデックス $k$ から始めて、対応する $T[SA[k]]$ を出力
3.  $T[SA[k]]$ の文字の $BW_T$ 上でのインデックスを求め、次の $k$ とする
4. \$を出力するまで2と3を繰り返す

$BW_T = g \$ c c a g g a a$



$T = a c a g c a g g \$$

同じようにして $SA[i]$ も復元できる！

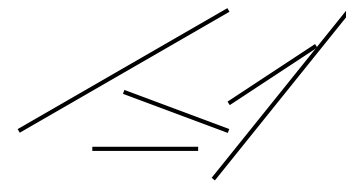


# 幾何的アルゴリズム

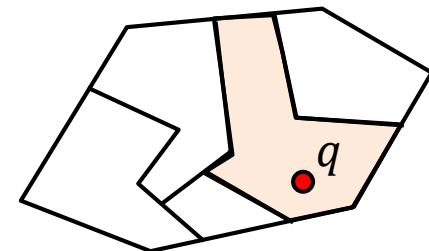
計算機で解きたい幾何問題はいっぱいある！

結構、難しい

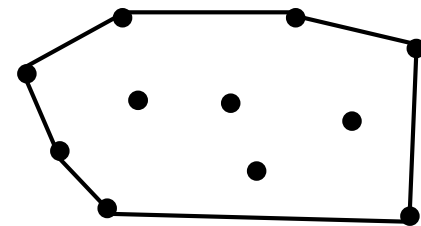
**線分の交差判定問題**: 「平面上に $n$ 本の線分が与えられたとき, これらの線分のあいだに交差があるかどうかを判定せよ」



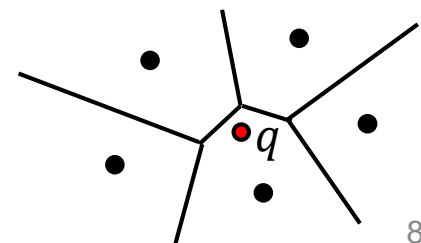
**点位置決定問題**: 「平面上に直線で描かれたサイズ $n$ の平面地図に対して, 質問点 $q$ を含む領域の名前を求めよ」



**凸包問題**: 「 $n$ 個の点を与えられたとき, これらすべての点を含む最小の凸多角形 (凸包, convex hull) を求めよ」



**最近点問題**: 「あらかじめ与えられた $n$ 個の点に対し, 任意に与えられた質問点にもっとも近い点を効率よく求めよ」





# 点・直線・線分

点 $p$ の表現: 平面の座標値 $p(x, y)$

点 $p_1(x_1, y_1)$ と点 $p_2(x_2, y_2)$ の距離: ユークリッド距離

$$\text{dist}(p_1, p_2) = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$

直線の表現: 直線の方程式  $ax + by + c = 0$

あるいは, 通る2点 $p_1(x_1, y_1)$ と $p_2(x_2, y_2)$ を指定した式を使う.

$$y = \frac{y_2 - y_1}{x_2 - x_1} (x - x_1) + y_1$$

ただし,  $x_1 \neq x_2$ である必要がある. よって一般には, 変形した式

$$(y_2 - y_1)x - (x_2 - x_1)y = x_1y_2 - x_2y_1$$

を使う. **線分**は2つの端点を指定する. 線分を含む直線の式は上のとおり.

# 符号付面積

3頂点 $p_1(x_1, y_1), p_2(x_2, y_2), p_3(x_3, y_3)$ からなる三角形の面積

$$\text{面積} S = \frac{1}{2} \left( (x_2 - x_1)(y_3 - y_1) - (x_3 - x_1)(y_2 - y_1) \right)$$

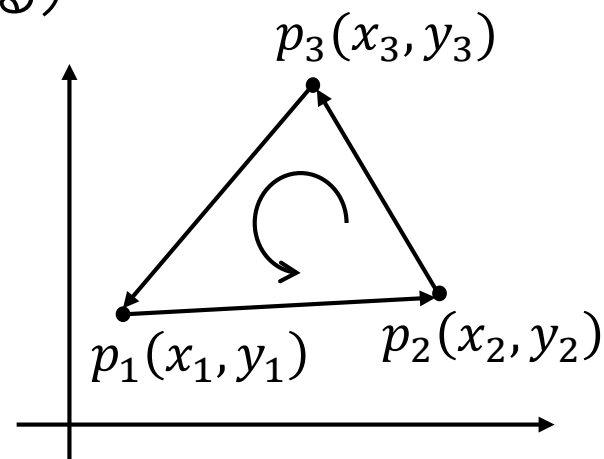
として求めることができる(通常は絶対値をとる)

実はこの式の符号には重要な意味がある！

$$\begin{aligned} \Delta(p_1, p_2, p_3) \\ &= (x_2 - x_1)(y_3 - y_1) - (x_3 - x_1)(y_2 - y_1) \end{aligned}$$

とすると,

$$\Delta(p_1, p_2, p_3) \begin{cases} > 0 : p_1, p_2, p_3 \text{ が反時計回りのとき} \\ = 0 : p_1, p_2, p_3 \text{ が一直線上にあるとき} \\ < 0 : p_1, p_2, p_3 \text{ が時計回りのとき} \end{cases}$$

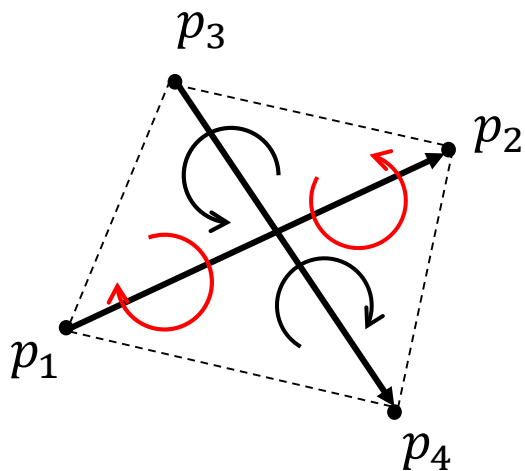


ベクトルの  
外積とほぼ  
等価な概念

# 線分の交差判定

2つの線分が与えられたとき、それらが互いに交差するかどうかを判定したい

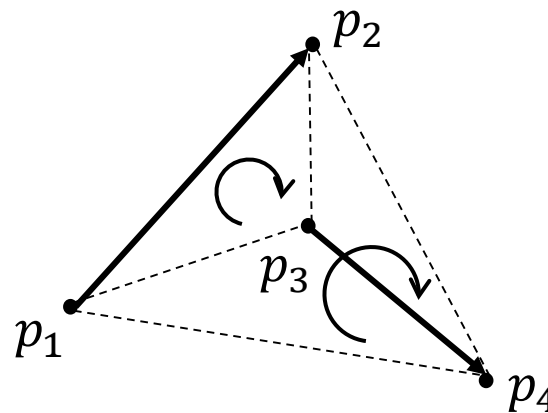
符号付面積をうまく使うと、簡単に判定できる！



$$\Delta(p_1, p_2, p_3) \times \Delta(p_1, p_2, p_4) < 0$$

かつ

$$\Delta(p_3, p_4, p_1) \times \Delta(p_3, p_4, p_2) < 0$$



$$\Delta(p_1, p_2, p_3) \times \Delta(p_1, p_2, p_4) > 0$$

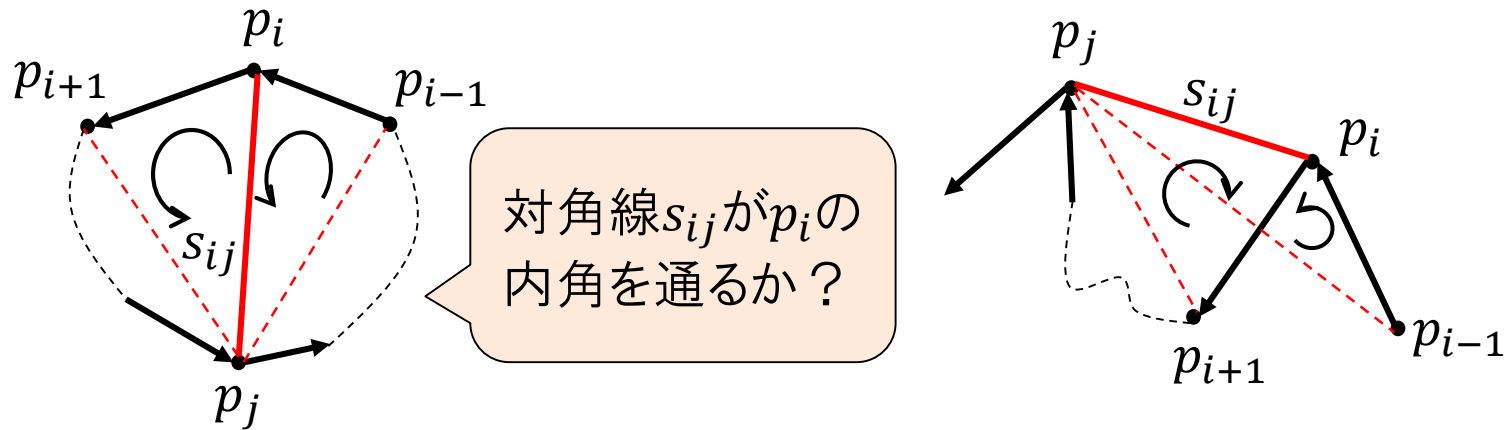
交差しない場合

# 凸多角形の判定

多角形の頂点列が反時計回りで与えられているとき, この多角形が凸多角形(すべての頂点が凸の多角形)かどうかを判定したい

## 定理

多角形 $P$ が凸多角形であるための必要十分条件は,  $P$ のすべての対角線が $P$ の内部にあるときである.



対角線の個数は頂点数が $n$ のとき,  $n(n-1)/2$ 個ある  $\Rightarrow O(n^2)$