

ПРИБЛИЗИТЕЛЬНАЯ ОЦЕНКА ПОДОБИЯ ДЕРЕВЬЕВ

Белорусский государственный университет информатики и радиоэлектроники
г. Минск, Республика Беларусь

Лычковский А. В.

Волорова Н. А. – к-т. техн. наук, доцент

Сравнение данных является распространенной операцией в вычислительных системах. Большие объемы информации требуют эффективных алгоритмов их сравнения. Зачастую сравниваемые данные берутся из разных источников, в которых они представлены в разных форматах. Простое сопоставление таких данных не даст приемлемого результата, т.к. их структура различается. Однако, в большинстве случаев, информация о структуре также может быть использована в процессе сравнения.

Реальным примером является задача слияния баз данных с адресами. Эта задача не может быть решена простым сопоставлением строк, т.к. имена улиц могут отличаться в зависимости от соглашений принятых при их сохранении. Более того, имена улиц могут быть записаны на другом языке. Для решения этих проблем можно использовать иерархическую структуру адреса. На рисунке 1 представлены таблицы с именами улиц из разных источников, а также таблицы с адресами на этих улицах. Представление в виде дерева записей 30 и 91 из таблиц показано на рисунке 2. Наличие меры, которая оценивает схожесть деревьев, позволило бы решить поставленную задачу слияния баз данных с адресами.

SRD		SLR		RO					LR				
id	street	id	street	id	num	entr	apt	resident	id	num	entr	apt	owner
30	Giuseppe-Cesare-Abba-Str.	91	CESARE ABBA STRASSE	30	1	-	1	Pichler	91	1	-	1	Maier
120	Sebastian-Altman-Str.	74	S. ALTMANN STRASSE	30	1	-	3	Rieder	91	1	-	2	Rossi
5220	Bozner-Boden-Str.	33	BOZNER BODENWEG	30	2	A	-	Maier	91	1	-	3	Sparber
3000	Hermann-von-Gilm-Str.	109	GILMWEG	30	2	B	1	Rossi	91	2	A	-	Maier
3030	Pater-Reginaldo-Giuliani-Str.	185	P. R. GIULIANI STR.	30	2	B	2	Woelk	91	2	B	1	Totti
3540	Italienallee	115	ITALIENSTRASSE	30	2	B	3	Verdi	91	2	B	2	Bracco
4440	Musterplatz	165	MUSTERPLATZ	30	2	B	4	Verdi	91	2	B	3	Mair
7180	Raffaello-Sernesi-Galerie	207	SERNESIDURCHGANG	30	2	C	-	Burger	91	2	B	4	Lun
7590	Telsergalerie	259	TELSERDURCHGANG	30	3	-	-	Hofer	91	2	D	-	Tribus
7620	Friedensplatz	139	SIEGESPLATZ	30	4	A	1	Tribus	91	3	-	-	Costanzi
7650	Turiner Str.	266	TURINER STRASSE	30	4	A	2	Palermo	91	4	A	-	Palermo
7740	Trienter Str.	262	TRIENTER STRASSE	30	4	A	3	Palermo	91	4	B	-	Abel
7860	Triester Str.	263	TRIESTER STRASSE	30	4	B	-	Abel	91	4	C	-	Rossi
8580	Walther-v.-d.-Vogelweide-Pl.	285	WALTHERPLATZ	30	4	C	-	Rossi	91	6	-	-	Spiro
3930	Giannantonio-Manci-Str.	86	MANCISTRASSE	30	6	-	-	Spiro	91	6	-	-	Spiro
...		...		120	3	A	1	Spiro	74	3	A	1	Spiro
				120	3	A	2	Barducci	74	3	A	2	Barducci
				120	3	A	3	Costanzi	74	3	A	3	Costanzi
				120	3	A	4	Pichler	74	3	A	4	Spiro
				120	3	A	5	Spiro	74	3	A	6	Hofer
				120	3	A	6	Raifer	74	4	-	-	Mueller
							

Рисунок 1 – Исходные данные в разных источниках

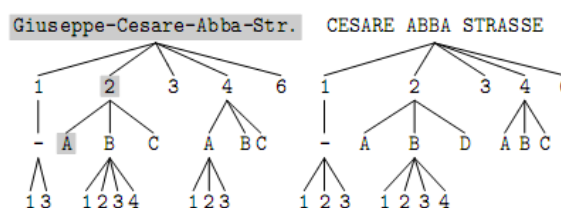


Рисунок 2 – Деревья с адресами для записей 30 и 91

Одна из хорошо известных мер – редакционное расстояние, которое определяется как минимальная по стоимости последовательность операций вставки узла, удаления узла, переименования метки, которая преобразует одно дерево в другое [1]. Существует алгоритм [2], вычисляющий редакционное расстояние за $O(n^2 \min^2(l, d))$ времени и $O(n^2)$ памяти для деревьев с n узлами, l листьями и глубиной d . Более ранние алгоритмы имеют сложность $O(n^2)$ и не могут быть использованы для больших деревьев.

Для решения некоторых прикладных задач, например такой, как описана выше, требуется лишь примерная оценка схожести деревьев. Rq-грамм расстояние позволяет сделать такую оценку за время $O(n \log(n))$ и требует $O(n)$ памяти.

Ниже дается определение rq-грамм и мере основанной на этом понятии. Если сформулировать коротко, то rq-граммы дерева – это все поддеревья определенной формы. Для того чтобы быть уверенным, что каждый узел дерева присутствует хотя бы в одной rq-грамме, начальное дерево расширяется нулевыми узлами. Исходя из этого, rq-граммы определяются, как поддеревья расширенного дерева.

Определение 1(расширенное дерево) Пусть T дерево, а $p > 0, q > 0$ два целых числа. Расширенное дерево T^{pq} строится из T добавлением $p - 1$ предка к корневому узлу, вставкой $q - 1$ ребенка перед

первым и после последнего ребенка каждого не листового узла, а также добавляем q детей к каждому листу T . Все новые узлы, должны быть нулевыми, такими, что не встречаются в T .

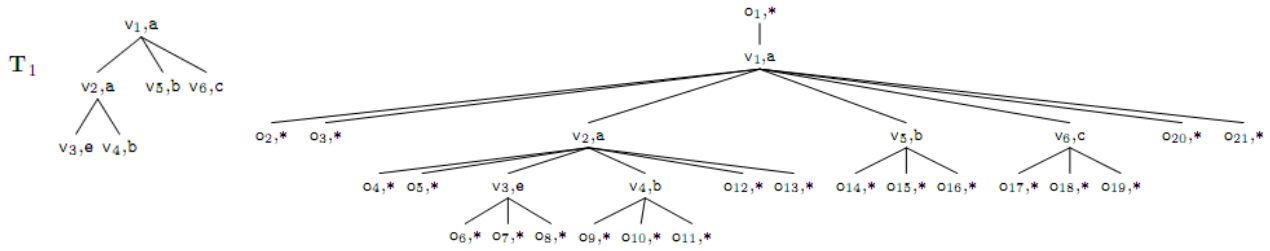


Рисунок 3 – Дерево T_1 и его расширенная версия $T_1^{2,3}$

Определение 2(rq-грамм шаблон) Для $p > 0, q > 0$, rq-грамм шаблон – дерево из основного узла с $p - 1$ предками и q детьми.

Определение 3(rq-грамм) Для $p > 0, q > 0$, rq-грамм дерева T - поддереву расширенного дерева T^{pq} , которое изоморфно rq-грамм шаблону.

Определение 4(кортеж меток) Пусть G – rq-грамма с узлами $V(G) = \{v_1, \dots, v_{p+q}\}$, где $v_i - i$ узел в начальной сортировке. Кортеж $l(G) = (l(v_1), \dots, l(v_{p+q}))$ называется кортежем меток дерева G .

Определение 5(rq-грамм профиль) Для $p > 0, q > 0$, rq-грамм профиль, $PP^q(T)$, дерева T – набор кортежей меток $l(G_i)$ всех rq-грамм G_i дерева T .

Определение 6(rq-грамм расстояние) Для $p > 0, q > 0$, rq-грамм расстояние, $\Delta^{p,q}(T_1, T_2)$, для деревьев T_1, T_2 вычисляется по следующей формуле:

$$\Delta^{p,q}(T_1, T_2) = 1 - 2 \frac{|PP^q(T_1) \cap PP^q(T_2)|}{|PP^q(T_1) \cup PP^q(T_2)|}$$

Rq-грамм расстояние равно единице, если два дерева не имеют одинаковых rq-грамм. Деревья с расстоянием равным нулю имеют одинаковые rq-грамм профили. Следует отметить, что нулевое расстояние не подразумевает эквивалентности деревьев. Rq-грамм расстояние может быть вычислено за $O(n \log(n))$ времени, нужное для вычисления количества общих rq-грамм в профилях размера $O(n)$. Существует оценка, которая показывает, что размер rq-грамм профиля равен $|PP^q(T)| = 2l + qi - 1$, где l – количество листьев, а q – количество узлов, не являющихся листьями [3].

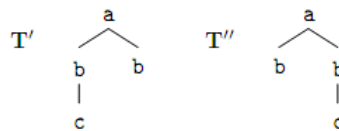


Рисунок 4 – Неравные деревья, имеющие одинаковый rq-грамм профиль

Рассмотренная мера может быть применена для приближенной оценки сходства деревьев. Такая оценка может понадобиться в задачах, пример которой приведен в начале статьи. Также эта оценка, благодаря скорости своей работы, может использоваться на предварительном этапе сравнения деревьев, при нахождении редакционного расстояния.

Список использованных источников:

1. K.-C.Tai. Tree to tree correction problem. Journal of the ACM (JACM), 1979
2. K.Zhang, D.Shasha. Simple fast algorithms for the editing distance between trees and related problems. SIAM J. Comput. 1989.
3. N. Augsten, M.Bochlen, J.Gamper. Approximate matching of hierarchical data using pq-grams. VLDB. 2005