

Adaptivity Analysis

Abstract

A data analysis is said to *generalize* when it allows one to draw conclusions from the data that are true of the population from which the data are sampled. Statistician and data scientist have devised several methods aimed to guarantee generalization in data analyses, and avoid in this way overfitting to the data. Guaranteeing generalization is more difficult when data analyses are *adaptive*: when the result of an analysis depends on the result of previous analyses.

A recent line of work focuses on methods aimed at guaranteeing generalization in adaptive data analysis through the addition of carefully calibrated statistical noise to the empirical results of the analysis on the sampled data. In these works, the confidence intervals on the generalization error that one can achieve for a given analysis usually depend on the *level of adaptivity* of the analysis: the number of adaptive steps that depend on each other.

In this work we introduce a programming framework, named **AdaptFun**, aimed at supporting the study of the generalization error for adaptive data analysis. Through its analysis system, an upper bound on the adaptivity – depth (the length of the longest chain of queries) of a program implementing an adaptive data analysis can be overestimated. We show how this language can help to analyze the generalization error between two data analyses with different adaptivity structures.

Contents

1	System Overview	1
2	Low Level Control Flow Based Language	2
2.1	Syntax and Semantics	3
2.2	Adaptivity of Programs in Low level language	5
3	High level Language	7
3.1	Syntax and Semantics	7
3.2	Rewriting from High Level Program into Low Level Program	7
4	Towards Probability	10
4.1	Syntax and Semantics	10
4.2	Extending Adaptivity onto Probabilistic Program	11
5	Analysis of Program Adaptivity	12

1 System Overview

In adaptive data analysis, a data analysis can depend on the results of previous analysis over the same data. This dependency may affect the *generalization properties of the data analysis*. To study this

phenomenon in a formal way, we consider the *statistical query model*. In this model, a dataset X consisting of d attributes (columns) and n individuals' data (rows) can be accessed only through an interface to which one can submit statistical queries. More precisely, suppose that the type of a row is R (as an example, a row with d binary attributes would have type $R = \{0, 1\}^d$). Then, in the statistical query model one can access the dataset only by submitting a query to the interface, in the form of a function $p : D \rightarrow [0, 1]$ where D represents dataset. The collected answer of the asked query is the average result of p on each row in the dataset D . For example, the result is the value $\frac{1}{n} \sum_{i=1}^n p(X_i)$ where X_i is the row of index i in X . While this model is rather simple, in fact it supports sufficient statistics one may be interested.

We support the model by providing a “high” level language, which is expressive to represent the adaptive mechanisms. In this language, the queries are allow to carry arguments, for example, the expression $q(e)$ simulates the process of submitting a query to the interface in the model while the argument e is consumed to construct the query. For instance, one submitted query may use the average of answers of previous queries, can be expressed $q(x)$ where x stores the expected average results in our language. In this sense, the result of the query from a specific D can vary under different contexts.

Nevertheless, in *statistical query model*, the dependency between two submitted queries is vague, especially when we study in adaptive scenarios. Our definition of the dependency lies in the independence between two queries: one query q_1 does not depend on another query q_2 when the result of q_1 remains the same regardless of the modification of the result of q_2 . Following such definition, the dependency between two queries comes from either the control flow or the argument of the function q_2 , which brings us difficulty to distinguish them at the semantic level.

An example to show how do control flow and query argument may bring challenges.

To simplify the model, we start from a low level language where the queries are atomic – q – given data base D , the result of the query from D is deterministic. Hence, we define the adaptivity of the program under this model based on only the control flow. The arguments used in queries in the high level language are rearranged and the program expressed in the high level one can be transformed into its low level version.

We give the definition of adaptivity of a low level program by graphs, called dependency graph. These graph is produced using a trace of queries which is generated along with the semantics of the program. The queries in the trace consists of the nodes in the graph while the edge represent dependency. If there is no dependency between two node(queries), there will be no edge.

Finally we extend the language to support the probabilistic program and extend the adaptivity definition accordingly.

The key component of the system is an program analysis procedure, which provides an upper bound on the adaptivity of the program.

2 Low Level Control Flow Based Language

We first consider a low level language where the queries are atomic and the dependency relations are caused only by controal flow.

2.1 Syntax and Semantics

Syntax.

Arithmetic Operators	$*_a ::= + \mid - \mid \times \mid \div$
Boolean Operators	$*_b$
Relational Operators	$*_r ::= < \mid \leq \mid =$
Label	l
While Map	$w \in \text{Label} \times \mathbb{N}$
AExpr	$a ::= n \mid x \mid a *_a a \mid [] \mid [a_0, \dots, a_i] \mid \text{uniform} \mid \text{bernoulli} \mid a \times a$
BExpr	$b ::= \text{true} \mid \text{false} \mid \neg b \mid b *_b b \mid a *_r a$
Expr	$e ::= a \mid b$
Command	$c ::= [x \leftarrow e]^l \mid [x \leftarrow q]^l \mid [\text{switch}(e, x, (v_i \rightarrow q_i))]^l \mid \text{loop } [e_N]^l (f) \text{ do } c \mid c; c \mid \text{if}([b]^l, c_1, c_2) \mid \text{while}([b]^l, c) \mid [\text{skip}]^l \mid \text{unfold}([b]^l, c)$
Memory	$m ::= [] \mid m[x^l \rightarrow v]$
Trace	$t ::= [] \mid [(q, v)^{(l, w)}] \mid t ++ t$

Operational Semantics.

$$\langle m, a \rangle \rightarrow_a a' : \text{Memory} \times \text{AExpr} \Rightarrow \text{AExpr}$$

$$\frac{}{\langle m, x \rangle \rightarrow_a m(x)} \quad \frac{\langle m, a_1 \rangle \rightarrow_a a'_1}{\langle m, a_1 *_a a_2 \rangle \rightarrow_a a'_1 *_a a_2} \quad \frac{\langle m, a_2 \rangle \rightarrow_a a'_2}{\langle m, n_1 *_a a_2 \rangle \rightarrow_a n_1 *_a a'_2}$$

$$\frac{n_3 = n_1 *_a n_2}{\langle m, n_1 *_a n_2 \rangle \rightarrow_a n_3}$$

$$\langle m, b \rangle \rightarrow_b b' : \text{Memory} \times \text{BExpr} \Rightarrow \text{BExpr}$$

$$\frac{}{\langle m, \text{true} \rangle \rightarrow_b \text{true}} \quad \frac{}{\langle m, \text{false} \rangle \rightarrow_b \text{false}} \quad \frac{\langle m, a_1 \rangle \rightarrow_a a'_1}{\langle m, a_1 *_r a_2 \rangle \rightarrow_b a'_1 *_r a_2}$$

$$\frac{\langle m, a_2 \rangle \rightarrow_a a'_2}{\langle m, n_1 *_r a_2 \rangle \rightarrow_b n_1 *_r a'_2} \quad \frac{b_3 = n_1 *_r n_2}{\langle m, n_1 *_r n_2 \rangle \rightarrow_b b_3} \quad \frac{\langle m, b_1 \rangle \rightarrow_b b'_1}{\langle m, b_1 *_b b_2 \rangle \rightarrow_b b'_1 *_b b_2}$$

$$\frac{\langle m, b_2 \rangle \rightarrow_b b'_2}{\langle m, \text{true} *_b b_2 \rangle \rightarrow_b \text{true} *_b b'_2} \quad \frac{\langle m, b_2 \rangle \rightarrow_b b'_2}{\langle m, \text{false} *_b b_2 \rangle \rightarrow_b \text{false} *_b b'_2} \quad \frac{\langle m, b \rangle \rightarrow_b b'}{\langle m, \neg b \rangle \rightarrow_b \neg b'}$$

$$\boxed{\langle m, c, D, t \rangle \rightarrow \langle m', c', D, t' \rangle}$$

$$Memory \times Com \times Db \times Trace \times WhileMap \Rightarrow Memory \times Com \times Db \times Trace \times WhileMap$$

$$\begin{array}{c}
\frac{q(D) = v}{\langle m, [x \leftarrow q]^l, D, t, w \rangle \rightarrow \langle m[v/x], \text{skip}, D, t++[(q, v)^{(l, w)}], w \rangle} \text{query} \\
\\
\frac{m, e \Rightarrow e'}{\langle m, [x \leftarrow e]^l, D, t, w \rangle \rightarrow \langle m, [x \leftarrow e']^l, D, t, w \rangle} \text{assn1} \quad \frac{}{\langle m, [x \leftarrow v]^l, D, t, w \rangle \rightarrow \langle m[v/x], [\text{skip}]^l, D, t, w \rangle} \text{assn2} \\
\\
\frac{\langle m, c_1, D, t, w \rangle \rightarrow \langle m', c'_1, D, t', w \rangle}{\langle m, c_1; c_2, D, t, w \rangle \rightarrow \langle m', c'_1; c_2, D, t', w \rangle} \text{seq1} \quad \frac{}{\langle m, [\text{skip}]^l; c_2, D, t, w \rangle \rightarrow \langle m, c_2, D, t, w \rangle} \text{seq2} \\
\\
\frac{\langle m, b \rangle \rightarrow_b b'}{\langle m, \text{if}([b]^l, c_1, c_2), D, t, w \rangle \rightarrow \langle m, \text{if}([b']^l, c_1, c_2), D, t, w \rangle} \text{if} \\
\\
\frac{}{\langle m, \text{if}([\text{true}]^l, c_1, c_2), D, t, w \rangle \rightarrow \langle m, c_1, D, t, w \rangle} \text{if-t} \quad \frac{}{\langle m, \text{if}([\text{false}]^l, c_1, c_2), D, t, w \rangle \rightarrow \langle m, c_2, D, t, w \rangle} \text{if-f} \\
\\
\frac{}{\langle m, \text{while}([b]^l, c), D, t, w \rangle \rightarrow \langle m, \text{unfold}([b]^l, \text{while}([b]^l, c)), D, t, w \rangle} \text{while} \\
\\
\frac{\langle m, b \rangle \rightarrow b'}{\langle m, \text{unfold}([b]^l, c), D, t, w \rangle \rightarrow \langle m, \text{unfold}([b']^l, c), D, t, w \rangle} \text{unfold} \\
\\
\frac{}{\langle m, \text{unfold}([\text{false}]^l, c), D, t, w \rangle \rightarrow \langle m, [\text{skip}]^l, D, t, (w \setminus l) \rangle} \text{unfold-f} \\
\\
\frac{}{\langle m, \text{unfold}([\text{true}]^l, c), D, t, w \rangle \rightarrow \langle m, c, D, t, (w + l) \rangle} \text{unfold-t} \\
\\
\frac{\langle m, e \rangle \rightarrow e'}{\langle m, [\text{switch}(e, x, (v_i \rightarrow q_i))]^l, D, t, w \rangle \rightarrow \langle m, [\text{switch}(e', x, (v_i \rightarrow q_i))]^l, D, t, w \rangle} \text{switch} \\
\\
\frac{}{\langle m, [\text{switch}(v_k, x, (v_i \rightarrow q_i))]^l, D, t, w \rangle \rightarrow \langle m, [x \leftarrow q_k]^l, D, t, w \rangle} \text{switch-v} \\
\\
\frac{\langle m, e_N \rightarrow e'_N \rangle}{\langle m, \text{loop}[e_N]^l(f) \text{ do } c, D, t, w \rangle \rightarrow \langle m, [\text{loop}[e'_N]^l(f) \text{ do } c]^l, D, t, w \rangle} \text{loop} \\
\\
\frac{v_N > 0}{\langle m, \text{loop}[v_N]^l(f) \text{ do } c, D, t, w \rangle \rightarrow \langle m, c; \text{loop}[(v_N - 1)]^l(f) \text{ do } c, D, t, (w + l) \rangle} \text{loop-t} \\
\\
\frac{v_N = 0}{\langle m, \text{loop}[v_N]^l(f) \text{ do } c, D, t, w \rangle \rightarrow \langle m, [\text{skip}]^l, D, t, (w \setminus l) \rangle} \text{loop-f}
\end{array}$$

where w_l refers to a map w without the key l .

$$\begin{aligned} w \setminus l &= w & l \notin \text{Keys}(w) \\ &= w_l & \text{Otherwise} \\ w + l &= w[l \rightarrow 0] & l \notin \text{Keys}(w) \\ &= w_l[l \rightarrow w(l) + 1] & \text{Otherwise} \end{aligned}$$

2.2 Adaptivity of Programs in Low level language

Definition 1 (Label Order). $<_w$ and $=_w$.

$$\begin{aligned} w_1 =_w w_2 &\triangleq \text{Keys}(w_1) = \text{Keys}(w_2) \wedge \forall k \in \text{Keys}(w_1). w_1(k) = w_2(k) \\ \emptyset =_w \emptyset & \end{aligned}$$

$$mk(w_i) = \text{MinKey}(w_i)$$

$$\begin{aligned} w_1 <_w w_2 &\triangleq & w_1 = \emptyset \\ &\triangleq mk(w_1) < mk(w_2) & w_1, w_2 \neq \emptyset \\ &\triangleq w_1(mk(w_1)) < w_2(mk(w_2)) & mk(w_1) = mk(w_2) \\ &\triangleq (w_1 \setminus mk(w_1)) <_w (w_2 \setminus mk(w_2)) & \text{Otherwise} \end{aligned}$$

Definition 2 (Query Direction). *Direction between two queries.*

$\forall Q_1, Q_2, l_1, l_2, w_1, w_2. (Q_1^{l_1, w_1}) \text{TO} (Q_2^{l_2, w_2})$, denoted as $\text{To}(Q_1^{l_1, w_1}, Q_2^{l_2, w_2})$

iff $(Q_1^{l_1, w_1}) <_q (Q_2^{l_2, w_2})$

where

$(Q_1^{l_1, w_1}) <_q (Q_2^{l_2, w_2})$ is defined:

$$\begin{aligned} l_1 < l_2 & \quad w_1 = \emptyset \vee w_2 = \emptyset \vee w_1 =_w w_2 \\ w_1 <_w w_2 & \quad \text{Otherwise} \end{aligned}$$

Independence between two queries in Low level language When two queries q_1, q_2 are independent in a program c , suppose q_1 appears before q_2 in the program c , we think the choice of queries starting from q_1 , ending with query q_2 should be fixed no matter the change of the result of q_1 .

Definition 3 (Query Independence). *Two queries q_i and q_j in a program c are independent, $\text{IND}(q_i^{l_1}, q_j^{l_2}, c)$.*

$$\begin{aligned} &\forall m, D. \left(\langle m, c, D, [] \rangle \rightarrow \langle m', \text{skip}, D, t \rangle \right. \\ &\wedge \left((q_i^{l_1}, v_i) \in t \wedge (q_j^{l_2}, v_j) \in t \implies \forall v \in \text{codomain}(q_i^{l_1}). \left(\langle m, c[v/q_i], D, [] \rangle \rightarrow \langle m', \text{skip}, D, t' \rangle \wedge (q_j^{l_2}, v_j) \in t' \right) \right) \\ &\left. \wedge \left((q_i^{l_1}, v_i) \in t \wedge (q_j^{l_2}, v_j) \notin t \implies \forall v \in \text{codomain}(q_i^{l_1}). \left(\langle m, c[v/q_i], D, [] \rangle \rightarrow \langle m', \text{skip}, D, t' \rangle \wedge (q_j^{l_2}, v_j) \notin t' \right) \right) \right) \end{aligned}$$

In following examples:

$$\begin{aligned} c_1 &\triangleq \begin{array}{l} [x \leftarrow q_1]^1; \\ \text{if } [x]^2 \\ \text{then } [y \leftarrow q_2]^3 \\ \text{else } [y \leftarrow 0]^4; \\ \text{if } [x]^5 \\ \text{then } [y \leftarrow 0]^6 \\ \text{else } [y \leftarrow q_2]^7; \end{array} & c_2 &\triangleq \begin{array}{l} [x \leftarrow q_1]^1; \\ [y \leftarrow q_2]^2; \\ \text{if } [x + y = 5]^3 \\ \text{then } [z \leftarrow q_3]^4 \\ \text{else } [\text{skip}]^5; \\ [w \leftarrow q_4]^6; \end{array} & c_3 &\triangleq \begin{array}{l} [x \leftarrow q_1]^1; \\ \text{while } ([x \leq 100])^2, \\ [y \leftarrow q_2]^3; \\ [z \leftarrow q_3]^4; \\ [x \leftarrow y + z + x]^5; \end{array} & c_4 &\triangleq \begin{array}{l} [x \leftarrow q_1]^1; \\ \text{while } ([x \leq 100])^2, \\ [y \leftarrow q_2 + x]^3; \\ \text{while } ([y \leq 50])^4, \\ [z \leftarrow q_3]^5; \\ [y \leftarrow y + z]^6; \\ [x \leftarrow y + x]^7; \end{array} \end{aligned}$$

We have the dependency as:

In program c_1 : $\text{IND}(q_2^3, q_2^7, c_1)$,

In program c_2 : $\text{IND}(q_1^1, q_2^2, c_4)$, $\text{IND}(q_3^3, q_4^4, c_4)$, $\text{IND}(q_1^1, q_4^4, c_4)$, $\text{IND}(q_2^2, q_4^4, c_4)$,

In program c_3 : $\text{IND}(q_2^{(3,l)}, q_3^{(4,l)}, c_2)$ for all l .

In program c_4 : No independent queries.

Dependency between multiple queries

Definition 4 (Dependency Graph). A dependency graph over a program P is defined as $G(P) = (V, E, \text{end})$, where V is set of verticals and E is the set of directed edges:

$V = \{q_1, q_2, \dots, q_n\}$, where q_1, q_2, \dots, q_n are reachable queries in the program P .

$E = \{(q_i, q_j) | \exists m. \text{Dep}(q_i, q_j, P, m)\}$, in the program P .

$\text{end} = \{q_{e_1}, \dots, q_{e_l}\}$ be the set of queries related to the return value.

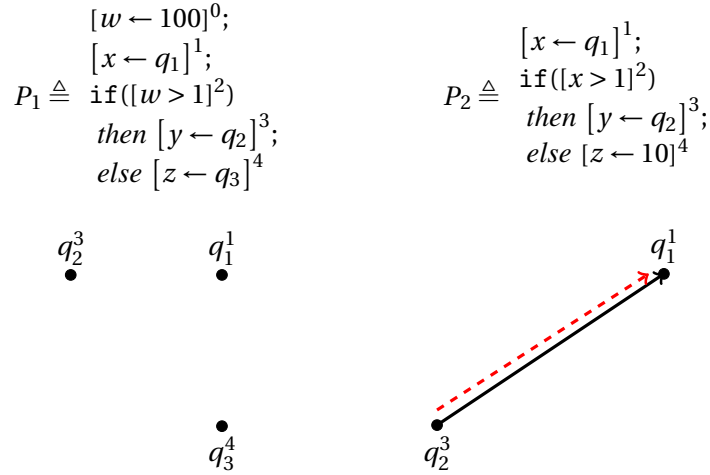
Definition 5 (Reachable Query). A query q in program P is reachable iff q may be executed.

Definition 6 (Adaptivity). Given a program P and its dependency graph $G(P) = (V, E, \text{end})$, the adaptivity of the program is defined as $A(P)$, s.t.: for every $q_i \in \text{end}$, let p_i be the longest path starting from q_i with length l_i ,

$$A(P) = \max_{q_i \in \text{end}} \{l_i \mid l_i = |p_i|\}$$

Example 2.1. Dependency graphs for programs containing 3 atomic queries

Let $q_1 \triangleq \lambda D. D_1 * D_j$, $q_2 \triangleq \lambda D. D_3 * D_4$ and $q_3 \triangleq \lambda D. D_3 * D_2$. in program P_1 and P_2 as follows:



3 High level Language

3.1 Syntax and Semantics

Syntax.

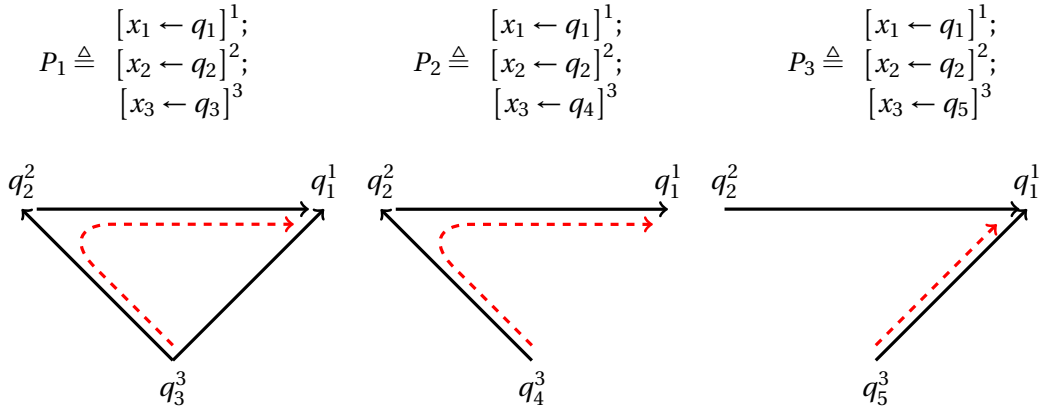
Arithmetic Operators	$*_a ::= + \mid - \mid \times \mid \div$
Boolean Operators	$*_b$
Relational Operators	$*_r ::= < \mid \leq \mid =$
Label	l
While Map	$w \in \text{Label} \times \mathbb{N}$
AExpr	$a ::= n \mid x \mid a *_a a \mid [] \mid [a_0, \dots, a_i] \mid \text{uniform} \mid \text{bernoulli} \mid a \times a$
BExpr	$b ::= \text{true} \mid \text{false} \mid \neg b \mid b *_b b \mid a *_r a$
Command	$c ::= [x \leftarrow e]^l \mid [x \leftarrow q(e)]^l \mid c; c \mid \text{if}([b]^l, c_1, c_2) \mid \text{while}([b]^l, c) \mid [\text{skip}]^l \mid \text{loop } [v_N]^l (f) \text{ do } c$
Memory	$m ::= [] \mid m[x^l \rightarrow v]$
Trace	$t ::= [] \mid [(q, v)^{(l, w)}] \mid t ++ t$

Example 3.1. Dependency graphs for high level programs containing non-atomic queries

Let $q_1 = \lambda D. D_i * D_j$,

Let $q_2(x_1) = \lambda D. D_i * D_j + x_1$.

Let $q_3(x_1 - x_2) = \lambda D. D_i * D_j + x_1 - x_2$, $q_4(x_2) = \lambda D. D_i * D_j + x_2$, and $q_5(x_1) = \lambda D. D_i * D_j + x_1$.
in program P_1 , P_2 and P_3 as following:



3.2 Rewriting from High Level Program into Low Level Program

The transformation $\langle e^h \rangle = e^l$ transfers the expression e^h in the high level language to an expression e^l in the low language. Let us look at the special cases: the query.

In the first transition, if a query in high level language isn't atomic, i.e., $q(e)$ depends on e with free variables, then it will be rewrite into a switch command. This rewriting will switch on the possible values v_i of e and convert the $q(e)$ into a series of atomic queries q_i .

In the second transition, if a query in high level language is atomic, $q()$ only depends on data base D and some constant values, then it will be rewrite into identity in our low level language.

Another special case is the sampling command in high level language. To exclude the dependency caused by the randomness, we will rewrite the sampling into an assignment command in low level

language. This will assign a constant value to the corresponding variable.
The resting commands will be rewrote identically.

$$\begin{aligned}
\langle [x \leftarrow q(e)]^l \rangle &\Rightarrow \left[\text{switch} \left(e, x, \begin{pmatrix} v_1 \rightarrow q_1, \\ \cdots, \\ v_m \rightarrow q_m \end{pmatrix} \right) \right]^l \\
\langle q() \rangle &\Rightarrow q \\
\langle [x \leftarrow \text{uniform}]^l \rangle &\Rightarrow [x \leftarrow c_u]^l \\
\langle [x \leftarrow e]^l \rangle &\Rightarrow [x \leftarrow e]^l \\
\langle c_1; c_2 \rangle &\Rightarrow \langle c_1 \rangle; \langle c_2 \rangle \\
\langle \text{if}([b]^l, c_1, c_2) \rangle &\Rightarrow \text{if}([b]^l, \langle c_1 \rangle, \langle c_2 \rangle) \\
\langle \text{while}([b]^l, c) \rangle &\Rightarrow \text{while}([b]^l, \langle c \rangle)
\end{aligned}$$

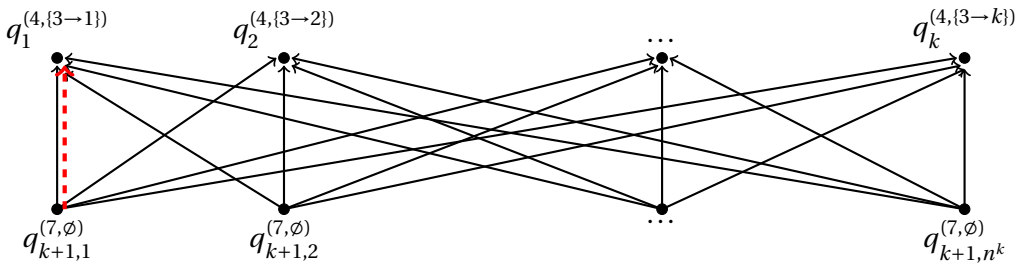
Example 3.2 (Two Round Algorithm).

$$\begin{aligned}
&\begin{aligned}
&[j \leftarrow 1]^1; \\
&[a \leftarrow []]^2; \\
&\text{while}([j \leq k]^3, \\
&\quad [x \leftarrow q_j()]^4; \\
&\quad [a \leftarrow x :: a]^5; \\
&\quad [j \leftarrow j+1]^6); \\
&[l \leftarrow q_{k+1}(a)]^7;
\end{aligned}
\end{aligned}
\Rightarrow
\begin{aligned}
&\begin{aligned}
&[j \leftarrow 1]^1; \\
&[a \leftarrow []]^2; \\
&\text{while}([j \leq k]^3, \\
&\quad [x \leftarrow q_j]^4; \\
&\quad [a \leftarrow x :: a]^5; \\
&\quad [j \leftarrow j+1]^6); \\
&\left[\text{switch} \left(a, x, \begin{pmatrix} [-n, -n, -n, \dots, -n] \rightarrow q_{k+1,1}, \\ \cdots \\ [n, n, n, \dots, n] \rightarrow q_{k+1, n^k} \end{pmatrix} \right) \right]^7
\end{aligned}
\end{aligned}$$

$$\langle \emptyset, TR^L, D, [], \emptyset \rangle \rightarrow \langle m, \text{skip}, D, t, w \rangle.$$

$$t = [(q_1^{(4, \{3 \rightarrow 1\})}, v_1), (q_2^{(4, \{3 \rightarrow 2\})}, v_2), \dots, (q_k^{(4, \{3 \rightarrow k\})}, v_k), (q_{k+1, i}^{(7, \emptyset)}, v_1)]$$

$$A(TR^L) = I$$



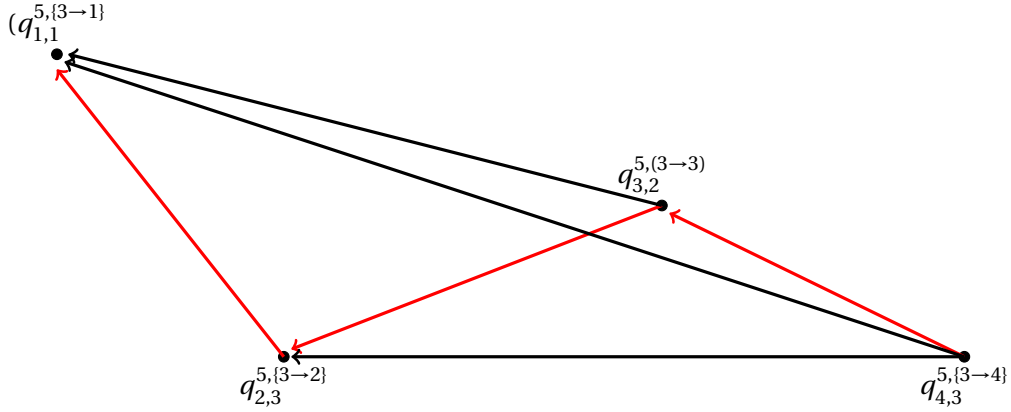
Example 3.3 (Multi-Round Algorithm).

$$\begin{aligned}
 MR^H \triangleq & \begin{aligned} & [j \leftarrow 0]^1; \\ & [I \leftarrow []]^2; \\ & \text{while}([j \leq k]^3, \\ & [p \leftarrow \text{uniform}(0,1)]^4; \\ & [a \leftarrow q_j(p,D)]^5; \\ & [I \leftarrow \text{update}(I, (a,p))]^6; \\ & [j \leftarrow j+1]^7); \end{aligned} \quad \Rightarrow \quad MR^L \triangleq \begin{aligned} & [j \leftarrow 0]^1; \\ & [I \leftarrow []]^2; \\ & \text{while}([j \leq k]^3, \\ & [p \leftarrow 0]^4; \\ & \left[\text{switch} \left(I, x \begin{pmatrix} [] \rightarrow q_{j,1}, \\ \dots \\ [1,2,3,\dots,n] \rightarrow q_{j,n!} \end{pmatrix} \right) \right]^5 \\ & [I \leftarrow \text{update}(I, (a,p))]^6; \\ & [j \leftarrow j+1]^7); \end{aligned}
 \end{aligned}$$

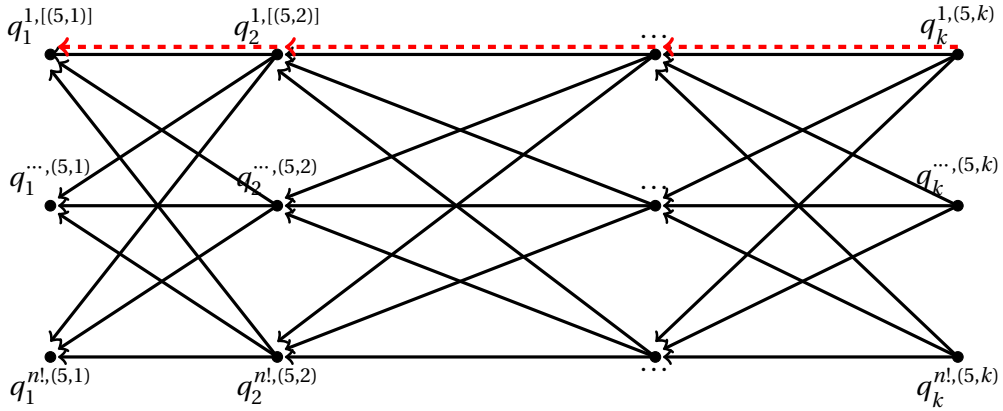
Let $k = 4$, given a specific database D , we have $\langle \emptyset, MR^L, D, [], \emptyset \rangle \rightarrow \langle m, \text{skip}, D, t, w \rangle$ and the trace as:

$$t = \left[(q_{1,1}^{5,\{3 \rightarrow 1\}}, v_1), (q_{2,3}^{5,\{3 \rightarrow 2\}}, v_2), (q_{3,2}^{5,\{3 \rightarrow 3\}}, v_3), (q_{4,3}^{5,\{3 \rightarrow 4\}}, v_4) \right]$$

$$A(TR^L) = 3$$



$\forall k. \forall D$, we have $A(TR^L) = (k - 1)$ given all possible execution traces.



4 Towards Probability

4.1 Syntax and Semantics

Syntax.

Arithmetic Operators	$*_a ::= + \mid - \mid \times \mid \div$
Boolean Operators	$*_b ::= \vee \mid \wedge$
Relational Operators	$*_r ::= < \mid \leq \mid =$
Label	l
While Map	$w \in \text{Label} \times \mathbb{N} \triangleq w + l \mid w \setminus l \mid w \oplus_\rho w$
AExpr	$a ::= n \mid x \mid a *_a a \mid [] \mid [a_0, \dots, a_i] \mid a \times a$
BExpr	$b ::= \text{true} \mid \text{false} \mid \neg b \mid b *_b b \mid a *_r a$
Deterministic Expr	$e_d ::= a \mid b$
Randomized AExpr	$a_r ::= x_r \mid n \mid a_r *_a a_r$
Randomized BExpr	$b_r ::= \neg b_r \mid b_r *_b b_r \mid a_r *_r a_r$
Randomized Expr	$e_r ::= e_d \mid a_r \mid b_r$
Command	$C ::= [x \leftarrow e_d]^l \mid [x_r \leftarrow e_r]^l \mid [x \leftarrow q]^l \mid [x_r \leftarrow \text{uniform}]^l \mid [x_r \leftarrow \text{bernoulli}]^l \mid P; P \mid \text{if}_D([b]^l, P_1, P_2) \mid \text{if}_R([b_r]^l, P_1, P_2) \mid \text{while}([b]^l, P) \mid [\text{skip}]^l \mid \text{unfold}([b]^l, P) \mid [\text{switch}(e, x, (v_i \rightarrow q_i))]^l \mid \text{loop } [v_N]^l (f) \text{ do } P$
Memory	$m ::= [] \mid m[x^l \rightarrow v]$
Trace	$t ::= [] \mid [(q, v)^{(l, w)}] \mid t ++ t \mid t \oplus_\rho t$

$$\begin{aligned}
 w \setminus l &= \begin{cases} w & l \notin \text{Keys}(w) \\ w_l & \text{Otherwise} \end{cases} \\
 w + l &= \begin{cases} w[l \rightarrow 0] & l \notin \text{Keys}(w) \\ w_l[l \rightarrow w(l) + 1] & \text{Otherwise} \end{cases}
 \end{aligned}$$

Semantics We have a countable set RV of random variables ($x_r \in \text{RV}$), a countable set Val of values. For any subset $S \subseteq \text{RV}$, we denote $\text{RanM}[S] \triangleq S \rightarrow \text{Val}$. We let DV as a countable set of deterministic variables (x) and the deterministic memory $\text{DetM} \triangleq \text{DV} \rightarrow \text{Val}$. Distribution over A as $D(A)$. The *distribution unit* $\text{unit} : A \rightarrow D(A)$. The *distribution bind* $\text{bind} : D(A) \rightarrow (A \rightarrow D(B)) \rightarrow D(B)$. $\llbracket C \rrbracket : (\text{DetM} \times D(\text{RandM}) \times \text{Trace} \times \text{WhileMap} \times \text{DB}) \rightarrow (\text{DetM} \times D(\text{RandM}) \times \text{Trace} \times \text{WhileMap} \times \text{DB})$

$$\begin{aligned}
 \langle (\sigma, \mu_1, t_1, w_1, D) \rangle \oplus_\rho \langle (\sigma, \mu_2, t_2, w_2, D) \rangle &\triangleq \langle \sigma, \mu_1 \oplus_\rho \mu_2, t_1 \oplus_\rho t_2, w_1 \oplus_\rho w_2, D \rangle \\
 (t_1 \oplus_\rho t_2) ++ t_3 &\triangleq (t_1 ++ t_3) \oplus_\rho (t_2 ++ t_3) \\
 t_3 ++ (t_1 \oplus_\rho t_2) &\triangleq (t_3 ++ t_1) \oplus_\rho (t_3 ++ t_2) \\
 (w_1 \oplus_\rho w_2) + l &\triangleq (w_1 + l) \oplus_\rho (w_2 + l) \\
 (w_1 \oplus_\rho w_2) \setminus l &\triangleq (w_1 \setminus l) \oplus_\rho (w_2 \setminus l)
 \end{aligned}$$

$$\begin{aligned}
\llbracket [\text{skip}]^l \rrbracket (\sigma, \mu, t, w, D) &\triangleq (\sigma, \mu, t, w, D) \\
\llbracket [x \leftarrow e_d]^l \rrbracket (\sigma, \mu, t, w, D) &\triangleq (\sigma[x \rightarrow \llbracket e_d \rrbracket (\sigma)], \mu, t, w, D) \\
\llbracket [x_r \leftarrow e_r]^l \rrbracket (\sigma, \mu, t, w, D) &\triangleq (\sigma, \text{bind}(\mu, m \rightarrow \text{unit}(m[x_r \rightarrow \llbracket e_r \rrbracket (\sigma, m)])), t, w, D) \\
\llbracket [x \leftarrow q]^l \rrbracket (\sigma, \mu, t, w, D) &\triangleq (\sigma[x \rightarrow v], \mu, t + [(q, v)]^{(l, w)}, w, D) \quad : v = q(D) \\
\llbracket [x_r \leftarrow \text{uniform}]^l \rrbracket (\sigma, \mu, t, w, D) &\triangleq (\sigma, \text{bind}(\mu, m \rightarrow \text{bind}(\text{uniform}, u \rightarrow m[x_r \rightarrow u])), t, w, D) \\
\llbracket [x_r \leftarrow \text{bernoulli}]^l \rrbracket (\sigma, \mu, t, w, D) &\triangleq (\sigma, \text{bind}(\mu, m \rightarrow \text{bind}(\text{bernoulli}, u \rightarrow m[x_r \rightarrow u])), t, w, D) \\
\llbracket [P; P'] \rrbracket (\sigma, \mu, t, w, D) &\triangleq \llbracket P' \rrbracket (\llbracket P \rrbracket \sigma, \mu, t, w) \\
\llbracket [\text{if}_D([b]^l, P_1, P_2)] \rrbracket (\sigma, \mu, t, w, D) &\triangleq \begin{cases} \llbracket P_1 \rrbracket (\sigma, \mu, t, w, D) & : \llbracket b \rrbracket (\sigma) = \text{true} \\ \llbracket P_1 \rrbracket (\sigma, \mu, t, w, D) & : \llbracket b \rrbracket (\sigma) = \text{false} \end{cases} \\
\llbracket [\text{if}_R([b_r]^l, P_1, P_2)] \rrbracket (\sigma, \mu, t, w, D) &\triangleq \begin{cases} \llbracket P_1 \rrbracket (\sigma, \mu | \llbracket b_r \rrbracket \sigma = \text{true}, t, w, D) \oplus_\rho \llbracket P_2 \rrbracket (\rho, \mu | \llbracket b_r \rrbracket \sigma = \text{false}, t, w, D) \\ \llbracket P_1 \rrbracket (\sigma, \mu | \llbracket b_r \rrbracket \sigma = \text{true}, t, w, D) & \rho = 1 \\ \llbracket P_2 \rrbracket (\sigma, \mu | \llbracket b_r \rrbracket \sigma = \text{false}, t, w, D) & \rho = 0 \end{cases} \\
&\text{where } \rho = \mu(\llbracket b_r \rrbracket \sigma = \text{true}) \\
\llbracket [\text{while}([b]^l, P)] \rrbracket (\sigma, \mu, t, w, D) &\triangleq \llbracket [\text{unfold}([b]^l, \text{while}([b]^l, P))] \rrbracket (\sigma, \mu, t, w, D) \\
\llbracket [\text{unfold}([b]^l, P)] \rrbracket (\sigma, \mu, t, w, D) &\triangleq \begin{cases} \llbracket P \rrbracket (\sigma, \mu, t, w + l, D) & : \llbracket b \rrbracket (\sigma) = \text{true} \\ \llbracket [\text{skip}] \rrbracket (\sigma, \mu, t, w - l, D) & : \llbracket b \rrbracket (\sigma) = \text{false} \end{cases} \\
\llbracket [\text{switch}(e, x, (v_i \rightarrow q_i))]^l \rrbracket (\sigma, \mu, t, w, D) &\triangleq \llbracket [x \leftarrow q_1]^l \rrbracket (\sigma, \mu, t, w, D) \quad : v_1 = \llbracket e \rrbracket (\sigma) \\
\llbracket [\text{loop } [e_N]^l \text{ (f) do } P] \rrbracket (\sigma, \mu, t, w, D) &\triangleq \begin{cases} \llbracket [f; P; \text{loop } [e_N - 1]^l \text{ (f) do } P] \rrbracket (\sigma, \mu, t, w + l, D) & : \llbracket e_N \rrbracket (\sigma) > 0 \\ \llbracket [\text{skip}] \rrbracket (\sigma, \mu, t, w - l, D) & : \llbracket e_N \rrbracket (\sigma) = 0 \end{cases}
\end{aligned}$$

Figure 1: Semantics of programs

$$\begin{aligned}
&[j \leftarrow 0]^1; \\
&[I \leftarrow []]^2; \\
&\text{while}([j \leq k]^3, \\
&[p \leftarrow 0]^4; \\
&\left[\text{switch} \left(I, x \begin{pmatrix} [] \rightarrow q_{j,1}, \\ \dots \\ [1, 2, 3, \dots, n] \rightarrow q_{j,n!} \end{pmatrix} \right) \right]^5 \Rightarrow \text{SSA} \triangleq \left[\text{switch} \left(I_2, x \begin{pmatrix} [] \rightarrow q_{j,1}, \\ \dots \\ [1, 2, 3, \dots, n] \rightarrow q_{j,n!} \end{pmatrix} \right) \right]^5 \\
&[I \leftarrow \text{update}(I, (x, p))]^6; \\
&[j \leftarrow j + 1]^7); \\
&[j \leftarrow 0]^1; \\
&[I_1 \leftarrow []]^2; \\
&\text{while}([j \leq k]^3, \phi : I_2 = f(I_1, I_3) \\
&[p \leftarrow 0]^4; \\
&\left[\text{switch} \left(I_2, x \begin{pmatrix} [] \rightarrow q_{j,1}, \\ \dots \\ [1, 2, 3, \dots, n] \rightarrow q_{j,n!} \end{pmatrix} \right) \right]^5 \\
&[I_3 \leftarrow \text{update}(I_2, (x, p))]^6; \\
&[j \leftarrow j + 1]^7);
\end{aligned}$$

4.2 Extending Adaptivity onto Probabilistic Program

Definition 7. *Dependency Forest.*

Given a program P , a database D , dependency forest $F(P, D) \triangleq \{G_1, G_2, \dots, G_m\}$, $G_k = (V_k, E_k)$ is defined as:

$$V_k = \{q^{l,w} | \forall m, w. \langle \emptyset, P, D, [], \emptyset \rangle \rightarrow \langle m, \text{skip}, D, T, \emptyset \rangle \wedge q^{l,w} \in \pi_k(T)\};$$

$$E_k = \{(q_i^{l,w}, q_j^{l',w'}) \mid \neg \text{IND}(q_i^{l,w}, q_j^{l',w'}, P) \wedge \text{To}(q_j^{l',w'}, q_i^{l,w})\}.$$

where $\pi_k(T) = t_k$ s.t. $T = t_1 \oplus_{\rho_1} t_2 \oplus_{\rho_2} \dots \oplus_{\rho_{m-1}} t_m$ and there is no \oplus in t_i .

Definition 8. *Dependency Graph.*

Given a program P , a database D , dependency graph $G(P, D) \triangleq (V, E)$ is defined as:

$$V = \bigcup \{V_k \mid G_k \in F(P, D)\};$$

$$E_k = \{(q_i, q_j) \mid (q_i, q_j) \in E_k, G_k \in F(P, D)\},$$

Definition 9. *Adaptivity.*

Given a program P and for all the database D in a set of DBS of databases, the total dependency graph G is the combination of all the dependent graphs over every single database $G(P, D) = (V, E)$, the adaptivity of the program is defined as $A(P)$, s.t.: for every pair (i, j) let $p_{(i,j)}$ be the longest path starting from $q_i^{l,w}$ to $q_j^{l',w'}$,

$$A(P) = \max_{q_i^{l,w}, q_j^{l',w'} \in V} \{l_i \mid l_i = |p_{(i,j)}|\}$$

5 Analysis of Program Adaptivity

We have the judgment of the form : $\vdash_{M,V}^{c_1, c_2} P : \Phi \implies \Psi$. Our grade is a combination of a matrix M , used to track the dependency of variables appeared in the statement S , and a vector V indicating the variables associated with results from queries q . The size of the matrix M is $L \times L$, and vector V of size N , where L is the total size of variables needed in the program, which is fixed per program. The superscript c_1, c_2 is a counter specifying the range of "living" or "active" variables in the matrix and vector. c_1 is the starting line (and column) in the matrix where the new generated variables in program P starts to show up. Likewise, c_2 states the ending position of active range by P . The new generated variables will be treated carefully when we handle recursive statement, in our case, the loop statement.

We assume the program P is in the static single assignment form. That is to say, $x \leftarrow e_1; x \leftarrow e_2$ will be rewritten as $x_1 \leftarrow e_1; x_2 \leftarrow e_2$. And the if condition `if e_b then $x \leftarrow e_1$ else $x \leftarrow e_2$` will look like `if e_b then $x_1 \leftarrow e_1$ else $x_2 \leftarrow e_2$` . As we have seen, ssa provides unique variables, and these newly generated variables will be recorded in the matrix M . Worth to mention, c_1, c_2 can be used to track the exact location of newly generated variables. For example, the assignment statement $x \leftarrow e$ or $x \leftarrow q$ holds $c_2 = c_1 + 1$, telling us the variable x is at the c_1 th line(column) of the matrix. As we can notice, the loop increases the matrix by $N \times a$ where N is the number of rounds of the loop and a is the size of the variables generated in the loop body P .

We will have a global map, which maps the variable name to the position in the vector. We call it $G : VAR \rightarrow \mathbb{N}$,

We give an example of M and V of the program P .

$$\begin{array}{ll}
 & x_1 \leftarrow q; \\
 P = & x_2 \leftarrow x_1 + 1; \\
 & x_3 \leftarrow x_2 + 2
 \end{array}
 \quad
 \begin{array}{ll}
 & \begin{matrix} 0 & 0 & 0 & 1 \end{matrix} \\
 M = & \begin{matrix} 1 & 0 & 0 \end{matrix}, V = 0 \\
 & \begin{matrix} 1 & 1 & 0 & 0 \end{matrix}
 \end{array}$$

Analysis Rules.

$$\boxed{\Gamma \vdash_{M,V}^{c_1, c_2} P : \Phi \Rightarrow \Psi}$$

$$\frac{M = L(x) * (R(e) + \Gamma)}{\Gamma \vdash_{M, V_\emptyset}^{(c, c+1)} x \leftarrow e : \Phi \Rightarrow \Psi}$$

$$\frac{M = L(x) * (R(\emptyset) + \Gamma) \quad V = L(x)}{\Gamma \vdash_{M, V}^{(c, c+1)} x \leftarrow q : \Phi \Rightarrow \Psi}$$

$$\frac{\Gamma + R(e_b) \vdash_{M_1, V_1}^{(c_1, c_2)} P_1 : \Phi \Rightarrow \Psi \quad \Gamma + R(e_b) \vdash_{M_2, V_2}^{(c_2, c_3)} P_2 : \Phi \Rightarrow \Psi}{\Gamma \vdash_{M_1 \uplus M_2, V_1 \uplus V_2}^{(c_1, c_3)} \text{if } e_b \text{ then } P_1 \text{ else } P_2 : \Phi \Rightarrow \Psi}$$

$$\frac{\Gamma \vdash_{M_1, V_1}^{(c_1, c_2)} P_1 : \Phi \Rightarrow \Psi' \quad \Gamma \vdash_{M_2, V_2}^{(c_2, c_3)} P_2 : \Psi' \Rightarrow \Psi}{\Gamma \vdash_{M_1 \cdot M_2, V_1 \uplus V_2}^{(c_1, c_3)} P_1; P_2 : \Phi \Rightarrow \Psi}$$

$$\frac{\forall 0 \leq z < N. \Gamma \vdash_{M, V}^{(c, c+a)} f; P : \{\Phi \wedge e_N = z + 1\} \Rightarrow \{\Phi \wedge e_N = z\}}{\Gamma \vdash_{M_{c,a}^N(f), V_{c,a}^N}^{(c, c+N*a)} \text{loop } e_N (f) \text{ do } P : \{\Phi \wedge e_N = N\} \Rightarrow \{\Phi \wedge e_N = 0\}}$$

$$\frac{\Gamma + R(e) \vdash_{M_i, V_i}^{(c_1, c_1+1)} x_i \leftarrow q_i \quad i \in \{1, \dots, N\}}{\Gamma \vdash_{\sum_{i=0}^N M_i, \sum_{i=0}^N V_i}^{(c_1, c_1+N)} \text{switch}(e, x, (v_i \rightarrow q_i))}$$

$$\begin{aligned} V_1 \uplus V_2 &:= \begin{cases} 1 & (V_1[i] = 1 \vee V_2[i] = 1) \wedge i = 1, \dots, N \wedge |V_1| = |V_2| \\ 0 & o.w. \end{cases} \\ M_1 \uplus M_2 &:= \begin{cases} 1 & (M_1[i][j] = 1 \vee M_2[i][j] = 1) \wedge i, j = 1, \dots, N \wedge |M_1| = |M_2| \\ 0 & (M_1[i][j] = 0 \wedge M_2[i, j] = 0) \wedge i, j = 1, \dots, N \wedge |M_1| = |M_2| \\ \perp & o.w. \end{cases} \\ V_{(c,a)}^N &:= \begin{cases} V[c + i * a, c + (i + 1) * a - 1] = V[c, c + a - 1] & i = 1, \dots, N - 1 \\ \perp & o.w. \end{cases} \\ M_{(c,a)}^N(f) &:= \begin{cases} M[c + i * a, c + (i + 1) * a - 1][c + i * a, c + (i + 1) * a - 1] \\ = M[c, c + a - 1][c, c + a - 1] & i = 1, \dots, N - 1 \\ M[c + i * a, c + (i + 1) * a - 1][0, c + (i) * a - 1] = 0 & i = 1, \dots, N - 1 \\ M[0, c + i * a - 1][c + i * a, c + (i + 1) * a - 1] \\ = M[0, c + (i - 1) * a - 1][c + (i - 1) * a, c + (i) * a - 1] & i = 1, \dots, N - 1 \\ M[l][k] = 1 & l \in [c + (i - 1) * a, c + (i) * a - 1], \\ & k \in [c + i * a, c + (i + 1) * a - 1] \\ & \wedge i = 1, \dots, N \wedge, l = G(x_l) \\ & \wedge k = G(x_k) \wedge x_l = f(\dots, x_k, \dots) \\ \perp & o.w. \end{cases} \end{aligned}$$

1	...	c-1	c	...	c+a-1	c+a	...	c+2a-1	...	c+N*a-1	c+N*a	...
...			0	0	0	0	0	0				
c-1			0	0	0	0	0	0				
c						0	0	0				
...						0	0	0				
c+a-1						0	0	0				
c+a												
...				f								
c+2a-1												
...												
c+N*a-1												
c+N*a												
...												

Definition 10. *Dependency Graph.*

Given a program P , a database D , dependency graph $G(P, D) = (V, E)$ is defined as:

$$V = \{q_i^{l,w} \mid \forall \sigma, \sigma', \mu, \mu', w, w', t, t'. \llbracket P \rrbracket(\sigma, \mu, t, w, D) \triangleq (\sigma', \mu', t', w', D) \wedge q_i^{l,w} \in (t' - t)\}.$$

$$E = \{(q_i^{l,w}, q_j^{l',w'}) \mid \neg \text{IND}(q_i^{l,w}, q_j^{l',w'}, P) \wedge \text{To}(q_j^{l',w'}, q_i^{l,w})\}.$$

Definition 11. Two queries q_i and q_j in a program c are independent, $\text{IND}(q_i^{l_1}, q_j^{l_2}, P)$.

$$\begin{aligned} & \forall m, D. \left(\llbracket P \rrbracket(\sigma, \mu, \mathbf{t}, w, D) \triangleq (\sigma', \mu', t', w', D) \right. \\ & \wedge \left((q_i^{l_1}, v_i) \in t' \wedge (q_j^{l_2}, v_j) \in t' \implies \forall v \in \text{Codom}(q_i^{l_1}). \left(\llbracket P[v/q_i] \rrbracket(\sigma, \mu, \mathbf{t}, w, D) \triangleq (\sigma', \mu', t'', w', D) \wedge (q_j^{l_2}, v_j) \in t'' \right) \right) \\ & \wedge \left((q_i^{l_1}, v_i) \in t \wedge (q_j^{l_2}, v_j) \notin t \implies \forall v \in \text{Codom}(q_i^{l_1}). \left(\llbracket P[v/q_i] \rrbracket(\sigma, \mu, \mathbf{t}, w, D) \triangleq (\sigma', \mu', t'', w', D) \wedge (q_j^{l_2}, v_j) \notin t'' \right) \right) \Big). \end{aligned}$$

Definition 12. *Adaptivity.*

Given a program P and for all the database D in a set of DBS of databases, the total dependency graph G is the combination of all the dependent graphs over every single database $G(P, D) = (V, E)$, the adaptivity of the program is defined as $A(P)$, s.t.: for every pair (i, j) let $p(i, j)$ be the longest path starting from $q_i^{l,w}$ to $q_j^{l',w'}$,

$$A(P) = \max_{q_i^{l,w}, q_j^{l',w'} \in V} \{l_i \mid l_i = |p(i, j)|\}$$

Definition 13 (Adapt). Given a program P s.t. $\cdot \vdash_{M,V}^{c_1, c_2} P : \Phi \implies \Psi$, there exists 1 and only one Graph $G(M, V) = (\text{Nodes}, \text{Edges}, \text{Weights})$ defined as:

$$\text{Nodes} = \{i \mid i \in V\}$$

$$\text{Edges} = \{(i, j) \mid M[i][j] \geq 1\}$$

$Weights = \{1 | i \in V \wedge V[i] = 1\} \cup \{-1 | i \in V \wedge V[i] = 0\}$

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

$$Adapt(M, V) := \max_{k, l \in Nodes} \{i | V[i] = 1 \wedge i \in p(k, l)\},$$

where $p(k, l)$ is the longest weighted path in graph $G(M, V)$ starting from k to l .

Definition 14 (Subgraph). Given two graphs $G_1 = (V_1, E_1)$, $G_2 = (V_2, E_2)$, $G_1 \subseteq G_2$ iff:

1. $V_1 \subseteq V_2$;
2. $\forall (q_i, q_j) \in E_1, \exists$ a path in G_2 from q_i to q_j .

Theorem 5.1 (Soundness). Given a program P , Γ , μ , c_1, c_2 and σ s.t. $FreeVar(P) \subseteq dom(\sigma) \cup dom(\mu) \wedge \Gamma \subseteq FreeVar(P) \wedge \Gamma \vdash_{M, V}^{c_1, c_2} P : \Phi \implies \Psi$, for all database D , $(\sigma, \mu) \models t$ s.t. $\llbracket P \rrbracket(\sigma, \mu, t, w, D) \triangleq (\sigma', \mu', t', w', D)$, then

$$A(P) \leq Adapt(M, V)$$

Lemma 1 (Subgraph). Given a program P , Γ , μ , c_1, c_2 and σ s.t. $FreeVar(P) \subseteq dom(\sigma) \cup dom(\mu) \wedge \Gamma \subseteq FreeVar(P) \wedge \Gamma \vdash_{M, V}^{c_1, c_2} P : \Phi \implies \Psi$, for all database D , $(\sigma, \mu) \models t$ s.t. $\llbracket P \rrbracket(\sigma, \mu, t, w, D) \triangleq (\sigma', \mu', t', w', D)$, then

$$G(P, D) \subseteq G(M, V)$$