

Program Analysis for Adaptive Data Analysis

ANONYMOUS AUTHOR(S)

An adaptive data analysis is based on multiple queries over a data set, in which some queries rely on the results of some other queries. The error of each query is usually controllable and bound independently, but the error can propagate through the chain of different queries and bring to high generalization error. To address this issue, data analysts are adopting different mechanisms in their algorithms, such as Gaussian mechanism, etc. To utilize these mechanisms in the best way one needs to understand the depth of chain of queries that one can generate in a data analysis. Unfortunately, this depth which relies on the program(implementation) itself is costly in human efforts, and how to statically obtain this information is not well studied to support data analysts.

In this work, we define a programming framework which can provide, through analysis on some transformation of the target program implementing an adaptive data analysis, an upper bound on the adaptivity depth (the length of the longest chain of queries) along with another upper bound on the total number of queries requested. We show how this framework can help to analyze the generalization error of two data analyses with different adaptivity structures.

Additional Key Words and Phrases: data analysis, program analysis, dependency graph

1 INTRODUCTION

Consider a dataset X consisting of n independent samples from some unknown population P . How can we ensure that the conclusions drawn from X *generalize* to the population P ? Despite decades of research in statistics and machine learning on methods for ensuring generalization, there is an increased recognition that many scientific findings generalize poorly (e.g.). While there are many reasons a conclusion might fail to generalize, one that is receiving increasing attention is *adaptivity*, which occurs when the choice of method for analyzing the dataset depends on previous interactions with the same dataset.

Adaptivity can arise from many common practices, such as exploratory data analysis, using the same data set for feature selection and regression, and the re-use of datasets across research projects. Unfortunately, adaptivity invalidates traditional methods for ensuring generalization and statistical validity, which assume that the method is selected independently of the data. The misinterpretation of adaptively selected results has even been blamed for a “statistical crisis” in empirical science.

A recent line of work initiated by Dwork *et al.* [Dwork et al. 2015] and Hardt and Ullman posed the question: Can we design *general-purpose* methods that ensure generalization in the presence of adaptivity, together with guarantees on their accuracy? This line of work has identified many new algorithmic techniques for ensuring generalization in adaptive data analysis, leading to algorithms with greater statistical power than all previous approaches. It has also identified problematic strategies for adaptive analysis, showing limitations on the statistical power one can hope to achieve.

A key development in this line of work is that the best method for ensuring generalization in an adaptive data analysis depends to a large extent on the number of *rounds of adaptivity*. That is, not only does the analysis of the generalization error depend on the number of rounds, but knowing the number of rounds actually allows one to choose methods that lead to the smallest possible generalization error.

It is obviously promising if this number of *rounds of adaptivity* can be estimated statically. Nevertheless, the corresponding study on static analysis over programs implementing adaptive

2018. 2475-1421/2018/1-ART1 \$15.00

<https://doi.org/>

analysis is not well explored, even the formal definition of this *rounds of adaptivity* is not clear. Our goal is to study adaptivity and conduct analysis over programs to find the number of rounds. To this end, we design a programming framework, which statically provides an upper bound of the number of rounds of adaptivity for programs implementing adaptive data analysis algorithms. The first question we need to decide on is what language the target programs to be analyzed are written, in a functional programming language or an imperative one?

We choose the imperative programming language, then the following question is what does this imperative language look like? A principle of the ideal language in our mind is supposed to be simple enough to alleviate the burden of analysis on complex components that is unnecessary to adaptivity, and expressive enough to support most of the adaptive data analysis algorithms. With this principle in mind, we introduce an imperative language called the loop language, with one finite loop construct, allowing data analysts recursively request queries in their programs. The query request is also supported in the language.

For the number of rounds of adaptivity of program in the loop language, a definition of "one query relies on another query", or "one query may depend on the other" is the next big thing. We choose to define this "may-dependency" with the help of a trace-based operational semantics. The trace is a list of queries generated along with the execution of the program, a history of the execution only caring about query requests.

Our framework includes transformations of the target programs in loop language, to "ssa-form" programs. The programs in single static assignment form, with all the variables assigned once (variables in the loop is also handled and will be discussed in later section), greatly helps our analysis avoid the complexity of handling variables overwriting. The direct analysis on program in ssa form then allows our framework to provide the upper bound on the number of rounds of adaptivity. To be more specific, our framework is equipped with algorithms to record the necessary information and construct a dependency graph to reach a final estimation of the adaptivity.

Our contribution lies in the static analysis on programs implementing adaptive analysis for a tight estimation of its number of rounds of adaptivity through a framework. Specifically, our contributions can be broken into the following:

- (1) A formal definition of dependency between queries in adaptive data analysis programs.
- (2) A formal definition of adaptivity through a query-based dependency graph, bases on a trace-based operational semantics for the loop language.
- (3) A transformation between programs in loop language and its ssa form language, with the soundness of the transformation shows that every program in loop language can be appropriately transformed to its ssa form.
- (4) An matrix-and-vector-based algorithm to track the dependency between unique variables of ssa form programs (including loop) and construct a variable-based dependency graph for our static estimation of rounds of adaptivity.
- (5) A soundness proof shows that the predicted upper bound from our algorithm strictly bounds the actual number of rounds of adaptivity in execution.

2 OVERVIEW

In this section, we start with a review of the prior work of adaptive data analysis, which motivates our work, a framework to statically give an upper bound on the rounds of adaptivity . We then show the architecture of the framework and give our readers a taste from simple examples.

2.1 Review of Prior Work

To explain the key concepts and motivate our work, we review the model of adaptive data analysis of [??]. We remark that many of the details of the model are immaterial for our work, and are included for illustrative purposes and concreteness.

In the model there is some distribution P over a domain \mathcal{X} that an analyst would like to study. For our purposes, the analyst wishes to answer *statistical queries* (also known as *linear queries* or *linear functionals*) on the distribution. A statistical query is defined by some function $f : \mathcal{X} \rightarrow [-1, 1]$. The analyst wants to learn the *population mean*, which (abusing notation) is defined as

$$f(P) = \mathbb{E}_{X \sim P} [f(X)].$$

However, the distribution P can only be accessed via a set of *samples* X_1, \dots, X_n drawn identically and independently from P . These samples are held by a mechanism $M(X_1, \dots, X_n)$ who receives the query f and computes an answer $a \approx f(P)$.

In this work we consider analysts that ask a sequence of k queries f_1, \dots, f_k . If the queries are all chosen in advance, independently of the answers then we say they are *non-adaptive*. If the choice of each query f_j may depend on the prefix $f_1, a_1, \dots, f_{j-1}, a_{j-1}$ then they are *fully adaptive*. An important intermediate notion is *r-round adaptive*, where the sequence can be partitioned into r batches of non-adaptive queries. Note that non-interactive queries are 1-round and fully adaptive queries are k rounds.

We now review what is known about the problem of answering r -round adaptive queries. Given the samples, the naïve way to approximate the population mean is to use the *empirical mean*, which (abusing notation) is defined as

$$f(X_1, \dots, X_n) = \frac{1}{n} \sum_{i=1}^n f(X_i).$$

THEOREM 2.1. *For any distribution P , and any k non-adaptive statistical queries, the naïve mechanism satisfies*

$$\max_{j=1, \dots, k} |a_j - f_j(P)| = O\left(\sqrt{\frac{\log k}{n}}\right)$$

For any $r \geq 2$ and any r -round adaptive statistical queries, it satisfies

$$\max_{j=1, \dots, k} |a_j - f_j(P)| = O\left(\sqrt{\frac{k}{n}}\right)$$

And there exists analysts that make these bounds tight (up to constant factors).

Note that even allowing one extra round of adaptivity leads to an exponential increase in the generalization error from $\log k$ to k .

Perhaps surprisingly, Dwork *et al.* [?] and Bassily *et al.* [?] showed that an alternative mechanism can actually achieve much stronger generalization error as a function of the number of queries, specifically.

THEOREM 2.2 ([??]). *For any k , there exists a mechanism such that for any distribution P , and any $r \geq 2$ any r -round adaptive statistical queries, it satisfies*

$$\max_{j=1, \dots, k} |a_j - f_j(P)| = O\left(\frac{\sqrt[k]{k}}{\sqrt{n}}\right)$$

And there is an analyst that make this bound tight (up to constant factors).

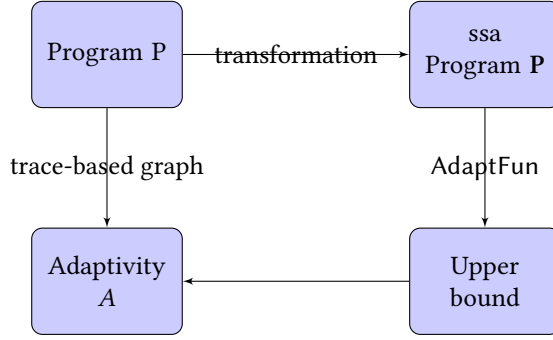


Fig. 1. High level architecture

Notice that Theorem 2.2 has different quantification in that the optimal choice of mechanism depends on the number of queries. Thus, we need to know the number of queries *a priori* to choose the best mechanism.¹

Later work by Dwork gave more refined bounds in terms of the number of rounds of adaptivity. Similarly, if one knows a good *a priori* upper bound on the number of rounds of adaptivity, one can get a much better guarantee of generalization error, but only by using an appropriate choice of the mechanism. Specifically,

THEOREM 2.3 ([?]). *For any r and k , there exists a mechanism such that for any distribution P , and any $r \geq 2$ any r -round adaptive statistical queries, it satisfies*

$$\max_{j=1,\dots,k} |a_j - f_j(P)| = O\left(\frac{r\sqrt{\log k}}{\sqrt{n}}\right)$$

And there is an analyst that make this bound tight (up to constant factors).

2.2 Static analysis framework for adaptivity

We show a framework that gives a good *a priori* upper bound on the number of rounds of adaptivity, to achieve the aforementioned refined bounds.

The architecture of the framework is shown in Figure 1. We are interested in the number of rounds of adaptivity (we use adaptivity as a short in the rest of the paper) of a program P in the loop language. We use a simplified adaptivity data analysis algorithm- two round strategy TR to illustrate how the framework works to produces an upper bound on adaptivity in Figure 2. The two round algorithm consists of two main steps. The first one is asking the database through query request $q()$ for k times and storing results in a list a . The second step is to construct another query $q_2(a)$ using the list a , and asks for the result from database by requesting the new query $q_2(a)$. In the loop language, the query request is assigned by $x \leftarrow q(e)$. In the above example, we simplify the query request, for example, query $q()$ in the first step does not have argument and the new constructed query is just $q(a)$ for the sake of brevity. In TR , we also add the line number. We do not go deep into it and more details about the syntax of the loop language are presented in Section 3.

As shown in the left hand side of Figure 1, our framework generates a query-based dependency graph to achieve the actual adaptivity A of program P . This query-based dependency graph requires two main components.

¹ One can, in principle, avoid knowing the number of queries and rounds *a priori* using a “guess-and-double” strategy, however this would weaken the bound on generalization error considerably.

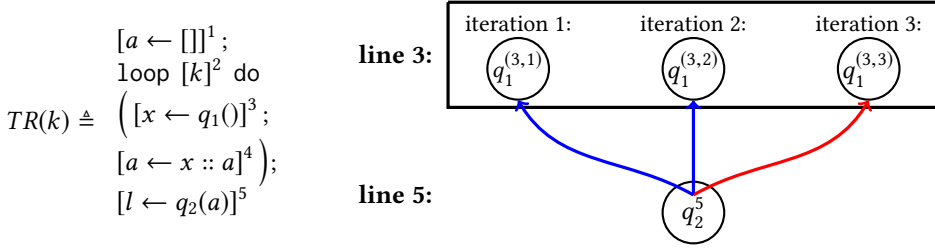


Fig. 2. A query-based dependency graph for simplified two round algorithm

- (1) A trace-based operational semantics to keep the history of the execution in record.
- (2) A formal definition of may-dependency between two query requests.

The trace-based operational semantics of the loop language have the shape of $\langle m, c, t, w \rangle \rightarrow \langle m', c', t', w' \rangle$. It works on a configuration of memory m , a program c , a trace t , a loop map w . The trace t is a list of annotated queries. One annotated query is a query with the annotation (l, w) , l is the line number and w is used for loop and tells the iteration number. For example, the query $q()$ at line 3 in the example TR at the first iteration is represented as the annotated query $q()^{(3, [2:1])}$. The loop map w is a map, from the line number of the loop counter $([k]^2)$ to its iteration number. The trace tracks the query asked during the execution, still use the TR as an example. The trace of TR starting with an empty loop map \emptyset has the following trace t_{tr} , supposing $k = 3$.

$$t_{tr} = [q()^{(3, [2:1])}, q()^{(3, [2:2])}, q()^{(3, [2:3])}, q(a)^{(5, \emptyset)}]$$

The second component is the definition of may-dependency between queries. In particular, what we need to be clear is what does it mean by saying one query may depend on another query. To this end, we look at trace. Now whether one query q depends on another query p can be checked, by witnessing q in the new generated trace upon change of p . For example, if we change the result of $q()^{(3, [2:1])}$, we have a new trace t'_{tr} , and a may change as well. In this case, $q(a)^{(5, \emptyset)}$ may not show up in the new generated trace t'_{tr} . So we think $q(a)^{(5, \emptyset)}$ may depend on $q()^{(3, [2:1])}$.

With the two components ready, a query-based dependency graph is built, in which the nodes are the annotated queries in the trace and the edge is the may dependency. We show a simple query-based graph below. The adaptivity is obtained by finding out the longest path in the graph.

Now, we look at the right hand side of Figure 1. We find the loop language is not straightforward for analysis due to the complexity of variables overwriting.

To solve this dilemma, we transform the program P to its single static assignment (short for ssa) form P . We call it ssa language. More details about the transformation as well as the ssa language is in Section ???. We briefly show our two round example in its ssa form TR^{ssa} to illustrate the ssa language.

$$TR^{ssa}(k) \triangleq \begin{aligned} & [a_1 \leftarrow []]^1; \\ & \text{loop } [k]^2 \text{ do} \\ & \quad (a_3 = \phi(a_1, a_2); \\ & \quad [x_1 \leftarrow q()]^3; \\ & \quad [a_2 \leftarrow x :: a_3]^4); \\ & [l_1 \leftarrow q(a_3)]^5 \end{aligned}$$

We notice there is one more command without label $a_3 = \phi(a_1, a_2)$ in TR^{ssa} compared to its loop variant. It adds a new variable a_3 whose value may come from a_1 or a_2 . In the syntax of our ssa

Arithmetic Expr.	$a ::= n \mid x \mid \chi \mid a \oplus_a a \mid [] \mid [a, \dots, a]$
Boolean Expr.	$b ::= \text{true} \mid \text{false} \mid \neg b \mid b \oplus_b b \mid a \sim a$
Expressions	$e ::= a \mid b$
commands	$c ::= \text{skip} \mid x \leftarrow e \mid x \leftarrow q(e) \mid \text{loop } a \text{ do } c \mid c; c \mid \text{if}(b, c, c)$

Fig. 3. Syntax of loop language

language, we use $[a_3, a_1, a_2]$ to represent this extra command in the loop command, will be discussed in Section ??.

Then we have our analysis algorithm AdaptFun to directly analyze the ssa program \mathbf{P} to get an upper bound. This algorithm constructs a variable-based dependency graph and shows the may-dependency between annotated variables. We show the graph of our two round example.

Finally, we show that our upper bound from the variable-based dependency graph is sound with respect to our adaptivity.

3 LOOP LANGUAGE

In this section, we formally introduce the language we will focus on for writing data analyses. This is a simple loop language with some primitives for calling queries. [MG: I would just call the language Loop language]. After defining the syntax of the language and showing an example, we will define its trace-based operational semantics. This is the main technical ingredient we will use to define program's adaptivity. We will conclude this section by discussing the limitation of this language with respect to static analysis for adaptivity.

3.1 Syntax

We introduce the syntax of the Loop language we will use to write our data analyses in Figure 3. It is standard that expressions can be either arithmetic expressions or boolean expressions. An arithmetic expression can be a constant n denoting integer, a variable x from some countable set Var , the special variable χ representing a row of the database, a combination of arithmetic expressions by means of the symbol \oplus_a , denoting basic operations including addition, product, subtraction, etc., the empty list $[]$, a list $[a_1, \dots, a_k]$ of arithmetic expressions. A boolean expression can be true or false, the negation of a boolean expression, or a combination of boolean expressions by means of \oplus_b , denoting basic boolean connectives, or the result of some basic comparison sym between arithmetic expressions, e.g. $\leq, =, <$, etc. A command c can either be skip, an assignment command $x \leftarrow e$, the composition of two commands $c; c$, an if statement $\text{if}(b, c, c)$, a loop statement $\text{loop } a \text{ do } c$.

The main novelty of the syntax is the query request command $x \leftarrow q(e)$. As a reminder, the aforementioned tight bound of adaptive data analysis in Section 2 focuses on linear queries. So our framework targets the linear queries as well. The linear query is specified by a function from rows to $[0, 1]$ or $[-1, +1]$. To express these functions, we introduce the special variable χ to represent the rows of the database in the arithmetic expressions. In this sense, a simple linear query which returns the first element of the row is written as $q(\chi(1))$, as a notation for a query $q(\chi) = \chi(1)$.

[MG: I don't think that this corresponds exactly to our approach. I think that we want to focus on linear queries, these are the ones for which we have the bounds. A linear query is specified by a function from rows to $[0,1]$ or $[-1,+1]$. So, in some sense, we want to have a language that describes these functions. Cannot we use $q(r) = e$ where r is a special variable denoting the given row, and e is an expression as we have right now? Example: $q_j(x) = x(i) \cdot x(j)$ can be written as $q(\chi(i) \cdot \chi(j))$ which is a notation for $q = \lambda \chi \chi(i) \cdot \chi(j)$.]

TRC(k)	TRC ^{ssa} (k)
$[a \leftarrow []]^1;$	$[a_1 \leftarrow []]^1;$
$[j \leftarrow 0]^2;$	$[j_1 \leftarrow 0]^2;$
loop $[k]^3$ do	loop $[k]^3, 0$, do $[(j_3, j_1, j_2), (a_3, a_1, a_2)]$
$\left([x \leftarrow q(\chi(j) \cdot \chi(k))]^4;$	$\left([x_1 \leftarrow q(\chi(j_3) \cdot \chi(k))]^4;$
$[j \leftarrow j + 1]^5;$	$[j_2 \leftarrow j_3 + 1]^5;$
$[a \leftarrow x :: a]^6);$	$[a_2 \leftarrow x_1 :: a_3]^6);$
$\left[l \leftarrow q(\text{sign}(\sum_{i \in [k]} \chi(i) \times \ln \frac{1+a[i]}{1-a[i]})) \right]^7$	$\left[l_1 \leftarrow q(\text{sign}(\sum_{i \in [k]} \chi(i) \times \ln \frac{1+a_3[i]}{1-a_3[i]})) \right]^7$

Fig. 4. Two round algorithm complete version

We have seen a simplified version of the two round algorithm in Section 2. We show its complete version *TRC* expressed in our loop language on the left hand side in Figure 4.

As described before, the complete version still has two steps. The query asked in the first step now depends on the iteration number so that the query ask at the j iteration is $(q_j(x) = x(j) \cdot x(k))$, expressed as $q(\chi(j) \cdot \chi(k))$. The iteration counter is initialized to 0, $j = 0$. In the second step, the final query is more complicated. It uses an auxiliary function $\text{sign}(y) = \begin{cases} +1 & \text{if } y \geq 0 \\ -1 & \text{otherwise} \end{cases}$. The input of this function is $\sum_{i \in [k]} \chi(i) \times \ln \frac{1+a[i]}{1-a[i]}$, which sums the product of $\chi(i)$ and $\ln \frac{1+a[i]}{1-a[i]}$ that uses the value at index i of the list a .

In a word, we go through a two-round adaptive analysis algorithm and shows how we can represent it in the loop language. It reveals the expressiveness of this language for most of adaptive data analysis algorithms.

3.2 Trace-based Operational semantics

In order to capture the dependency relations between queries we will evaluate programs in our Loop language by means of a trace-based operational semantics. For distinguishing among different (occurrences of) queries, our trace-based semantics will generate traces that record also the line and loop numbers for different queries. For this reason we will consider a labelled version of the Loop language. This can be described as follows: **[MG: Change "while map" to "Loop maps" everywhere.]**

Label	l	$::=$	\mathbb{N}
Loop Maps	w	\in	$\text{Label} \times \mathbb{N}$
Labelled commands	c	$::=$	$[x \leftarrow e]^l \mid [x \leftarrow q(e)]^l \mid \text{loop } [a]^l \text{ do } c \mid c; c \mid \text{if } ([b]^l, c, c) \mid [\text{skip}]^l$
Memory	m	$::=$	$[] \mid m[x \rightarrow v]$
Trace	t	$::=$	$[] \mid q(v)^{(l,w)} :: t$
Annotated Query	AQ	$::=$	$\{q(v)^{(l,w)}\}$

The commands are now labelled — we assume that labels are unique and correspond to the line of code where they appear. We denote the label l , which is the natural number standing for the line number of command in a program. Here, we put the label l to the conditional b in the if statement and to the loop counter a in the loop command, which is helpful when we present the Loop Maps w .

A memory is standard, a map from variables to values. An important component of the labelled loop language is w , what we call "Loop Maps", designed for loop specifically as an complementary

annotation to the label. With the label and Loop Maps defined, we can show annotated queries \mathcal{AQ} have unique annotations.

A trace t is supposed to accumulate along with the execution of the program, it is a list of annotated queries \mathcal{AQ} , in which the query $q(v)$ has the unique annotation (l, w) . We can understand the annotation in such a way. The label l locates arbitrary query request when we look at the execution path of a program with no loop. However, when loop repeats queries in its body, these query requests from varied iterations share the same line number. Hence, the label l is not enough to distinguish queries in loops. For this purpose, another symbol is needed to represent the iteration number of a loop. A simplified approach is to use natural number n for the iteration number so that a pair (l, n) can help, but it fails to support nested loop. This accounts for the appearance of loop maps w into the annotation (l, w) . As a map from label l to the iteration number n (a natural number), w gives accurate information on which loop specified by its label and which iteration n the query belongs to. We define some operations on w below.

$$\begin{aligned} w \setminus l &= w & l \notin \text{Keys}(w) \\ &= w_l & \text{Otherwise} \\ w + l &= w[l \rightarrow 1] & l \notin \text{Keys}(w) \\ &= w[l \rightarrow w(l) + 1] & \text{Otherwise} \end{aligned}$$

We denote $w \setminus l$ to remove the mapping of the key l in the loop maps w , used when exiting the loop at the line l . We register a label to w in the first iteration of a loop marked by label l and assign l to iteration 1, and increase its iteration number by one mapped by the label l when going into another iteration of the loop at line l . We denote $\text{Keys}(w)$ to return all the keys of the loop maps w .

A trace can be regarded as the history of query requests during the execution, generated by the trace-based small-step operational semantics. The evaluation of labelled commands are presented in Figure 5, of the form $\langle m, c, t, w \rangle \rightarrow \langle m', \text{skip}, t', w' \rangle$, stating a configuration $\langle m, c, t, w \rangle$ evaluates to another configuration with the trace and while map updated and the command c evaluated to skip. The configuration contains a memory m , a starting trace t , a starting while map w . Most of the time, the while map remains empty until the evaluation goes into the loop. We also have the evaluation of arithmetic expressions $\langle m, a \rangle \rightarrow_a a'$, evaluating an arithmetic expression a in the memory m , and similar for the boolean expressions $\langle m, b \rangle \rightarrow_b b'$.

The rule **low-query-e** evaluates the argument of a query request. After the query arguments in the normal form, the rule **low-query-v** modifies the starting memory m to $m[v_q/x]$ using the return results v_q of the query $q(v)$ request from a hidden database protected by a mechanism. We also notice that the trace expands by appending this query request $q(v)$ with the current annotation (l, w) . The rule for assignment is standard and the trace remains unchanged. The sequence rule keeps tracking the modification of the trace, and the evaluation rules for if conditional goes into one branch based on the result of the conditional b . The rules for loop are interesting. In rule **low-loop**, the while map w is updated by $w + l$ because the execution goes into another iteration when $v_N > 0$. When v_N reaches 0, the loop exits and the while map w eliminates the label l of this loop command by $w \setminus l$ in the rule **low-loop-exit**.

3.3 Query-based may-dependency graph

we give our attempt to define the number of rounds of adaptivity in our loop language, through a query-based dependency graph. A query-based dependency graph of a program consists of nodes, which are the queries during the execution of the program, and directed edges showing "may dependency" between two nodes(queries).

To this end, the definition of one query may depend on the other is fundamental to construct the query-based dependency graph. We first look at two possible candidates:

$$\begin{array}{c}
\boxed{\langle m, c, t, w \rangle \rightarrow \langle m', c', t', w' \rangle} \quad \frac{\langle m, e \rangle \rightarrow e'}{\langle m, [x \leftarrow q(e)]^l, t, w \rangle \rightarrow \langle m, [x \leftarrow q(e')]^l, t, w \rangle} \text{low-query-e} \\
\frac{q(v) = v_q}{\langle m, [x \leftarrow q(v)]^l, t, w \rangle \rightarrow \langle m[v_q/x], \text{skip}, t ++ [q(v)]^{(l, w)}, w \rangle} \text{low-query-v} \\
\frac{}{\langle m, [x \leftarrow v]^l, t, w \rangle \rightarrow \langle m[v/x], [\text{skip}]^l, t, w \rangle} \text{low-assn} \\
\frac{\langle m, c_1, t, w \rangle \rightarrow \langle m', c'_1, t', w' \rangle}{\langle m, c_1; c_2, t, w \rangle \rightarrow \langle m', c'_1; c_2, t', w' \rangle} \text{low-seq1} \quad \frac{}{\langle m, [\text{skip}]^l; c_2, t, w \rangle \rightarrow \langle m, c_2, t, w \rangle} \text{low-seq2} \\
\frac{\langle m, b \rangle \rightarrow_b b'}{\langle m, [\text{if}(b, c_1, c_2)]^l, t, w \rangle \rightarrow \langle m, [\text{if}(b', c_1, c_2)]^l, t, w \rangle} \text{low-if} \\
\frac{}{\langle m, [\text{if}(\text{true}, c_1, c_2)]^l, t, w \rangle \rightarrow \langle m, c_1, t, w \rangle} \text{low-if-t} \\
\frac{}{\langle m, [\text{if}(\text{false}, c_1, c_2)]^l, t, w \rangle \rightarrow \langle m, c_2, t, w \rangle} \text{low-if-f} \\
\frac{v_N > 0}{\langle m, [\text{loop } v_N \text{ do } c]^l, t, w \rangle \rightarrow \langle m, c; [\text{loop } (v_N - 1) \text{ do } c]^l, t, (w + l) \rangle} \text{low-loop} \\
\frac{v_N = 0}{\langle m, [\text{loop } v_N \text{ do } c]^l, t, w \rangle \rightarrow \langle m, [\text{skip}]^l, t, (w \setminus l) \rangle} \text{low-loop-exit}
\end{array}$$

Fig. 5. Trace-based Operational semantics

- (1) One query depends on the other if and only if the change of the return results of one query may also change the results of the other query.
- (2) One query depends on the other if and only if the change of the return results of one query may also change the appearance of the other query.

The first one looks reasonable to witnessing the result of one query according to the change of return result of another query. We can easily find that the two queries have nothing to do with each other in a simple example $p = x \leftarrow q(\chi(1)); y \leftarrow q(\chi(2))$. However, it does not consider that one query may not even being executed according the change of another query, as shown in program p_1 , which means we may fail to witness the change of result.

$$p_1 = x \leftarrow q(\chi(1)); \text{if}(x > 2, y \leftarrow q(\chi(2)), \text{skip})$$

We choose the second definition, witnessing the appearance of one query $q(\chi(2))$ upon the change of results of another query $q(\chi(1))$. One benefit of focusing on the appearance of queries forces us to study the effect of one query on the control flow, and then the execution. Still use $p = x \leftarrow q(\chi(1)); y \leftarrow q(\chi(2))$ as an example, we can still easily conclude that $q(\chi(1))$ has nothing to do with $q(\chi(2))$.

because the appearance of $q(\chi(2))$ is always there during the execution(or in the trace). Let us look at p_1 again, in which control flow plays its role. We think the two queries may have dependency because the change of the result of the previous query $q(\chi(1))$ affects the control flow, hence the appearance of $q(\chi(2))$ in the corresponding trace.

Another benefit of the appearance definition is it also takes the arguments of the queries into account. Let us look at another variant of program p , marked as p_2 , in which the queries equipped with functions using previous assigned variables.

$$p_2 = x \leftarrow q(\chi(2)); y \leftarrow q(x + \chi(3))$$

As a reminder, in the loop language, the query request is composed by two components: a symbol q representing a linear query type and the argument e , which represents the function specifying what the query asks. So we do think $q(\chi(1))$ is different from $q(\chi(2))$. Informally, we think $q(x + \chi(3))$ may depend on the query $q(\chi(2))$, because equipped function of the former $x + \chi(3)$ is affected by the latter $q(\chi(2))$. The appearance definition catches this case, since $q(x + \chi(2))$ will be different queries if x is changed along with the change of the $q(\chi(2))$.

We give a formal definition of query may dependency based on the trace-based operational semantics as follows.

DEFINITION 1 (QUERY MAY DEPENDENCY). *One query $q(v_1)$ may depend on another query $q(v_2)$ in a program c , with a starting while map w , denoted as $\text{DEP}(q(v_1)^{(l_1, w_1)}, q(v_2)^{(l_2, w_2)}, c, w, m, D)$ is defined as below.*

$$\begin{aligned} & \forall t. \exists m_1, m_3, t_1, t_3. \langle m, c, t, w \rangle \rightarrow^* \langle m_1, [x \leftarrow q(v_1)]^{l_1}, c_2, t_1, w_1 \rangle \rightarrow \\ & \langle m_1[q(v_1)(D)/x], c_2, t_1 + [q(v_1)^{(l_1, w_1)}], w_1 \rangle \rightarrow^* \langle m_3, \text{skip}, t_3, w_3 \rangle \\ & \wedge (q(v_1)^{(l_1, w_1)} \in (t_3 - t) \wedge q(v_2)^{(l_2, w_2)} \in (t_3 - t) \implies \exists v \in \text{codom}(q(v_1)), m'_3, t'_3, w'_3. \\ & \langle m_1[v/x], c_2, t_1 + [q(v_1)^{(l_1, w_1)}], w_1 \rangle \rightarrow^* \langle m'_3, \text{skip}, t'_3, w'_3 \rangle \wedge (q(v_2)^{(l_2, w_2)} \notin (t'_3 - t)) \\ & \wedge (q(v_1)^{(l_1, w_1)} \in (t_3 - t) \wedge q(v_2)^{(l_2, w_2)} \notin (t_3 - t) \implies \exists v \in \text{codom}(q(v_1)), m'_3, t'_3, w'_3. \\ & \langle m_1[v/x], c_2, t_1 + [q(v_1)^{(l_1, w_1)}], w_1 \rangle \rightarrow^* \langle m'_3, \text{skip}, t'_3, w'_3 \rangle \wedge (q(v_2)^{(l_2, w_2)} \in (t'_3 - t)) \end{aligned}$$

We give a formal definition of the query-based dependency graph with the formal definition of may-dependency between two queries.

DEFINITION 2 (QUERY-BASED DEPENDENCY GRAPH). *Given a program c , a database D , a starting memory m , an initial while map w , the query-based dependency graph $G(c, D, m, w) = (V, E)$ is defined as:*

$$\begin{aligned} V &= \{q(v)^{l, w} \mid \forall t. \exists m', w', t'. \langle m, c, t, w \rangle \rightarrow^* \langle m', \text{skip}, t', w' \rangle \wedge q(v)^{l, w} \in (t' - t)\}. \\ E &= \{(q(v)^{(l, w)}, q(v')^{(l', w')}) \in \mathcal{AQ} \times \mathcal{AQ} \mid \text{DEP}(q(v)^{(l, w)}, q(v')^{(l', w')}, c, w, m, D) \wedge \text{To}(q(v')^{(l', w')}, q(v)^{(l, w)})\}. \end{aligned}$$

The function $\text{To}(q(v')^{(l', w')}, q(v)^{(l, w)})$ tells that the query request $q(v')^{(l', w')}$ appears after the query request $q(v)^{(l, w)}$ in the trace, by comparing the annotation (l', w') and (l, w) . It helps to decide on the direction of one edge.

The query-based dependency graph only considers the newly generated annotated queries during the execution of the program c , so we see the nodes comes from the trace $t' - t$. The previous trace before the execution of c is excluded when constructing the graph. To summary, for every execution of a program c staring with different starting configuration, we can construct a corresponding dependency graph.

Finally, we reach the definition of what we want – the number of rounds of adaptivity, by means of a query-based dependency graph.

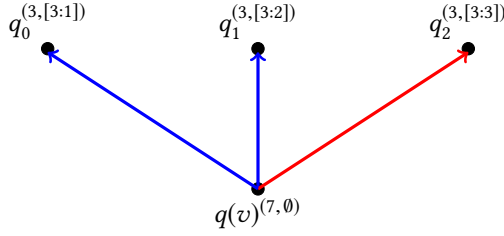


Fig. 6. A query-based dependency graph for two round algorithm complete version

DEFINITION 3 (ADAPTIVITY). Given a program c , and a meory m , a database D , a starting while map w , the adaptivity of the dependency graph $G(c, D, m, w) = (V, E)$ is the length of the longest path in this graph. We denote the path from $q(v)^{(l,w)}$ to $q(v')^{(l',w')}$ as $p(q(v)^{(l,w)}, q(v')^{(l',w')})$. The adaptivity denoted as $A^*(c, D, m, w)$.

$$A(c, D, m, w) = \max_{q(v)^{(l,w)}, q(v')^{(l',w')} \in V} \{|p(q(v)^{(l,w)}, q(v')^{(l',w')})|\}$$

We give an example to illustrate the aforementioned contents to collect the trace, build the dependency graph and get the adaptivity, via our familiar adaptive data analysis algorithm - two round strategy.

Given a specific database $D = [[1, 1], [0, 0], [1, 1], [1, 1]]$, the execution trace is generated along with the operational semantics as follows:

$$t = \{q(\chi(0) \cdot \chi(3))^{(4,[3:1])}, q(\chi(1) \cdot \chi(3))^{(4,[3:2])}, q(\chi(2) \cdot \chi(3))^{(4,[3:3])}, q(v)^{(7,0)}\}$$

For the sake of brevity, we use $q(v)$ to represent $q(\text{sign}(\sum_{i \in [k]} \chi(i) \times \ln \frac{1+a[i]}{1-a[i]}))$ where $a = [1.0.1]$. We use q_0, q_1, q_2 to represent $q(\chi(0) \cdot \chi(3)), q(\chi(1) \cdot \chi(3)), q(\chi(2) \cdot \chi(3))$. Then we have the dependency graph generated in Figure 6.

We can notice the complete version generates the similar query-based dependency graph as its simplified version in Figure 2.

4 THE ANALYSIS ALGORITHM ON SSA PROGRAMS

In this section, we clarify the algorithm AdaptFun that analyzes the adaptivity of a target program in ssa language, which consists of three auxiliary algorithms: AG and AD for generating a weighted variable-based dependency graph and AP to find the most weighted path in the graph. We do not show details of AP, it is quite standard path-finding algorithm in a weighted graph. We go through the two round example to illustrate the algorithm. Finally, we show the soundness of our algorithm with respect to the adaptivity.

4.1 The ideas behind the algorithm

In consideration of the definition of adaptivity from the query-based dependency graph, our analysis whose goal is a sound upper bound on the adaptivity, is supposed to take care of paths(possible adaptivity candidates) in all the possible dependency graphs (per configuration). To this end, our algorithm aims to syntactically construct a dependency graph, in which the nodes are annotated ssa variables and the directed edges showing one annotated variable may depend on the other if there exists an edge between them. Intuitively, query requests in the query-based dependency graph is assigned to variables that appears in the predicted ssa-variable-based dependency graph.

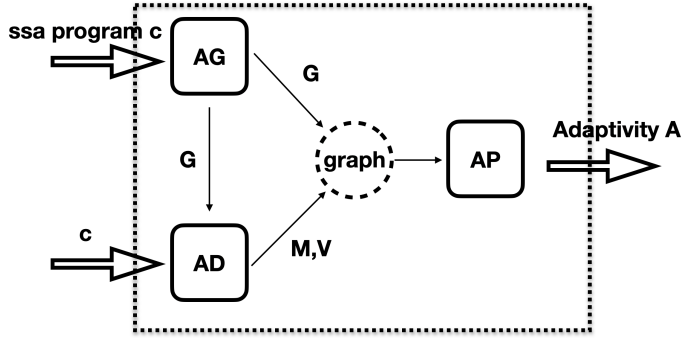


Fig. 7. The overview of AdaptFun

The algorithm AdaptFun estimates the adaptivity from the weighted variable-based dependency graph, whose structure is shown in Figure 7. We have a look at the structure of AdaptFun. We start with the dependency graph, represented in the form of a matrix M . To know which variable is associated with a query request in the matrix, an extra vector V is produced in the algorithm, of the size of all estimated assigned variables in the target program c . Naturally, the first question comes to us, what is the size of the matrix M and V , or the estimated assigned variables? To solve this, the analysis goes through the ssa program c twice, estimating the assigned variables to be tracked and adding the appropriate annotation (l, w) , in the first time scan. Those annotated variables are collected in a global list G , whose size determines the size of the matrix and vector. This first scan generating G is conducted by the algorithm AG in Figure 7, which takes the target program c as input. The global list G is fed into the second scan, which decides the size of the matrix and vector. For example, if G of size N , then the matrix used in the second round is of size $N \times N$, and vector of size N . It maintains a unique mapping from variables in G to the matrix M and the vector V . For example, the i th row, j th column of the matrix $M[i][j] = 1$ represents the may-dependency from variable $G[i]$ to $G[j]$, $M[i][j] = 0$ means no dependency. In a similar way, $V[i] = 1$ means the variable $G[i]$ is assigned to a query request. The second scan of c , implemented by the algorithm AD, then records the may-dependency in M , the weight is tracked in V . We think variable associated with a query request of weight 1 in the graph and others of weight 0. After the analysis of AG and AD, a ssa-variable-based weighted dependency graph is constructed. We use another auxiliary algorithm AP to find the path with the most weights in the graph. Then the weight is the adaptivity A AdaptFun estimates.

4.2 Variable Estimation algorithm

We first show how G will be collected, through a variable estimating algorithm AG of the form $AG(G; w; c) \rightarrow (G'; w')$ in Figure 8. The input of AG is a list of estimated annotated variables G collected before the program c , a while map w consistent with previous estimation, a program c . The output of the algorithm is the updated global list G' , along with the updated while map w for later estimation. The assignment is the source of variables AG estimates, in the case **ag-asgn** and **ag-query**, the output global list is extended by $x^{(l, w)}$. When it comes to the if command in the rule **ag-if**, variables assigned in the then branch c_1 , as well as the variables assigned in the else branch c_2 , and the new generated variables $\bar{x}, \bar{y}, \bar{z}$ in $[\bar{x}, \bar{x}_1, \bar{x}_2], [\bar{y}, \bar{y}_1, \bar{y}_2], [\bar{z}, \bar{z}_1, \bar{z}_2]$. The sequence is standard by accumulating the predicted variables in the two commands c_1 and c_2 in a

$$\begin{array}{c}
\text{589} \quad \overline{\text{AG}(G; w; [\mathbf{x} \leftarrow \mathbf{e}]^l) \rightarrow (G + +[\mathbf{x}^{(l,w)}]; w)} \quad \text{ag-asgn} \\
\text{590} \\
\text{591} \quad \overline{\text{AG}(G; w; [\mathbf{x} \leftarrow q(\mathbf{e})]^l) \rightarrow (G + +[\mathbf{x}^{(l,w)}]; w)} \quad \text{ag-query} \\
\text{592} \\
\text{593} \quad \overline{\begin{array}{c} \text{AG}(G; w; \mathbf{c}_1) \rightarrow (G_1; w_1) \quad \text{AG}(G_1; w; \mathbf{c}_2) \rightarrow (G_2; w_2) \\ G_3 = G_2 + +[\bar{\mathbf{x}}^{(l,w)}] + +[\bar{\mathbf{y}}^{(l,w)}] + +[\bar{\mathbf{z}}^{(l,w)}] \end{array}} \quad \text{ag-if} \\
\text{594} \\
\text{595} \quad \overline{\text{AG}(G; w; [\text{if}(\mathbf{b}, [\bar{\mathbf{x}}, \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2], [\bar{\mathbf{y}}, \bar{\mathbf{y}}_1, \bar{\mathbf{y}}_2], [\bar{\mathbf{z}}, \bar{\mathbf{z}}_1, \bar{\mathbf{z}}_2], \mathbf{c}_1, \mathbf{c}_2)]^l) \rightarrow (G_3; w)} \\
\text{596} \\
\text{597} \quad \overline{\begin{array}{c} \text{AG}(G; w; \mathbf{c}_1) \rightarrow (G_1; w_1) \quad \text{AG}(G_1; w_1; \mathbf{c}_2) \rightarrow (G_2; w_2) \\ \text{AG}(G; w; (\mathbf{c}_1; \mathbf{c}_2)) \rightarrow (G_2; w_2) \end{array}} \quad \text{ag-seq} \\
\text{598} \\
\text{599} \quad \overline{\begin{array}{c} G_0 = G \quad w_0 = w \quad \forall 0 \leq z < N. \text{AG}(G_z + +[\bar{\mathbf{x}}^{(l, w_z+1)}]; (w_z + l); \mathbf{c}) \rightarrow (G_{z+1}; w_{z+1}) \\ G_f = G_N + +[\bar{\mathbf{x}}^{(l, w_N \setminus l)}] \quad \mathbf{a} = N \end{array}} \quad \text{ag-loop} \\
\text{600} \\
\text{601} \quad \overline{\text{AG}(G; w; [\text{loop } \mathbf{a}, n, [\bar{\mathbf{x}}, \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2] \text{ do } \mathbf{c}]^l) \rightarrow (G_f; w_N \setminus l)} \\
\text{602} \\
\text{603} \\
\text{604} \\
\text{605} \\
\text{606} \\
\text{607} \\
\text{608} \\
\text{609}
\end{array}$$

Fig. 8. The algorithm AG

sequence $\mathbf{c}_1; \mathbf{c}_2$. The loop considers the loop iterations as well by assuming the loop counter \mathbf{a} to be certain natural number N in the rule **ag-loop**. The algorithm counts the assigned variables in every iteration, including those new assigned variables in $\bar{\mathbf{x}}$, those variables assigned in the body \mathbf{c} , with the appropriate annotation showing the iteration number of the variables.

4.3 Matrix and Vector based algorithm

We have a global list G after the first scan of the ssa programs \mathbf{c} . We develop a matrix and vector based approach based on the global list G to get an estimated upper bound on the adaptivity of the program. In this approach, a matrix is constructed according to those estimated annotated variables from the global list, in which every row and every column corresponds to the unique variable in G by position. To be precise, in this matrix, 0 means no dependency while non-zero value in the matrix shows may-dependency between corresponding variables. If the value of $M[i][j]$ in the matrix M is greater than zero, we know annotated variable $G[i]$ may depend on $G[j]$. A vector V has the same size as G and it records whether the corresponding variable $G[i]$ is assigned with a query request. $V[i] = 1$ means assigned with a query request and otherwise $V[i] = 0$. We call this algorithm which tracks may-dependency in the ssa programs and query information, AD. The algorithm of the form $\text{AD}(\Gamma; \mathbf{c}; i) \rightarrow (M; V; i')$, the input is a tuple consisting of three elements: (1) a one-row- N -column matrix Γ storing the dependency from previous program, it is used when handling the if command. (2) the ssa program \mathbf{c} to be analyzed (3) an index i specifying the location of the first assigned variable of the program \mathbf{c} in the global list G . The output of AD also consists of three elements; (1) A matrix M showing the may-dependency in \mathbf{c} (2) A vector V records the query requests in \mathbf{c} (3) the index i' that refers to the next position of the last assigned variable in \mathbf{c} , if exists. The existence of the index i' helps to locate the first assigned variable if we need to analyze programs after \mathbf{c} .

We first define some functions which use the indices in G . The function $L(i)$ generates one-column- N -rows matrix, where only the i -th row is 1 and all the other rows are 0. This function is used to locate the right row when calculate the matrix when analyzing assignment and query.

$$\begin{array}{c}
\text{638} \quad \frac{M = L(i) * (R(e, i) + \Gamma)}{\text{AD}(\Gamma; [\mathbf{x} \leftarrow \mathbf{e}]^l; i) \rightarrow (M; V_0; i + 1)} \quad \text{ad-assign} \quad \frac{M = L(i) * (R(e, i) + \Gamma) \quad V = L(i)}{\text{AD}(\Gamma; [\mathbf{x} \leftarrow q(\mathbf{e})]^l; i) \rightarrow (M; V; i + 1)} \quad \text{ad-query} \\
\text{639} \\
\text{640} \\
\text{641} \\
\text{642} \quad \frac{\begin{array}{c} \text{AD}(\Gamma + R(\mathbf{b}, i_1); \mathbf{c}_1; i_1) \rightarrow (M_1; V_1; i_2) \quad \text{AD}(\Gamma + R(\mathbf{b}, i_1); \mathbf{c}_2; i_2) \rightarrow (M_2; V_2; i_3) \\ \text{AD}(\Gamma; [\bar{\mathbf{x}}, \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2]; i_3) \rightarrow (M_x; V_0; i_3 + |\bar{\mathbf{x}}|) \quad \text{AD}(\Gamma; [\bar{\mathbf{y}}, \bar{\mathbf{y}}_1, \bar{\mathbf{y}}_2]; i_3 + |\bar{\mathbf{x}}|) \rightarrow (M_y; V_0; i_3 + |\bar{\mathbf{x}}| + |\bar{\mathbf{y}}|) \\ \text{AD}(\Gamma; [\bar{\mathbf{z}}, \bar{\mathbf{z}}_1, \bar{\mathbf{z}}_2]; i_3 + |\bar{\mathbf{x}}| + |\bar{\mathbf{y}}|) \rightarrow (M_z; V_0; i_3 + |\bar{\mathbf{x}}| + |\bar{\mathbf{y}}| + |\bar{\mathbf{z}}|) \\ M = (M_1 + M_2) + M_x + M_y + M_z \end{array}}{\Gamma \vdash_{M, V_1 \uplus V_2}^{(i_1, i_3 + |\bar{\mathbf{x}}| + |\bar{\mathbf{y}}| + |\bar{\mathbf{z}}|)} [\text{if}(\mathbf{b}, [\bar{\mathbf{x}}, \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2], [\bar{\mathbf{y}}, \bar{\mathbf{y}}_1, \bar{\mathbf{y}}_2], [\bar{\mathbf{z}}, \bar{\mathbf{z}}_1, \bar{\mathbf{z}}_2], \mathbf{c}_1, \mathbf{c}_2)]^l} \quad \text{ad-if} \\
\text{643} \\
\text{644} \\
\text{645} \\
\text{646} \quad \frac{\text{AD}(\Gamma; \mathbf{c}_1; i_1) \rightarrow (M_1; V_1; i_2) \quad \text{AD}(\Gamma; \mathbf{c}_2; i_2) \rightarrow (M_2; V_2; i_3)}{\text{AD}(\Gamma; (\mathbf{c}_1; \mathbf{c}_2); i_1) \rightarrow ((M_1; M_2); V_1 \uplus V_2; i_3)} \quad \text{ad-seq} \\
\text{647} \\
\text{648} \\
\text{649} \\
\text{650} \\
\text{651} \quad \frac{\begin{array}{c} B = |\bar{\mathbf{x}}| \quad A = |\mathbf{c}| \\ \forall 0 \leq j < N. \text{AD}(\Gamma; [\bar{\mathbf{x}}, \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2]; i + j * (B + A)) \rightarrow (M_{1j}; V_{1j}; i + B + j * (B + A)) \\ \text{AD}(\Gamma; \mathbf{c}; i + B + j * (B + A)) \rightarrow (M_{2j}; V_{2j}; i + B + A + j * (B + A)) \\ \text{AD}(\Gamma; [\bar{\mathbf{x}}, \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2]; i + N * (B + A)) \rightarrow (M; V; i + N * (B + A) + B) \\ \mathbf{a} = N \quad M' = M + \sum_{0 \leq j < N} (M_{1j} + M_{2j}) \quad V' = V \uplus \sum_{0 \leq j < N} (V_{1j} \uplus V_{2j}) \end{array}}{\Gamma \vdash_{M', V'}^{(i, i + N * (B + A) + B)} [\text{loop } \mathbf{a}, 0, [\bar{\mathbf{x}}, \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2] \text{ do } \mathbf{c}]^l} \quad \text{ad-loop} \\
\text{652} \\
\text{653} \\
\text{654} \\
\text{655} \\
\text{656} \\
\text{657} \\
\text{658} \\
\text{659} \\
\text{660} \\
\text{661} \\
\text{662}
\end{array}$$

Fig. 9. The algorithm AD

The function $R(e, i)$ generates a one-row- N -column matrix. For every variable used in e , it finds the corresponding index i in G so that $G[i]$ maps to the variable and mark the i th column as 1. If it is not found, we do not mark. When we say $G[i]$ maps to a target variable, we take off the annotation of $G[i]$ and check if the left variable is the same as the target variable. To handle loop, for instance, a variable y appears many times in G , but with different annotations (iteration numbers), the argument i helps to find most recent assignment variable y before the index i in G . It is still used when analyzing assignment and query. Thanks to our ssa language, our choice of the most recent assigned variable is reasonable because the variable used in the loop refers to the most recent assignment and every variable is uniquely assigned in its ssa form.

We define $M_1; M_2$ to combine two matrix, where $M_1 + M_2$ is the standard sum of two matrix.

$$M_1; M_2 := M_2 \cdot M_1 + M_1 + M_2$$

And define the operator \uplus to combine two vectors.

$$V_1 \uplus V_2 := \begin{cases} 1 & (V_1[i] = 1 \vee V_2[i] = 1) \wedge i = 1, \dots, N \wedge |V_1| = |V_2| \\ 0 & \text{o.w.} \end{cases}$$

For the sake of brevity, we also define some annotations as follows. We show how the algorithm AD handles the extra part $[\bar{\mathbf{x}}, \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2]$ in the if and loop commands. First, we give a unique name for variables in lists $\bar{\mathbf{x}}, \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2$ respectively, as follows: $\forall 0 \leq z < |\bar{\mathbf{x}}|. \bar{\mathbf{x}}(z) = \mathbf{x}_z, \bar{\mathbf{x}}_1(z) = \mathbf{x}_{1z}, \bar{\mathbf{x}}_2(z) = \mathbf{x}_{2z}$. And then we treat every tuple $(\mathbf{x}_z, \mathbf{x}_{1z}, \mathbf{x}_{2z})$ in $[\bar{\mathbf{x}}, \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2]$ as the simple may dependency case: \mathbf{x}_z may depend on both \mathbf{x}_{1z} and \mathbf{x}_{2z} , just like $\mathbf{x}_z \leftarrow \mathbf{x}_{1z} + \mathbf{x}_{2z}$, defined as follows. $\text{AD}(\Gamma; [\bar{\mathbf{x}}, \bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2]; i) \rightarrow (M; V_0; i + |\bar{\mathbf{x}}|) \triangleq \forall 0 \leq z < |\bar{\mathbf{x}}|. \text{AD}(\Gamma; \mathbf{x}_z \leftarrow \mathbf{x}_{1z} + \mathbf{x}_{2z}; i + z) \rightarrow (M_{x_z}; V_0; i + z + 1)$ where $M = \sum_{0 \leq z < |\bar{\mathbf{x}}|} M_{x_z}$. One of the key idea under algorithm AD is to track the indices i, i' both in the input and output

to synchronize with its previous algorithm AG. The index in AD increases as the same way as the global list expands after the analysis of a program c , which helps AD record the dependency relation from the program c in the right place of the matrix. For example, in the case **ad-asgn** and **ad-query**, the index increases by 1, which corresponds to their counterparts of algorithm **AG**. The if and loop commands have the extra part $[\bar{x}, \bar{x}_1, \bar{x}_2]$ and we find that the output index increases by also considering this part as we do in collecting the global list.

Another interesting point is the construction of the matrix. The fundamental case is the assignment and query cases. We use a function $L(i)$ to generate a N-row-one-column matrix L to guarantee the resulting matrix only has non-zeros at row i . The intuition behind is that one single assignment or query request can only reveal the dependency of its assigned variable (corresponding to one row of the matrix) to the variables used on the right hand sides. Thanks to the index i , we know which row this assignment should be in the matrix. The function $R(e, i)$ gets a one-row-N-column matrix marking the variables used in the right hand side. The Γ is designed for the if command and we will discuss it later. We can see one simple example sa to get a taste.

$$sa \triangleq \begin{cases} [x_1 \leftarrow 2]^1; \\ [x_2 \leftarrow x_1 + 2]^2; \\ [x_3 \leftarrow x_1 + x_2]^3 \end{cases}$$

In the program sa , only simple assignment is involved. When we assume Γ is empty, for the assignment at line 3, the matrix is built as follows.

$$\text{line3: } [x_3 \leftarrow x_1 + x_2]^3 : \begin{matrix} x_1 \\ x_2 \\ x_3 \end{matrix} \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} * \begin{bmatrix} x_1 & x_2 & x_3 \\ 1 & 1 & 0 \end{bmatrix} = \begin{matrix} x_1 \\ x_2 \\ x_3 \end{matrix} \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 1 & 0 \end{bmatrix}$$

We use a one-row-N-column matrix Γ as one of the input of AD, to handle the cases when the control flow diverges in the labelled if command $[\text{if}(\mathbf{b}, [\bar{x}, \bar{x}_1, \bar{x}_2], [\bar{y}, \bar{y}_1, \bar{y}_2], [\bar{z}, \bar{z}_1, \bar{z}_2], \mathbf{c}_1, \mathbf{c}_2)]^l$, where the execution of either branch may depend on the conditional guard \mathbf{b} . Follow this intuition, the analysis of either branch is supposed to consider the variables used in the conditional \mathbf{b} , tracked in Γ . In the case **ad-if**, we can see the analysis of the two branches \mathbf{c}_1 and \mathbf{c}_2 share the same input $\Gamma + R(\mathbf{b}, i)$ and $R(\mathbf{b}, i)$ tells the variables assigned before the if command and used in the conditional \mathbf{b} .

We compose the matrix and vectors in the case of sequence in **ad-seq**. The non-zeros or we call it effect range of the matrix and vector is decided by its input and output indices. From the case **ad-seq**, two programs \mathbf{c}_1 and \mathbf{c}_2 have disjoint effect ranges $[i_1, i_2)$ and $[i_1, i_3)$, it is safe to combine them without lose information.

The loop case is handled in a well organized way. We use $B = |\bar{x}|$ and $A = |c|$ to estimate the size of variables assigned in \bar{x} and c . And $|c|$ is defined by the help of the algorithm AG, defined as $|c| = |G|$ when $\text{AG}([]; \mathbf{c}; \emptyset) \rightarrow (G; \emptyset)$. The algorithm then gets how many iterations N the loop may executes from the loop counter \mathbf{a} . For every iteration, it first records the dependency relations between variables in $[\bar{x}, \bar{x}_1, \bar{x}_2]$ by constructing a corresponding matrix M_{1j} (j is the iteration number) and an empty vector V_{1j} , and analyze the loop body c with a resulting matrix M_{2j} and vector V_{2j} . We give an extra analysis of those new assigned variables as what AG does, it works well when the loop is executed ($N = 0$) or not. We know that for all the possible iteration number j , M_{1j} and M_{2j} have disjoint effect ranges so we combine them, similar as the vectors V_{1j} and V_{2j} .

Finally, we are able to construct a variable-based weighted dependency graph based on G, M and V generated by the framework. The definition of the estimated adaptivity is the weight of the most weighted path in the graph defined as follows.

DEFINITION 4 (ESTIMATED ADAPTIVITY). Given a program c , the global list G , and $AD(\Gamma; c; i_1) \rightarrow (M, V, i_2)$, the weighted dependency graph $G_{ssa}(M, V, G, i_1, i_2) = (Nodes, Edges, Weights)$ is defined as:

$Nodes \ Vt = \{G(j) \in \mathcal{LV} \mid i_1 \leq j < i_2\}$

$Edges \ E = \{(G(j_1), G(j_2)) \in \mathcal{LV} \times \mathcal{LV} \mid M[j_1][j_2] \geq 1 \wedge i_1 \leq j_1, j_2 < i_2\}$

$Weights \ Wt = \{(G(j), 1) \in \mathcal{LV} \times \mathcal{N} \mid i_1 \leq j < i_2 \wedge V[j] = 1\} \cup \{(G(j), 0) \in \mathcal{LV} \times \mathcal{N} \mid i_1 \leq j < i_2 \wedge V[j] = 0\}$.

Adaptivity of the program defined according to the graph is as:

$$Adapt(M, V, i_1, i_2) := \max_{v_{t1}, v_{t2} \in Vt} \{Weight(p(v_{t1}, v_{t2}), Wt)\},$$

where $p(k, l)$ is the path in graph $G_{ssa}(M, V, i_1, i_2)$ starting from k to l , $Weight(p(v_{t1}, v_{t2}), Wt)$ get the total sum of weights along the path $p(v_{t1}, v_{t2})$.

4.4 Analysis on two round algorithm

We show how AdaptFun analyze the two round algorithm. For the sake of brevity, we conduct the analysis on the simplified two round algorithm TR^{ssa} .

AdaptFun first runs the algorithm AG to generate the global list G . We assume the input $k = 2$ and have the following.

$$G_{k=2} = \left[a_1^{(1,0)}, a_3^{(2,[2:1])}, x_1^{(3,[2:1])}, a_2^{(4,[2:1])}, a_3^{(2,[2:2])}, x_1^{(3,[2:2])}, a_2^{(4,[2:2])}, a_3^{(2,0)}, l_1^{(5,0)} \right]$$

We denote a_1^1 short for $a_1^{(1,0)}$ and $a_3^{(2,1)}$ short for $a_3^{(2,[2:1])}$. Then the resulting matrix M_{tr} and V_{tr} of the algorithm AD as follows.

$$M_{tr} = \begin{bmatrix} a_1^1 & a_1^1 & a_3^{(2,1)} & x_1^{(3,1)} & a_2^{(4,1)} & a_3^{(2,2)} & x_1^{(3,2)} & a_2^{(4,2)} & a_3^2 & l_1^5 \\ a_3^{(2,1)} & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ x_1^{(3,1)} & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ a_2^{(4,1)} & 0 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ a_3^{(2,2)} & 1 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ x_1^{(3,2)} & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ a_2^{(4,2)} & 0 & 0 & 0 & 0 & 1 & 1 & 0 & 0 & 0 \\ a_3^2 & 1 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ l_1^5 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \end{bmatrix}, V_{tr} = \begin{bmatrix} a_1^1 & 0 \\ a_3^{(2,1)} & 0 \\ x_1^{(3,1)} & 1 \\ a_2^{(4,1)} & 0 \\ a_3^{(2,2)} & 0 \\ x_1^{(3,2)} & 1 \\ a_2^{(4,2)} & 0 \\ a_3^2 & 0 \\ l_1^5 & 1 \end{bmatrix}$$

4.5 Soundness of AdaptFun

We would like to show that the query-based dependency graph generated from the trace of the execution of the target ssa program is a subgraph of the variable-based dependency graph predicted from our algorithm AdaptFun, and the query requested during the execution is also bounded by an estimation from our algorithm.

We first give a definition of subgraph of a query-based dependency graph with respect to a variable-based dependency graph.

DEFINITION 5 (SUBGRAPH). Given a query-based dependency graph $G_s = (V_1, E_1)$, a variable-based dependency graph $G_{ssa} = (V_2, E_2)$, $G_s \subseteq G_{ssa}$ iff:

$\exists f$ so that

1. for every $v \in V_1$, $f(v) \in V_2$.

2. $\forall e = (v_i, v_j) \in E_1$, there exists a path from $f(v_i)$ to $f(v_j)$ in G_{ssa} .

Then we show the soundness of AdaptFun. In the theorem, we use some definition. $G \models M, V$ says that G and M, V have the corresponding size. $G; w \models (c, i_1, i_2)$ checks if the variables assigned in c calculated by AG matches the variables in G from index i_1 to i_2 .

THEOREM 4.1 (SOUNDNESS OF ADAPTFUN). *Given $AD(\Gamma; c; i_1) \rightarrow (M; V; i_2)$, for any global list G , loop maps w such that $G; w \models (c, i_1, i_2) \wedge G \models (M, V)$. K is the number of queries inquired during the execution of the piece of program c and $|V|$ gives the number of non-zeros in V . Then,*

$$K \leq |V| \wedge \forall D, \mathbf{m}. G_s(c, D, \mathbf{m}, w) \subseteq G_{ssa}(M, V, G, i_1, i_2)$$

5 MORE EXAMPLES

The two round strategy works well in our framework, we explore further to look at an advanced adaptive data analysis algorithm - multiple round algorithm.

Algorithm 1 A multi-round analyst strategy for random data (Algorithm 5 in ...)

Example 5.1 (Multiple Round Algorithm). **Require:** Mechanism M with a hidden state $X \in [N]^n$ sampled u.a.r., control set size c

Define control dataset $C = \{0, 1, \dots, c-1\}$

Initialize $Nscore(i) = 0$ for $i \in [N]$, $I = \emptyset$ and $Cscore(C(i)) = 0$ for $i \in [c]$

for $j \in [k]$ **do**

let $p = \text{uniform}(0, 1)$

define $q(x) = \text{bernoulli}(p)$.

define $qc(x) = \text{bernoulli}(p)$.

let $a = M(q)$

for $i \in [N]$ **do**

$Nscore(i) = Nscore(i) + (a - p) * (q(i) - p)$ if $i \notin I$

for $i \in [c]$ **do**

$Cscore(C(i)) = Cscore(C(i)) + (a - p) * (qc(i) - p)$

let $I = \{i | i \in [N] \wedge Nscore(i) > \max(Cscore)\}$

let $X = X \setminus I$

return X .

$$\begin{aligned}
 & \begin{aligned}
 & [I \leftarrow []]^1; \\
 & [ns \leftarrow 0]^2; \\
 & [cs \leftarrow 0]^3; \\
 & \text{loop } [k]^4 \text{ do} \\
 & \left([p \leftarrow c]^5; \right. \\
 & [a \leftarrow q(f(p, I))]^6; \\
 & [ns \leftarrow \text{update_nscore}(p, a)]^7; \\
 & [cs \leftarrow \text{update_cscore}(p, a)]^8; \\
 & \left. [I \leftarrow \text{update}(I, ns, cs)]^9 \right)
 \end{aligned}
 & \Rightarrow MR^{ssa}(k) \triangleq \begin{aligned}
 & [I_1 \leftarrow []]^1; \\
 & \text{loop } [k]^2, 0, [I_3, I_1, I_2] \\
 & \text{do} \\
 & \left([p_1 \leftarrow c]^3; \right. \\
 & [a_1 \leftarrow q(p, I_3)]^4; \\
 & \left. [I_2 \leftarrow \text{update}(I_3, (a_1, p_1))]^5 \right)
 \end{aligned}
 \end{aligned}$$

Description: The multiple round algorithm starts from an initialized empty tracking list I , a score called $Nscore$ $ns = 0$, another score $Cscore$ $cs = 0$. There is a hidden database X as well. It goes k rounds and every round, the two scores ns and cs are updated by a query result. Then the list I is updated by the two scores for every round. After the r rounds, the algorithm returns the columns of the hidden database X not specified in the tracking list I , which is $X \setminus I$.

The algorithm is written in the loop language as MR . It uses a loop for the k rounds computation and. We use functions $\text{update_nscore}(p, a), \text{update_cscore}(p, a), \text{update}(I, ns, cs)$ to simplify the complex update computation of $Nscore$, $Cscore$ and the tracking list I . It will not change our analysis



Fig. 10. A query-based dependency graph for multi round algorithm

because these functions provides enough information through their arguments. In comparison with the two round algorithm, the query asked in each iteration is not independent in the multiple round one any more. The query in one iteration j now depends on the tracking list I from its previous iteration $j - 1$, which is updated by the query result at the same iteration $j - 1$. We can easily see the connection between queries from different iterations.

In *MR*, the tracking list I is initialized to an empty list. It appears inside the function of query $q(f(p, I))$ and updated in each iteration. We choose $k = 3$, and assume the tracking list becomes I_a, I_b, I_c after 3 rounds, respectively. the execution trace is generated as:

$$t_{mr} = \{q(p, I_a)^{4,[2:1]}, q(p, I_b)^{4,[2:2]}, q(p, I_c)^{4,[2:3]}\}$$

Then we have the dependency graph generated as following, with the graph shown in Figure 10.

$$V = \{q(p, I_a)^{(6,[4:1])}, q(p, I_b)^{(6,[4:2])}, q(p, I_c)^{(6,[4:3])}\}$$

$$E = \{(q(p, I_b)^{(6,[4:2])}, q(p, I_a)^{(6,[4:1])}), (q(p, I_c)^{(6,[4:3])}, q(p, I_b)^{(6,[4:2])}), (q(p, I_c)^{(6,[4:3])}, q(p, I_a)^{(6,[4:1])}), \}$$

Then we have the adaptivity calculated from the graph as follows:

$$\begin{aligned} A^*(c, D, m, w) &= \max_{q(v)^{(l,w)}, q(v')^{(l',w')} \in V} \{|p(q(v))^{(l,w)}, q(v')^{(l',w')}|\} \\ &= \{|(q(p, I_c)^{(6,[4:3])} \rightarrow q(p, I_b)^{(6,[4:2])} \rightarrow q(p, I_a)^{(6,[4:1])})|\} \\ &= 3 \end{aligned}$$

Using *AdaptFun*, we first generate a global list G from an empty list $[]$ and empty whlemap \emptyset .

$$[]; \emptyset; MR^{ssa} \rightarrow G; w \wedge w = \emptyset$$

$$G_{k=2} = [I_1^{(1,\emptyset)}, I_3^{(2,[2:1])}, p_1^{(3,[2:1])}, a_1^{(4,[2:1])}, I_2^{(5,[2:1])}, I_3^{(2,[2:2])}, p_1^{(3,[2:2])}, a_1^{(4,[2:2])}, I_2^{(5,[2:2])}, I_3^{(2,\emptyset)}]$$

We denote I_1^1 short for $I_1^{(1,\emptyset)}$ and $I_3^{(2,1)}$ short for $I_3^{(2,[2:1])}$, where the label $(2, 1)$ represents at line number 2 and in the 1 st iteration.

$$M = \begin{bmatrix} I_1^1 & I_3^{(2,1)} & p_1^{(3,1)} & a_1^{(4,1)} & I_2^{(5,1)} & I_3^{(2,2)} & p_1^{(3,2)} & a_1^{(4,2)} & I_2^{(5,2)} & I_3^{(2,\emptyset)} \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \end{bmatrix}, V = \begin{bmatrix} I_1^1 & 0 \\ I_3^{(2,1)} & 0 \\ p_1^{(3,1)} & 0 \\ a_1^{(4,1)} & 1 \\ I_2^{(5,1)} & 0 \\ I_3^{(2,2)} & 0 \\ p_1^{(3,2)} & 0 \\ a_1^{(4,2)} & 1 \\ I_2^{(5,2)} & 0 \\ I_3^{(2,\emptyset)} & 0 \end{bmatrix}$$

The adaptivity is 1 computed from the graph. The query-based dependency graph is a subgraph of the variable dependency graph for multi round algorithm.

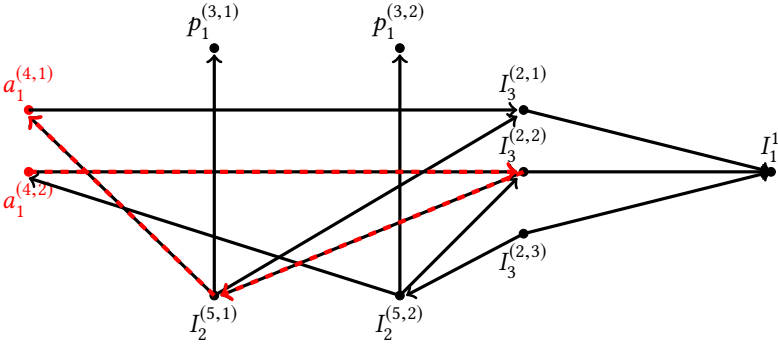


Fig. 11. the variable dependency graph for multi round algorithm

6 RELATED WORKS

REFERENCES

Cynthia Dwork, Vitaly Feldman, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Aaron Leon Roth. 2015. Preserving statistical validity in adaptive data analysis. In *Proceedings of the forty-seventh annual ACM symposium on Theory of computing*. 117–126.

A APPENDIX

Text of appendix ...