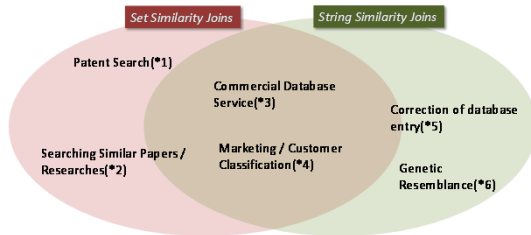
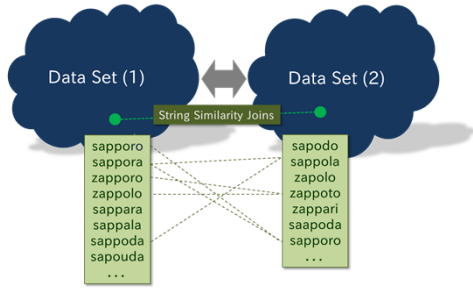


文字列類似結合におけるハイブリッド探索手法

JST-ERATO 湊離散構造処理系プロジェクト
白井康之, 高嶋宏之

文字列類似結合

一定の条件下(編集距離等)でマッチする類似ペアをすべて列挙する。



- (*1) Find similar patents in previous years, by other organizations
- (*2) Find similar research papers in previous years, in other countries, or in other fields
- (*3) Find similar records for cleaning and integration
- (*4) Find potentially good customers, fraud data (person / transaction)
- (*5) Find the misspelling of the person name, city name, ...
- (*6) Find similar parts of gene sequences

フィルタリング手法

- 分割・枝刈り
- q-gram インデクシング
- prefix / suffix によるフィルタリング

データの構造化に基づく手法

- Trie構造を用いた手法 (Trie-Join)
- ZDDを用いた手法 (ZDD-Join)

	フィルタリング手法	Trie-Join (ZDD-Join)
レコード長	長い(ほうが得意)	短い(ほうが得意)
アイテム数	多い	少ない
制約(編集距離)	大きい	小さい

問題定義(文字列類似結合)

Input

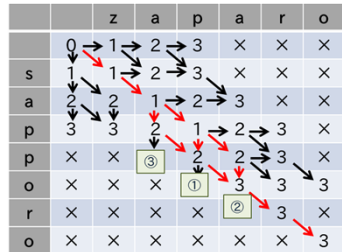
- シーケンスデータ集合: R, S
ex.) $R = \{ "abc", "bc", "def" \}, S = \{ "ab", "cd", "efg" \}$
- $r \in R, s \in S$ に対する類似尺度: $\text{sim}(r,s)$ (編集距離)
ex.) $r_1 = "abc", s_1 = "adc", \text{sim}(r_1, s_1) = 1$
 $r_2 = "abc", s_2 = "acd", \text{sim}(r_2, s_2) = 2$
- 閾値: N (編集距離の最大値)

Output

- R の部分集合 R' , s.t., $\forall r \in R', \exists s \text{ sim}(r,s) \leq N$

動的計画法によるストリングマッチング

sapporo
zaparo ($\delta = 3$)

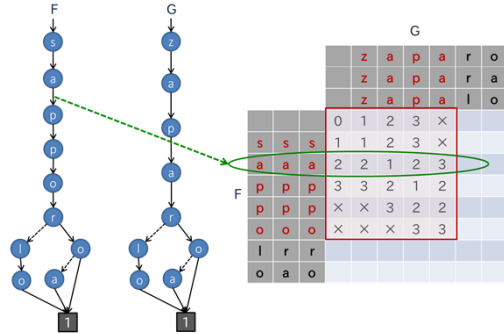


- ① sapporo
zaparo
- ② sapporo
zaparo
- ③ sapporo
zaparo

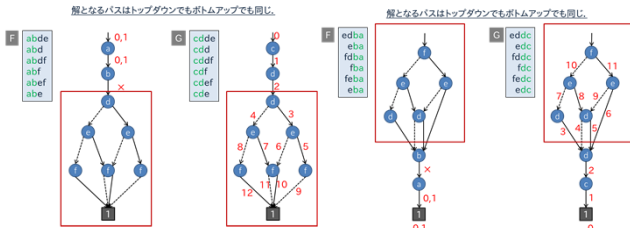
(赤字は replace, 青字は insert/delete)

ZDD (Seq-BDD) によるストリング類似結合

- Sequence BDD [Loekito 2009, Denzumi 2011, ほか] を用いて, ストリング集合間の類似結合を行う。
- f の各ステップでは, 動的計画法に基づき, 対応する g のパスとの編集距離を更新する。



ハイブリッド探索

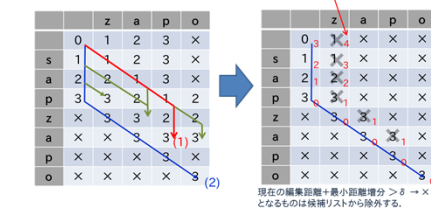


トップダウンの方が効率的な例 ($\delta = 1$)

ボトムアップの方が効率的な例 ($\delta = 1$)

Length Pruning

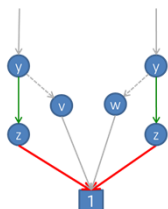
- ex. 編集距離 = 3, {zap, sapzapo}
- 長さが大きく異なるデータセットでは, 有効



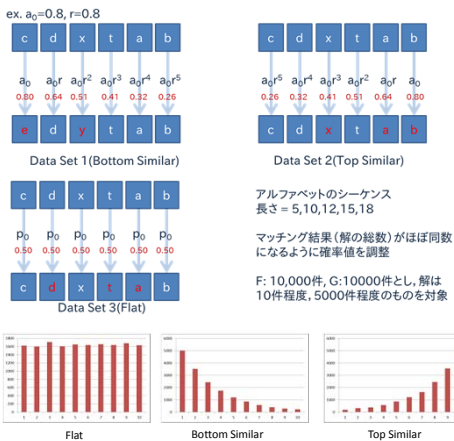
現在の編集距離 + 最小距離増分 $> \delta \rightarrow \times$ となるものは候補リストから除外する。

冗長計算の削除

キャッシュを利用して, 冗長計算を排除。
 g でパスは異なるが, アイテム集合は同じもの



実験データ



実験結果

10000件×10000件(解が少ないケース)

	編集距離	#Results	F側ノード探索数
random	top-down	1	159
	bottom-up	1	159
	hybrid	1	159
Bottom Similar	top-down	1	183
	bottom-up	1	183
	hybrid	1	183
Top Similar	top-down	1	206
	bottom-up	1	206
	hybrid	1	206

10000件×10000件(解が比較的多いケース)

	編集距離	#Results	F側ノード探索数
random	top-down	1	4971
	bottom-up	1	4971
	hybrid	1	4971
Bottom Similar	top-down	1	4863
	bottom-up	1	4863
	hybrid	1	4863
Top Similar	top-down	1	4886
	bottom-up	1	4886
	hybrid	1	4886

Length Pruning の効果

Length Pruning	編集距離	L=10-20		L=12-18		L=15	
		#Results	F側ノード探索数	#Results	F側ノード探索数	#Results	F側ノード探索数
無	1	3	305048	3	308074	3	311024
有	1	3	229610	3	243563	3	288917

今後の予定

- ボトムアップ手法とトップダウン手法の効率的な融合方法の検討(問題の類型化と制御方法の検討)
- 高速化