



Efficiently Detecting Disguised Web Spambots (with Mismatches) in a Temporally Annotated Sequence



Authors:

Hayam Alamro & Costas S. Iliopoulos

Presenter: [\(Hayam Alamro\)](#)

- Department of Informatics, King's College London, UK
- Department of Information Systems, Princess Nourah bint Abdulrahman University, Riyadh, KSA.
- (hayam.alamro@kcl.ac.uk)

Short resume of the presenter :

- Hayam Alamro is a Ph.D. student in Computer Science (Algorithms & Data analysis Research Group) in the Department of Informatics at King's College London.
- Hayam's research focuses on the analysis and advanced design of string algorithms, approximate pattern matching, Cybersecurity, and data privacy.
- Hayam received her M.Sc. and her B.Sc. (with second class Honour) in Computer Science and Information Systems from King Saud University, Riyadh, Kingdom of Saudi Arabia.
- Before starting her Ph.D. in the UK, Hayam was working as a Lecturer in Computer Science and Information Systems College in Princess Nora University, Riyadh, Kingdom of Saudi Arabia. Hayam also has an experience working as a Computer Assistance in the Ministry of Planning, Interest of Public Statistics, Riyadh, Kingdom of Saudi Arabia.

- ❑ Introduction
- ❑ Our contributions
- ❑ Background
- ❑ Disguised Actions Definition
- ❑ Problems Definition
- ❑ Disguised Actions with K Mismatches
- ❑ Illustration by Example
- ❑ Conclusion

Introduction

- A *spambot* is a computer program designed to do repetitive actions on websites, servers or social media communities.



<https://images.app.goo.gl/a5Yreu3X7MSCHmvU7>

Activities

- Carrying out certain attacks on websites/ servers.
- Involving irrelevant links to increase a website ranking in search engine results.
- Using web crawlers for planting unsolicited material.
- Collecting email addresses from different sources (*phishing emails*).

Introduction

Existing spam detection techniques

- ***Content-based*** : Inject repetitive keywords in meta tags to promote a website in search engines.
- ***Link-based*** : Add links to a web page to increase its ranking score in search engines.
- ***Supervised machine learning***: to identify the source of spambot, rather than detecting the spambot.

Nowdays: The spammers try to manipulate spambots' actions behaviour to appear as it were coming from a legitimate user to bypass the existing spam-filter tools



<https://images.app.goo.gl/a5Yreu3X7MSCHmvU7>

Introduction

More relevant to our work



<https://images.app.goo.gl/a5Yreu3X7MSCHmvU7>

- *String pattern matching-based techniques* (Hayati et al. and Ghanaei et al.),

But:

- They are inapplicable because they do not take into account temporal information of neither the sequence (i.e., the user log) nor the pattern (i.e., the spambot actions).

- P. Hayati, V. Potdar, A. Talevski, and W. Smyth, “Rule-based on-the-fly web spambot detection using action strings,” in CEAS, 2010.
- V. Ghanaei, C. S. Iliopoulos, and S. P. Pissis, “Detection of web spambot in the presence of decoy actions,” in IEEE International Conference on Big Data and Cloud Computing, 2014, pp. 277–279.

Our Contributions

1. We introduce an efficient algorithm that can detect one or more sequences of indeterminate (non solid) actions in text T in linear time.

→ Our algorithm can compute all occurrences of indeterminate sequence \tilde{S} in text T in $O(m + \log n + occ)$, where:

m is the $|\tilde{S}|$, n is the $|T|$, and occ is the number of the occurrences of the sequence \tilde{S} in T .

2. We propose an efficient algorithm for solving (f, c, k, W) -***Disguised Spambots Detection with indeterminate actions and mismatches***. It is a generalization of the previous problem (1).

➔ Our algorithm takes into account temporal information, because it considers:

- Time-annotated sequences, and
- Requires a match to occur within a time window t .

Background

- Let $T = a_0 \dots a_{n-1}$ be a string of length $|T| = n$ over an alphabet Σ of size $|\Sigma| = \sigma$
- For $1 \leq i \leq j \leq n$, $T[i]$ denotes the i th symbol of T , and $T[i, j]$ the contiguous sequence of symbols (called *factor* or *substring*)
- A substring $T[i, j]$ is a **suffix** of T if $j = n$ and it is a **prefix** of T if $i = 1$
- A string p is a **repeat** of T iff p has at least two occurrences in T
- A **degenerate or indeterminate string**, is defined as a sequence $\tilde{X} = \tilde{x}_0 \tilde{x}_1 \dots \tilde{x}_{n-1}$, where $\tilde{x}_i \subseteq \Sigma$ for all $0 \leq i \leq n - 1$
- A **degenerate symbol** \tilde{x} over an alphabet Σ is a non-empty subset of Σ

Background

- $|\tilde{x}|$ denotes the size of \tilde{x} , and we have $1 \leq \tilde{x} \leq |\Sigma|$.
- If $|\tilde{x}_i| = 1$, that is $|\tilde{x}_i|$ repeats a single symbol of Σ , we say that \tilde{x}_i is a **solid symbol** and i is a **solid position**. Otherwise, \tilde{x}_i and i are said to be a **non-solid symbol** and **non-solid position** respectively.
- A **conservative degenerate string** is a degenerate string where its number of non-solid symbols is upper-bounded by a fixed position constant c .

Example

$$X = ab[ac]a[bcd]bac$$

Is a degenerate string of length 8 over the alphabet $\Sigma = \{a, b, c, d\}$, and conservative degenerate string with $c = 2$.

Background

- A ***suffix array*** of T is the lexicographical sorted array of the suffixes of a string T i.e., the suffix array of T is an array $SA[1...n]$ in which $SA[i]$ is the i th suffix of T in ascending order.
- ***LCP***(T_1, T_2) is the length of the longest common prefix between strings T_1 and T_2 , and it is usually used with SA such that $LCP[i] = \text{lcp}(T_{SA[i]}, T_{SA[i-1]})$ for all $i \in [1..n]$.

Disguised (Indeterminate) Actions

- Some spambots might attempt to disguise their actions by varying certain actions.

Example

a spambot takes the actions **ABCDEF**, then **ACCDEF**, then **ABDDEF**

can be described as $\rightarrow A[BC][CD]DEF$

The action **[BC]** and **[CD]** are variations of the same sequence

- A, C, D, E, F \rightarrow (solid symbols)
- [BC] [CD] \rightarrow (indeterminate or non-solid symbols)
- A[BC][CD]DEF \rightarrow (degenerate string)

Problems Definitions

A. Given a sequence $T = a_1 \dots a_n$, and an action sequence $\tilde{S} = s_1 s_2 \dots s_m$, find all occurrences of \tilde{S} in T where s_i might be solid or indeterminate.

B. Given a temporal annotated sequence $T = (a_1, t_1) \dots (a_n, t_n)$, and an action sequence $\tilde{S} = s_1 s_2 \dots s_m$, find all occurrences of \tilde{S} in T in time window t , where s_i might be solid or indeterminate with hamming distance between \tilde{S} and T is no more than k mismatches.

Our Main Problem

Since Problem B is a generalization of Problem A, we will focus on Problem B in this presentation.

(f, c, k, W) -Disguised Spambots Detection with indeterminate actions and mismatches:

Given a temporal annotated sequence $T = (a_1, t_1) \dots (a_n, t_n)$, a dictionary \bar{S} containing sequences \hat{S}_i , each has a c non-solid symbol (represented by #), associated with a time window W_i , a minimum frequency threshold f , and a maximum Hamming distance threshold K , find all occurrences of each $\hat{S}_i \in \bar{S}$ in T , such that each \hat{S}_i occurs:

- I. At least f times within its associated time window W_i , and
- II. With at most K mismatches according to Hamming distance.

Disguised Actions with K Mismatches

Preprocessing Stage

Our algorithms require as input sequences temporally annotated actions. These temporally annotated sequences are produced from *user logs* consisting of a collection of *http requests*.

```
125.127.33.125 - smith [14/Oct/2019: 10:12:26 -0500] "GET /blog/page-address.htm HTTP/1.1" 200 1043  
"-" "Mozilla/5.0 Chrome/80.0.3134.311"
```

Date – Time Stamp

Action Request

Time Point t

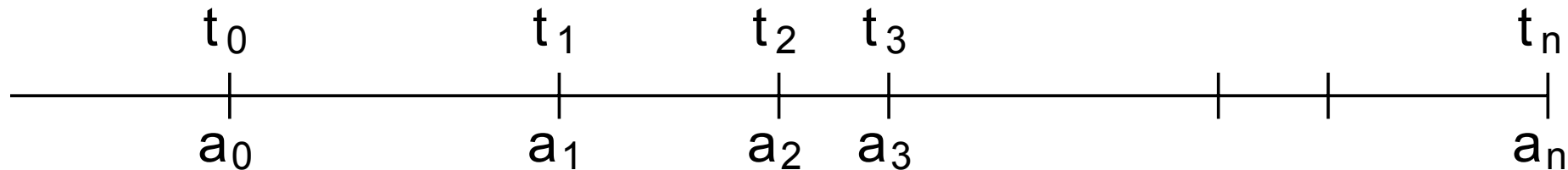
Index Key

Disguised Actions with K Mismatches

Definition

A Temporally Annotated Action Sequence: is a sequence

$T = (a_0, t_0), (a_1, t_1), \dots, (a_n, t_n)$, where $a_i \in A$, with A set of actions, and t_i represents the time that action a_i took place. Note that $t_i < t_{i+1}$, $\forall i \in [0, n]$.



Disguised Actions with K Mismatches

Our Spambot Dictionary Representation

i	S_i	W_i (sec)
1	cbbx	20
2	byadc	25
...

Disguised Actions with K Mismatches

Algorithm Steps:

Step 1: For each *non-solid* s_j occurring in degenerate sequence $\tilde{S} = s_1 s_2 \dots s_m$, we substitute each s_j with '#' symbol, where '#' is not in Σ . Let \hat{S} be the resulting pattern after substitution.

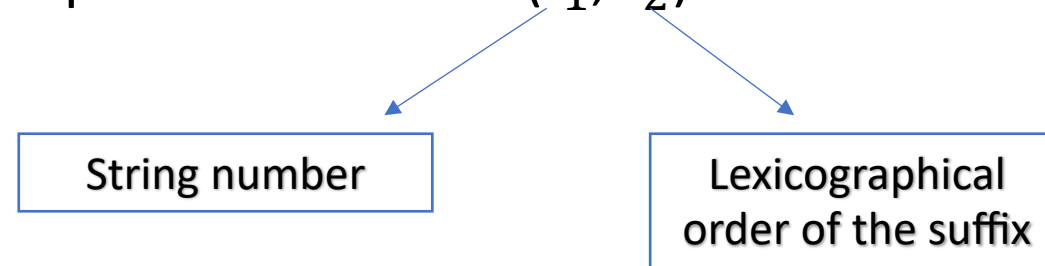
\tilde{S}	A	B	[GX]	C	[AD]	F
\hat{S}	A	B	# ₁	C	# ₂	F

Disguised Actions with K Mismatches

Step 2: Extract the actions of the temporally annotated sequence T into T_a such that it contains only the actions $a_1 \dots a_n$ from T .

Step 3: Build **Generalized Enhanced Suffix Array (GESA)**:

it is an enhanced suffix array for a set of strings, each one ending with a special character and usually is built to find the *Longest Common Sequence (LCS)* of two strings or more. GESA is indexed as a pair of identifiers (i_1, i_2) .



Disguised Actions with K Mismatches

Generalized Enhanced Suffix Array for a collection of texts (T_a and $\bar{S}_{\widehat{S_i}}$):

$$GESA(T_a, \bar{S}_{\widehat{S_i}}) = Ta!_0 \hat{S}_1!_1 \hat{S}_2!_2 \dots \hat{S}_r!_r$$

- $\hat{S}_1, \dots, \hat{S}_r$ are spambots sequences belong to dictionary \bar{S}
- $!_0, \dots, !_r$ are special delimiters not in Σ , and smaller than any alphabetical letter in T_a and smaller than '#'

Disguised Actions with K Mismatches

- Our algorithm will refer to a collection of tables $(GES A, GES A^R, LCS, T, \bar{S}_{\widehat{S}_i}) \rightarrow$ to find disguised spambots within a time window t as follows:
 - Given a temporally annotated action sequence $T = (a_0, t_0)(a_1, t_1) \dots (a_{n-1}, t_{n-1})$, an action sequence $\widehat{S} = s_1 s_2 \dots s_m$, and an integer t , we will compute j_1, j_2, \dots, j_m such that $a_{j_i} = s_i, 1 \leq i \leq m$ and $\sum_{i=1}^m t_{j_i} < t$ or $t_{j_m} - t_{j_1} < t$ with **Hamming distance between T_a and \widehat{S} no more than k mismatches.**

Disguised Actions with K Mismatches

- Our algorithm, also includes:
 - Initialization for *hashMatchTable* to do *bit masking* operation
 - *Kangaroo method* to find the *Longest Common Sequence LCS* between a sequence of actions in T and an action sequence \hat{S}_i with at most K mismatches in linear time.

Disguised Actions with K Mismatches

- **Step 4:** For each \hat{S}_i in the spambots dictionary $\bar{S}_{\hat{S}_i}$, the algorithm calculates the *Longest Common Sequence LCS* between \hat{S}_i and T_a starting at position 0 in sequence \hat{S}_i and at position j in sequence T_a , such that the common substring starting at these positions is maximal as follows:
 - Find the suffix index i of \hat{S}_i in $GESA$ using $GESA^R$ table (retains all the lexicographical ranks of the suffixes of $GESA$).
 - Find the closest suffix j (belongs to a sequence in T_a) to the suffix i (of \hat{S}_i) in $GESA$.

Disguised Actions with K Mismatches

- Compute the **Longest Common Sequence LCS** between $GESA(i)$ and $GESA(j)$ as follows:

$$LCS(\hat{S}_i, T_a) = \max(LCP(GESA(i_1, i_2), GESA(j_1, j_2))) = l_0$$

Where l_0 is the maximum length of the *longest common subsequence* matching characters between $GESA(i_1, i_2)$ and $GESA(j_1, j_2)$ until the first mismatch (or one of the sequences terminates).

Disguised Actions with K Mismatches

- Next, find the length of the *longest common subsequence* matching characters after previous mismatch position l_0 using **Kangaroo method** until the second mismatch (or one of the sequences terminates) as follows:

$$\mathit{max}(\mathit{LCP}(\mathit{GESA}(i_1, i_2 + l_0 + 1), \mathit{GESA}(j_1, j_2 + l_0 + 1))) = l_1$$

Disguised Actions with K Mismatches

- Once our algorithm encounters '#' in the sequence \hat{S}_i , it will get into the **verification process**:
 - Query whether the corresponding action a_i (in T_a) belongs to the set of actions in '#', to do that:
 - The algorithm uses a **bit masking** operation (**And** bit wise operation) between the two sets ('#' and a_i) such that (a_i is represented by a bit '1', and each action in '#' is represented by '1' and '0' otherwise using **hashMatchTable**).

Disguised Actions with K Mismatches

hashMatchTable

- The columns are indexed by the (ascii code) of each character ($a_i \in \Sigma$) (for directly access)
- The rows are indexed by the number of the spambots sequence \hat{S}_i and the number of ' $\#_l$ '
- The algorithm applies the following formula:

$$1 \wedge hashMatchTable[\hat{S}_r \#_l][ascii[a_i]]$$

Disguised Actions with K Mismatches

hashMatchTable

$$\widetilde{S}_1 = AB[GX]C[AD]F \rightarrow \widehat{S}_1 = AB\#_1C\#_2F$$

Ascii(a_i)	65 A	66 B	67 C	68 D	...	71 G	...	88 X	89 Y	90 Z
$\widehat{S}_1\#_1$	0	0	0	0	...	1	...	1	0	0
$\widehat{S}_1\#_2$	1	0	0	1	...	0	...	0	0	0
...
$\widehat{S}_r\#_l$

Disguised Actions with K Mismatches

- **Step 5:** Finally, at each occurrence of $\hat{\mathcal{S}}_i$ in the sequence T_a , our algorithm checks its time window W_i using the dictionary $\overline{\mathcal{S}}_{\hat{\mathcal{S}}_i}$ in T , such that it sums up each time t_i associated with its action a_i starting at the position j_2 in $GESA(j_1, j_2)$ until the length of the spambot $|\hat{\mathcal{S}}_i|$, and then compare it to its time window W_i . If the resultant time is less than or equal to W_i , the algorithm considers that the sequence $\hat{\mathcal{S}}_i$ is spambots and terminates.
- The algorithm will continue to find other occurrences of the spambot sequence $\hat{\mathcal{S}}_i$ using the adjacent suffixes to the suffix index of $\hat{\mathcal{S}}_i$ in GESA.

Illustration by Example

Example

- Suppose : $T_a = ABBABGCDFCBACAF AABGDFF$, $\hat{S}_1 = B \#_1 C \#_2 F$
- $\#_1 = [GX]$, $\#_2 = [AD]$, $K = 2$, $f = 2$
- Concatenation strings of $Ta!_0 \hat{S}_1!_1$

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28
A	B	B	A	B	G	C	D	F	C	B	A	C	A	F	A	A	B	G	D	F	F	$!_0$	B	$\#1$	C	$\#2$	F	$!_1$

^

Illustration by Example

i	$GESA[i]$	Suffix	$GESA^R[i]$
0	(1,28)	$!_1$	5
1	(0,22)	$!_0 b \#_1 c \#_2 f !_0$	13
2	(1,24)	$\#_1 c \#_2 f !_1$	11
3	(1,26)	$\#_2 f !_1$	6
4	(0,15)	$a a b g d f f !_0 b \#_1 c \#_2 f !_1$	14
5	(0,0)	$a b b a b g c d f c b a c a f a a b g d f f !_0 b \#_1 c \#_2 f !_1$	27
6	(0,3)	$a b g c d f c b a c a f a a b g d f f !_0 b \#_1 c \#_2 f !_1$	19
7	(0,16)	$a b g d f f !_0 b \#_1 c \#_2 f !_1$	20
8	(0,11)	$a c a f a a b g d f f !_0 b \#_1 c \#_2 f !_1$	25
9	(0,13)	$a f a a b g d f f !_0 b \#_1 c \#_2 f !_1$	18
10	(1,23)	$b \#_1 c \#_2 f !_1$	12
11	(0,2)	$b a b g c d f c b a c a f a a b g d f f !_0 b \#_1 c \#_2 f !_1$	8
12	(0,10)	$b \underline{a} c a f a a b g d f f !_0 b \#_1 c \#_2 f !_1$	17
13	(0,1)	$b b a b g c d f c b a c a f a a b g d f f !_0 b \#_1 c \#_2 f !_1$	9
14	(0,4)	$b \underline{g} c d f c b a c a f a a b g d f f !_0 b \#_1 c \#_2 f !_1$	24
15	(0,17)	$b \underline{g} d f f !_0 b \#_1 c \#_2 f !_1$	4
16	(1,25)	$c \#_2 f !_1$	7
17	(0,12)	$\underline{c} a f a a b g d f f !_0 b \#_1 c \#_2 f !_1$	15
18	(0,9)	$c b a c a f a a b g d f f !_0 b \#_1 c \#_2 f !_1$	28
19	(0,6)	$c d f c b a c a f a a b g d f f !_0 b \#_1 c \#_2 f !_1$	21
20	(0,7)	$d f c b a c a f a a b g d f f !_0 b \#_1 c \#_2 f !_1$	26
21	(0,19)	$d f f !_0 b \#_1 c \#_2 f !_1$	23
22	(1,27)	$f !_1$	1
23	(0,21)	$\underline{f} !_0 b \#_1 c \#_2 f !_1$	10
24	(0,14)	$f a a b g d f f !_0 b \#_1 c \#_2 f !_1$	2
25	(0,8)	$f c b a c a f a a b g d f f !_0 b \#_1 c \#_2 f !_1$	16
26	(0,20)	$\underline{f} f !_0 b \#_1 c \#_2 f !_1$	3
27	(0,5)	$g c d f c b a c a f a a b g d f f !_0 b \#_1 c \#_2 f !_1$	22
28	(0,18)	$g d f f !_0 b \#_1 c \#_2 f !_1$	0

- There are three occurrences for \hat{S}_1 in T at $i = 12, 14$ and 15

Disguised Actions with K Mismatches

Experimental Evaluation

- Our experiments showed that our algorithm is efficient and applicable to large action sequences.
- See our paper for details.

Conclusion

- We introduced our algorithm (f, c, k, W) -Disguised Spambots Detection with indeterminate actions and mismatches.
- Our algorithm takes into account temporal information, because it considers time-annotated sequences, and because it requires a match to occur within a time window.
- The problem seeks to find all occurrences of each conservative degenerate sequence corresponding to a spambot that occurs at least f times within a time window and with up to k mismatches.

Conclusion

- For this problem, we designed a linear time and space inexact matching algorithm, which employs the ***generalized enhanced suffix array*** data structure, ***bit masking*** and ***Kangaroo method*** to solve the problem efficiently.

Thank You