

Optimal Discovery of Subword Associations in Strings

(Extended abstract)

Alberto Apostolico^{1*}, Cinzia Pizzi^{2**}, and Giorgio Satta^{2***}

¹ University of Padova & Purdue University

² University of Padova

Abstract. Given a textstring x of n symbols and an integer constant d , we consider the problem of finding, for any pair (y, z) of subwords of x the number of times that y and z occur in tandem (i.e., with no intermediate occurrence of either one of them) within a distance of d symbols of x . Although in principle there might be n^4 distinct subword pairs in x , we show that it suffices to consider a family of only n^2 such pairs, with the property that for any neglected pair (y', z') , there is a corresponding pair (y, z) contained in our family and such that: (i) y' is a prefix of y and z' is a prefix of z , and (ii) the tandem index of (y', z') equals that of (y, z) . We show that an algorithm for the construction of the table of all such tandem indices can be built to run in optimal $O(n^2)$ time and space.

Keywords: Pattern Matching, String Searching, Subword Tree, Substring Statistics, Tandem Index, Association Rule.

1 Introduction

The problem of characterizing and detecting unusual events such as recurrent subsequences and other streams or over/under-represented words in sequences arises ubiquitously in diverse applications and is the subject of much study and interest in fields ranging from Computer and Network Security to Data Mining,

* Dipartimento di Ingegneria dell' Informazione, Università di Padova, Padova, Italy and Department of Computer Sciences, Purdue University, Computer Sciences Building, West Lafayette, IN 47907, USA. Work Supported in part by an IBM Faculty Partnership Award, by the Italian Ministry of University and Research under the National Projects FIRB RBNE01KNFP, and PRIN "Combinatorial and Algorithmic Methods for Pattern Discovery in Biosequences", and by the Research Program of the University of Padova, by NSF Grant CCR-9700276 and by NATO Grant CRG 900293. axa@dei.unipd.it

** Dipartimento di Ingegneria dell' Informazione, Università di Padova, Via Gradenigo 6/A, 35131 Padova, Italy. cinzia.pizzi@dei.unipd.it

*** Dipartimento di Ingegneria dell' Informazione, Università di Padova, Via Gradenigo 6/A, 35131 Padova, Italy. satta@dei.unipd.it

from Speech and Natural Language Processing to Computational Molecular Biology. It also gives rise to interesting modeling and algorithmic questions, some of which have displayed independent interest.

Among the problems in this class we find, for instance, the detection of all squares or palindromes in a string, for which optimal $O(n \log n)$ algorithms have long been known (refer, e.g., [3, 7] and references therein). It is not difficult to extend those treatments to germane problems such as the discovery of pairs of occurrences, within a given distance, of a same string, or of a string and its reverse, and so on.

In this paper, we concentrate on the problem of detecting repetitive phenomena that consist of unusually frequent *tandem* occurrences, within a pre-assigned distance in a string, of two distinct but otherwise unspecified substrings. By the two strings occurring in tandem, we mean that there is no intermediate occurrence of either one in between.

Specifically, we consider the following problem. Let x be a string of n symbols over some alphabet Σ and d some fixed non-negative integer. For any pair (y, z) of subwords of x , their *tandem index* $I(y, z)$ relative to x is the number of times that z has a closest occurrence in x within a distance of d from a corresponding, closest occurrence of y to its left. We are interested in finding pairs of subwords with surprisingly high tandem index.

The problem can be cast in the emerging contexts of *data mining* and *information extraction*. As is well known, while traditional data base queries aim at retrieving records based on their isolated contents, these contexts focus on the identification of patterns occurring across records, and aim at the retrieval of information based on the discovery of interesting rules present in large collection of data. Central to these developments is the notion of *association rule*, which is an expression of the form $S_1 \rightarrow S_2$ where S_1 and S_2 are sets of data attributes endowed with sufficient *confidence* and *support*. Sufficient support for a rule is achieved if the number of records whose attributes include $S_1 \cup S_2$ is at least equal to some pre-set minimum value. Confidence is measured instead in terms of the ratio of records having $S_1 \cup S_2$ over those having S_1 , and is considered sufficient if this ratio meets or exceeds a pre-set minimum. Clearly, a statistic of the number of records endowed with the given attributes must be computed as a preliminary step, and this is often a bottleneck for the process of information extraction. We refer to [1], [8] and [14] for a broader discussion of these concepts.

Back to our problem, we observe that, in principle, there might be n^4 distinct pairings of subwords of in x . We show, however, that it suffices to restrict attention to a family containing only only n^2 pairs, after which for any neglected pair (y', z') , there is a pair (y, z) in the family such that: (i) y' is a prefix of y and z' is a prefix of z , and (ii) the tandem index of (y', z') equals that of (y, z) . We show that an algorithm for the construction of the table of all such tandem indices can be built to run in optimal $O(n^2)$ time and space for a string x of n symbols.

A generalization of the notion of tandem, called *proximity word-association pattern*, has been proposed in [4] with similar motivations. A proximity word-

association pattern is defined as a tuple of k string components, each matching the input text at a distance smaller or equal than a given d from the matching of the previous component.³ Among several results, [4] implicitly defines an algorithm for counting matches of such patterns in an input string of length n , reporting a worst case time complexity of $O(d^k n^{k+1} \log n)$. A related algorithm was also proposed in [18], running in time $O(n^{k+1})$ but requiring $O(n^k)$ scans of the input string. In the present proposal we are therefore concerned with the case $k = 2$, and offer an asymptotical improvement over the algorithm in [4] with a worst case running time of $O(n^2)$ and a constant number of passes over the input string.

As mentioned, a tandem statistics of the kind collected by our algorithm finds application in many domains. In Section 4, we will briefly describe its uses in the discovery of dyadic structures in genome analysis and Part-of-Speech Tagging in Natural Language Processing.

2 Preliminaries

We begin by recalling an important “left-context” property from [5]. Given two words x and y , the *end-set* of y in x is the set of ending positions of *occurrences* of y in x , i.e., $endpos_x(y) = \{i : y = x_{i-|y|+1} \dots x_i\}$ for some i , $|y| \leq i \leq n$. Two strings y and z are equivalent on x if $endpos_x(y) = endpos_x(z)$. The equivalence relation instituted in this way is denoted by \equiv^x and partitions the set of all strings over Σ into equivalence classes. Thus, $[y]$ is the set of all strings that have occurrences in x ending at the same set of positions as y . In the example of the string *abaababaabaababaababa*, for instance, $\{ba, aba\}$ forms one such equivalence class and so does $\{aa, baa, abaa\}$. A symmetric equivalence relation can be defined in terms of *start-set* and denoted by \equiv_x , with obvious meaning. Recall that the *index* of an equivalence relation is the number of equivalence classes in it.

Fact 1 *The index k of the equivalence relation \equiv^x (alternatively, \equiv_x) obeys $k < 2n$.*

Fact 1 suggests that we might only need to look among $O(n^2)$ substring pairs of a string of n symbols in order to find unusually frequent pairs. The following considerations show that this statement can be made precise giving an indirect proof of it. We recall the notion of the *suffix tree* T_x associated with x . This is essentially a compact trie with $n + 1$ leaves and at most n internal nodes that collects all suffixes of $x\$$, where $\$$ a symbol not in Σ . We assume familiarity of the reader with the structure and its clever $O(n \log |\Sigma|)$ time and linear space constructions such as in [12, 16, 19]. The word ending precisely at vertex α of T_x is denoted by $w(\alpha)$. The vertex α is called the *proper locus* of $w(\alpha)$. The *locus* of

³ In [4] a different notation is used, with d representing the number of string components and k representing the maximum distance.

w is the unique vertex of T_x such that w is a prefix of $w(\alpha)$ and $w(\text{FATHER}(\alpha))$ is a proper prefix of w .

Having built the tree, some simple additional manipulations make it possible to count and locate the distinct (possibly overlapping) instances of any pattern w in x in $O(|w|)$ steps. For instance, listing all occurrences of w in x is done in time proportional to $|w|$ plus the total number of such occurrences, by reaching the locus of w and then visiting the subtree of T_x rooted at that locus. Alternatively, a trivial bottom-up computation on T_x can weigh each node of T_x with the number of leaves in the subtree rooted at that node. This weighted version serves then as a statistical index for x , in the sense that, for any w , we can find the frequency of w in x in $O(|w|)$ time. Note that the counter associated with the locus of a string reports its correct frequency even when the string terminates in the middle of an arc. This is, indeed, nothing but a re-statement and a proof of Fact 1, where the equivalence classes are represented by the nodes of the tree. From now on, we assume that a tree with weighted nodes has been produced and is available for our constructions.

3 Algorithms

For simplicity of exposition we assume that Σ is the binary alphabet $\Sigma = \{a, b\}$, but it should be clear that this assumption can be removed without penalty on our constructions.

In view of Fact 1 we may restrict attention to the $O(n)$ equivalence classes of \equiv_x represented by strings that end precisely at the internal nodes and leaves of T_x . Even these latter, each being formed by some consecutive prefixes of a singleton string, may be neglected as uninteresting.

As a warmup for the discussion, consider the problem of computing the following relaxed notion of a tandem index, which we denote by $\bar{I}(y, z)$ and consists of the number of instances of z that fall within d positions of a closest occurrence of y to its left. The difference with respect to $I(y, z)$ is that we still forbid intervening occurrences of y , but now allow possibly intervening occurrences of z .

It is easy to compute $\bar{I}(y, z)$ for a fixed word y and all z 's with a proper locus in T_x in overall linear time. This is done as follows. We can assume that, as a trivial by-product of the construction of T_x , there is access from each leaf of T_x to the corresponding position of x and *vice versa*. Now, let ν be the node of T_x such that $w(\nu) = y$. Visit the subtree of T_x rooted at ν and place a special mark on the positions of x that correspond to leaves in this subtree. This marks all the starting positions of occurrences of y in x . Next, scan the positions of x in ascending order: tag those positions that fall within d positions of an occurrence of y , and assign a weight of 1 to every leaf that corresponds to such a position. Visit now T_x bottom up, and assign to each node a weight equal to the sum of the weights of its children. Clearly, this accomplishes computation of $\bar{I}(y, z)$ for all z 's. Iterating the process for all y 's fills the table of $\bar{I}(y, z)$ values for all pairs

of words with a proper locus in T_x , and this takes (optimal) $O(n^2)$ time and space.

Consider now computing the $I(y, z)$ values for a given y and all strings z . If, in the bottom up computation above, we wanted to compute $I(y, z)$ instead of $\bar{I}(y, z)$, then we should prevent the weighting process from counting more than one occurrence of z within the d -range of each closest occurrence of y .

Let \mathcal{L} be sorted list of positions of x , with occurrences of y suitably marked and the positions falling within a distance of d from such occurrences tagged as like above. Define a *clump* in \mathcal{L} as a maximal run of tagged positions falling between two consecutive occurrences of y (or following the last occurrence of y). Let us say further that a clump is *represented* at a node μ of T_x if at least one element of the clump is found in the subtree of T_x rooted at μ . Clearly, we want each clump to contribute precisely one unit weight at all nodes where the clump itself is represented.

It is possible to achieve this by the following preprocessing of T_x . With k the total number of clumps, initialize k empty *clump lists*. Visit the (leaves of) tree from left to right. For each leaf encountered, check on its corresponding entry in \mathcal{L} whether this leaf belongs to a clump. In this case, assign to the leaf a *rank* equal to the ordinal number of arrival of the leaf in its clump according to the visit. At the end, the concatenation of the clump lists constitutes a sorted list of clumps such that leaves within each clump are met in order of ascending rank. This means that scanning each clump list now meets the leaves in the clump in the same order as they would be encountered when visiting T_x from left to right. At this point we invoke as a subroutine the following well known Lowest Common Ancestor (l.c.a.) Algorithm (see, e.g., [15]): given a tree T with n leaves, it is possible to preprocess T in time linear in the number of nodes and in such a way that, after preprocessing, for any two leaves i and j it is possible to give in constant time the lowest common ancestor of i and j .

For a given y , we pre-process all clumps in succession as follows. Singleton clumps are not touched. With reference now to the generic non-singleton clump, consider its leaves in order of ascending rank. For each leaf, find its l.c.a.'s relative to the leaf and its successor and give the weight of this node a -1 handicap. At the end of the process, weigh the tree by the bottom up weighting procedure as before.

We claim that the weights thus assigned to any node μ represent $I(y, w(\mu))$. This is seen by induction on the sparsification \bar{T} of T_x intercepted by the leaves of some arbitrary clump and the associated l.c.a. nodes. The assertion is true on any deepest internal node of \bar{T} . In fact, since T_x is binary then so is \bar{T} . Considering any deepest internal node of \bar{T} , the leftmost one of its two leaves will have caused the algorithm to give the node a handicap of -1. This would combine with the weight of 2 resulting from the bottom up weighting to yield a weight of 1. Assume now the assertion true for all descendants of a node μ of \bar{T} , we have that the two children α and β of μ will be assigned a weight of 1 in the bottom up process. Considering the rightmost leaf of the subtree rooted at α and the leftmost one in the subtree rooted at β , we have that μ , their l.c.a.,

has been given a handicap of -1 during pre-processing of the clump. Therefore, the sum of weights at μ will again be 1. The following theorem summarizes our discussion.

Theorem 2. *Given a string x of n symbols, the tandem indices for all pairs of subwords of x can be computed in $O(n^2)$ time and space.*

3.1 Implementation

In the tasks at hand, briefly discussed in Section 4, it is reasonable to assume that $|y|, |z|$, and d are much smaller than $|x|$. Therefore the algorithm was implemented on a *truncated* [13] rather than standard suffix tree data structure. This choice allows us to save a considerable amount of the working space at each run of the experiments. Truncated suffix trees are suffix trees in which the maximum length of a label from the root to a leaf is bounded [13] by some parameter N . The algorithm described in the previous sections needs some adjustments in order to work still correctly on a truncated suffix tree. The main problem is due to the fact that in a truncated suffix tree we lose the suffix tree property for which each leaf corresponds to one and only one string suffix. Therefore it can happen that the same leaf maps to several marked positions in the input string. To deal with this, we modify the algorithm as follows. Initially, each leaf is weighted with the sum of the corresponding number of marked positions. If these positions belong to different clumps, then no problem will arise in the total weight calculation, since each position will be treated separately at each clump analysis. On the other hand, if two (or more) positions belong to the same clump, then their associated node appears two or more consecutive times in the ordered clump list. In this case the LCA procedure is not executed. Instead, we decrement the value of the weight of that leaf. Then we proceed by comparing the next pair of nodes in the clump, and follow the above procedure if the nodes are again the same, or the standard algorithm procedure otherwise.

4 Applications

Here we describe two example applications of our algorithm, the first one to Bioinformatics, the second to Natural Language Processing.

4.1 Dyad Analysis

Dyad analysis for discovering cis-acting regulatory patterns in gene promoter regions are based on statistics on co-occurrences of short solid patterns. The method is based on the observation that many regulatory sites can be modeled by a structure $w_1 n_s w_2$, where w_1 and w_2 are solid words of length 3, s the distance between the two, and n_s a sequence of s unspecified nucleotides. We refer to [17] for a detailed description and validation of the method. Although experimental results proved the method to be efficient for the detection of sites bound by

C_6Zn_2 binuclear cluster proteins and other transcription factors, the exhaustive search through all the possible trinucleotides and all the possible pairs among them is time consuming and prevents a more efficient use of the tool. Moreover, especially for longer patterns, reasoning in terms of the equivalence classes of \equiv_x instead of fixed length words is more meaningful in so far as the dyads detected in this way are *maximal* in an obvious sense.

4.2 Part-of-Speech Tagging

Another target applications of the algorithms developed in this paper is computational learning for natural language processing. More specifically, we focus on the task of part-of-speech tagging [10], where words in the text must be classified for their lexical category (Noun, Verb, Adjective, etc.).

A common solution in the design of part-of-speech taggers is to exploit sets or cascades of rules that are automatically induced from examples of correct classifications. Usually, these rules classify words on the basis of finite size contextual windows, centered around the word of interest, as described for instance in [6]. However, such systems do not achieve 100% accuracy in classification, as compared with classifications provided by human experts. One of the reasons for such a gap is due to the limited size of the adjacency windows that are used by the rule schemata described above. It has been shown in [9] that certain contextual patterns called *barriers*, that are highly effective in classification, can appear at variable distances from the word of interest, demanding for more powerful rule patterns than those described above.

The tandem patterns that are discovered by the algorithm of Section 3 can be applied to the task of automatic extraction of contextual rules for part-of-speech tagging, by simply modifying the algorithm to count the end-to-start distance instead than the previously described start-to-start distance between tandem components. Preliminary results on the CommonNoun/Adjective confusion class showed that this method could be helpful when a set of fixed single window rules fail to give a clear classification.

5 Conclusions

We have investigated the problem of collecting counts on the number of times any pair of strings occur in tandem (i.e., with no intermediate occurrence of either one of them) in a given input text and within a given distance. We have provided an algorithm for the construction of the table of all such tandem counts, running in optimal $O(n^2)$ time and space. This improves on previous results that are found in the literature. We have also discussed two possible application of the algorithm: dyad analysis in bioinformatics, and part-of-speech tagging in natural language processing.

References

1. Agrawal, R., Imielinski, T., Swami, A.: Mining Association Rules between Sets of Items in Large Databases. In *Proc. ACM SIGMOD*, 207–216, Washington DC May 1993.
2. Aho, A.V., Hopcroft, J.E., Ullman, J.D.: *The Design and Analysis of Computer Algorithms*, Addison-Wesley, Reading, Mass (1974)
3. Apostolico, A., Galil, Z. (Eds.): *Pattern Matching Algorithms*, Oxford University Press, New York (1997)
4. Arimura, H., Arikawa, S.: Efficient Discovery of Optimal Word-Association Patterns in Large Text Databases. *New Generation Computing*, 18: 49–60 (2000)
5. Blumer, A., Blumer, J., Ehrenfeucht, A., Haussler, D., Chen, M.T. and Seiferas, J.: The Smallest Automaton Recognizing the Subwords of a Text. *Theoretical Computer Science*, 40: 31–55 (1985).
6. Brill, E.: Transformation-Based Error-Driven Learning and Natural Language Processing: A Case Study in Part of Speech Tagging, *Computational Linguistics* (1995)
7. Crochemore, M., Rytter, W.: *Text Algorithms*, Oxford University Press, New York (1994)
8. Han, J., and Kamber, M.: *Data Mining: Concepts and Techniques*, Morgan Kaufmann Publishers, 2000.
9. Karlsson, F., Voutilainen, A., Heikkilä, F., and Anttila, A.: Constraint Grammar. A Language Independent System for Parsing Unrestricted Text. Mouton de Gruyter (1995)
10. Manning C.D., Schütze, H.: *Foundations of Statistical Natural Language Processing*, MIT Press (1999)
11. Marcus, M. P., Santorini, B., and Marcinkiewicz, M.A.: Building a Large Annotated Corpus of English: The Penn Treebank, *Computational Linguistics*, 19(2):313–330 (1993)
12. McCreight, E.M.: A Space-Economical Suffix Tree Construction Algorithm. *Journal of the ACM*, 23(2): 262–272, April 1976.
13. Na, J.C., Apostolico A., Iliopoulos C.S., Park K., Truncated Suffix Trees and their Application to Data Compression *Theoretical Computer Science*, Vol. 304, Issue 1-3, pp.87–101 (2003)
14. Piatesky-Shapiro, G., Frawley W.J. (eds.): *Knowledge Discovery in Databases*. AAAI Press/MIT Press (1991).
15. Schieber, B., Vishkin, U.: On Finding Lowest Common Ancestors: Simplifications and Parallelizations, *SIAM Journal on Computing*, 17:1253–1262 (1988)
16. Ukkonen, E.: On-line Construction of Suffix Trees. *Algorithmica*, 14(3): 249–260 (1995)
17. van Helden, J., Rios, A.F., Collado-Vides, J.: Discovering Regulatory Elements in Non-coding Sequences by Analysis of Spaced Dyads. *Nucleic Acid Research*, Vol.28, No.8: 1808–1818, 2000.
18. Wang, J. T.-L., Chirn, G.-W., Marr, T.G., Shapiro, B., Shasha, D., Zhang, K.: Combinatorial Pattern Discovery for Scientific Data: Some Preliminary Results. In *Proceedings of 1994 SIGMOD*, 115–125 (1994)
19. Weiner, P.: Linear Pattern Matching algorithm. In *Proceedings of the 14th Annual IEEE Symposium on Switching and Automata Theory*, pages 1–11, Washington, DC (1973)