

# Checks and Balances: Monitoring Data Quality Problems in Network Traffic Databases

Flip Korn  
AT&T Labs-Research  
Florham Park, NJ 02932  
flip@research.att.com

S. Muthukrishnan\*  
Rutgers University and AT&T Labs-Research  
Piscataway, NJ 08854  
muthu@cs.rutgers.edu

Yunyue Zhu†  
New York University  
New York, NY 10012  
yunyue@cs.nyu.edu

## Abstract

Internet Service Providers (ISPs) use real-time data feeds of aggregated traffic in their network to support technical as well as business decisions. A fundamental difficulty with building decision support tools based on aggregated traffic data feeds is one of data quality. Data quality problems stem from network-specific issues (irregular polling caused by UDP packet drops and delays, topological mislabelings, etc.), and make it difficult to distinguish between artifacts and actual phenomena, rendering data analysis based on such data feeds ineffective.

In principle, traditional integrity constraints and triggers may be used to enforce data quality. In practice, data cleaning is done outside the database and is ad-hoc. Unfortunately, these approaches are too rigid and limited for the subtle data quality problems arising from network data where existing problems morph with network dynamics, new problems emerge

over time, and poor quality data in a local region may itself indicate an important phenomenon in the underlying network. We need a new approach – both in principle and in practice – to face data quality problems in network traffic databases.

We propose a continuous data quality monitoring approach based on probabilistic, approximate constraints (PACs). These are simple, user-specified rule templates with open parameters for tolerance and likelihood. We use statistical techniques to instantiate suitable parameter values from the data, and show how to apply them for monitoring data quality. In principle, our PAC-based approach can be applied to data quality problems in any data feed. We present PAC-Man, which is the system that manages PACs for the entire aggregate network traffic database in a large ISP, and show that it is very effective in monitoring data quality problems.

## 1 Introduction

In the internet world, much has been discussed about various databases: database of customers and their purchases, web clicks, DNS lookups, etc. Our focus here is on network measurement databases that are in the bowels of the internet, namely, those that capture the IP traffic in the routing elements of the internet.

Unfortunately, there is a fundamental data management problem that arises in (aggregated) IP traffic logs. That is one of *data quality*. The mechanism of automatically gathering data from management stations and setting up data feeds over the internet causes

---

Work supported in part by NSF CCR 00-87022, NSF ITR 0220280 and EIA 02-05116.

†This work was started while the author was a DIMACS visitor at AT&T Labs.

*Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the VLDB copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Very Large Data Base Endowment. To copy otherwise, or to republish, requires a fee and/or special permission from the Endowment.*

*unique* data quality problems; these are not problems that arise, for example, from manual entry of customer care records. (We describe this in more detail in Section 2.) Data quality problems that stem from network-specific issues make it difficult to distinguish between artifacts and actual phenomena. In turn, this renders data analysis based on such data feeds ineffective. This is of serious concern to ISPs.

There are ways to deal with data quality problems in practice and in principle. In practice, the data is cleaned outside the database. For example, missing values may be interpolated, or the database may get populated with default values, etc. This approach is often ad-hoc, executed in scripting environments, and does not scale as new data quality problems emerge and old problems morph over time and due to network dynamics. In addition, developing scripts is time-consuming. Furthermore, synthetically cleaning the database may not be desirable because data quality problems may by themselves indicate important network phenomenon. For example, missing values may indicate loss of UDP packets which may in turn indicate high congestion levels in the network. In principle, data quality problems can be monitored using traditional integrity constraints (ICs) and triggers to enforce database integrity. Unfortunately, this approach is too rigid and limited for the subtle data quality problems arising from network data; it is not always possible to make a binary decision regarding the correctness of network data at insertion time and these tools do not offer the capability required to detect them. We need a new approach – both in principle and in practice – to face data quality problems in network traffic databases.

We propose a novel approach for detecting and monitoring data quality problems in network traffic databases. Our approach is inspired by database integrity constraints but is quite different since ICs are inadequate for network data feeds. We propose using *probabilistic and approximate constraints* (PACs). One imagines the exact constraints one would like to impose on the logical view of the network database. Then, one defines various templates that are probabilistic, approximate versions of these constraints with parametrized tolerance levels and likelihoods of violation. PACs are therefore user-driven rule templates which are specified by the DBA or analyst; the parameters are instantiated based on training data. We use statistical techniques to derive these parameters. Our data monitoring system now manages various PACs that are learned from the database and adapted over time, and which are used to monitor the quality of the data feed.

Our PAC-based approach is quite general, and can be applied in principle to any database. However, in this paper, we use the aggregate IP traffic database example to demonstrate the effectiveness of PACs. We

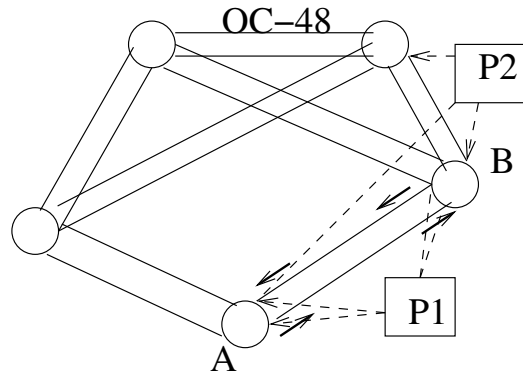


Figure 1: Network and poller.

have built a data quality detection and monitoring system for a large ISP on the entire aggregated IP traffic in their network, based on our PAC-based approach. All our experiments here are based on this real data. They show PACs are highly effective in both detecting data quality problems and in monitoring them over time. Our approach can learn stable and effective PAC rules from these real data.

The rest of the paper is organized as follows. In Section 2, we present a detailed view of network traffic databases and data quality problems. In Section 3, we provide detailed description of various types of PACs – domain, unique key and functional dependency-based PACs – all instantiated in the context of the aggregated IP traffic database. In Section 4, we describe our system, PAC-Man, that manages PACs and monitors data quality problems in network traffic databases. We also discuss the issues in integrating PAC-Man into a database in section 5. We present our experimental study in Section 6, related work in Section 7 and concluding remarks in Section 8.

## 2 Data Quality in Network Databases

In this section we describe an example network database and the associated data generation and data quality issues.

### 2.1 The Network

We focus on a large IP network such as an Internet Service Provider (ISP) where there are typically thousands of network elements configured in a topological hierarchy consisting of routing devices, each containing multiple interfaces that are endpoints for incoming and outgoing links connecting the devices. These links can have high-speed capacities, up to OC-48 (2.4 Gbps) or even OC-192 (9.6 Gbps). There can be multiple links between device pairs; traffic on any link is unidirectional. Thus, the network can be abstracted as a directed multigraph  $G(V, E)$  containing a multiset

of nodes (devices)  $V$  and a multiset of weighted edges (links)  $E$  connecting these devices; weights correspond to link speed capacities. Figure 1 depicts a toy example with links in both directions (illustrated between devices  $A$  and  $B$ ). The topology data also contains device and interface names, geographic locations, etc. We assume that snapshots of the topology data are obtained from router configuration files and provided on a regular basis (e.g., daily).

## 2.2 Network Databases: An Example

IP traffic consists of packets of data sent on links between routers, and can be captured in multiple levels. For example, there is a *physical level* traffic log which may comprise of the packet headers on the packets on each link: source, destination IP addresses, number of bytes in the packets, type of traffic, etc. Then, there is a *logical level* obtained by grouping packets on a link into logical “flows” and assembling packet level information into flow level information containing for example, source, destination IP addresses, number of packets and number of bytes in each flow. Finally, there is the *aggregated level* in which we only capture the amount of traffic — number of bytes — in each link per time interval, without regard to the source and/or destination IP addresses, type of traffic etc. As we zoom in from the aggregated level to the physical level, the amount of information in the logs increases because we get a refined look of the traffic data; but, the data sizes increase too, rather dramatically. Thus, while it is clear that finer traffic logs will reveal more about the network dynamics, few vendors, if any, provide data feeds at the physical level of packet traffic.<sup>1</sup> On the contrary, almost every router vendor supports aggregated level data logs. In fact, such aggregated traffic information is collected network-wide and used crucially in everyday network operations in today’s ISPs.

We focus on network databases that store aggregated traffic logs of ISPs. These logs are collected by network management stations using the Simple Network Management Protocol (SNMP) that is ubiquitous. SNMP data is collected by polling IP network elements. Figure 1 illustrates *pollers P1* and *P2*, which periodically request measurements at various interfaces; inbound as well as outbound traffic get polled at both endpoints of each link. These requests are sent and received using SNMP, which is implemented over UDP. UDP has less overhead than TCP because it lacks a protocol for verifying packet receipts; this can result in packet drops and delays. For simplicity, we consider only the (cyclic) count of bytes for each

interface along with a timestamp (e.g., in GMT) and an ID for the interface.

We say network traffic databases have a multigraph of topology and a data feed. The topology data is contained in two tables with the following schemas:

`DEVICE (date, ifaceID, devID, lat, long)`

specifies which device an interface is attached to and the geographic location of the device;

`LINK (date, ifaceID1, name, speed, ifaceID2)`

specifies the two endpoints(interfaces) of a link and the name and speed of the link. The data feed contains periodically-pollled inbound and outbound traffic rates and is recorded with a timestamp as well as an interface ID. The schema is

`TRAFFIC (date, time, ifaceID, inrate, outrate).`

We assume that this data is insert-only and arrives in time order, or is only slightly out of order so that it can be reordered within a reasonably sized window.

## 2.3 Data Quality Problems

The quality of SNMP-pollled data is poor for a number of reasons. The data management systems that gather aggregated traffic statistics use the underlying IP network to collect the data. So if there are problems in the network — minor problems with unreachable elements are common — then the data feed is delayed or disrupted. Similarly, these stations use the UDP protocol to collect SNMP traffic, which is unreliable: congestion in links can lead to lost UDP packets which disrupts data feeds. Also, these stations poll the network elements periodically: there are variations in the polling times because polling proceeds in some sequence, different polling stations may not be synchronized, etc. This leads to some fundamental difficulties in logically modelling the network database. For example, trivially, one would like to believe that the traffic that leaves one end of the link is the same as the traffic that reaches the other end, for any given time interval. This is difficult to ensure if the two endpoints are polled at slightly different time intervals, as frequently happens due to the sequencing delays in polling network elements. Also, there are cases where elements get polled multiple times because different pollers end up having common polling responsibility due to errors. Furthermore, configuration tables at various pollers may be out of date or erroneous, which leads to further data quality problems in the polled data. Finally, there are various counters in routers for aggregating traffic data in the links. At backbone routers where the links are exceedingly fast, these counters can wrap-around very quickly, and the reading on the counters may have to be adjusted to the proper traffic count. To summarize, the mechanism of automatically gathering data from management stations and setting up data feeds over the internet causes *unique* data quality problems;

<sup>1</sup>Special packet sniffers and monitoring tools are developed in research environments and deployed in test mode within ISPs, or special data processing systems are built in order to perform query processing in the router itself because ISP-wide datafeed of packet level traffic log is simply infeasible.

these are not problems that arise, for example, from manual entry of customer care records.

We focus on SNMP data quality problems that arise from the architecture described in Section 2.2, such as the following:

- **Missing polls:** There are either continuous or sporadic chunks of missing data which may be due to several factors: faults in a network device (e.g., failures, memory over-utilization), dropped UDP packets from SNMP polls, incorrect topology data, etc.
- **Irregular polls:** The polling rates are irregular due to UDP delays or to high CPU utilization at the network elements.
- **Extraneous polls:** A mislabeling of interfaces in the topology metadata and multiple pollers can accidentally result in double-polling.
- **Repeated values:** This often results from pre-set defaults (typically zeros), truncations resulting from exceeded capacity, or counters not being reset.
- **Violation of network traffic laws:** Disparity in traffic received from and transmitted to a link or device (à la Kirchoff’s Laws from circuits stating that the sum of the electric currents entering a node in the circuit must equal the sum of the currents exiting a node).

These observed symptoms can serve as the basis for ‘sanity’ checks for data quality control in that they indicate underlying problems. Indeed, we will employ these observations in Section 3 for designing example rules. These are many other data quality problems which can arise that are not listed here, and new ones are continually encountered.

### 3 Specific Constraints of Our Interest

Our proposed approach is based on *probabilistic, approximate constraints* (PACs). Whereas traditional ICs exactly define a set of legal values, in contrast PACs are flexible, indicating the likelihood of correctness of a given value. These constraints are expressed as a quantity over a single or multiple attributes to be satisfied within a tolerance of  $\epsilon$ . If  $\epsilon$  is zero, then the constraint is said to be *exact* and the quantity is expressed as an equality; otherwise, if  $\epsilon$  is greater than zero, then the constraint is said to be approximate. A PAC summarizes the likelihood  $\delta$  of a constraint being satisfied within tolerance  $\epsilon$  as a *cumulative (probability) distribution function* (CDF). Unlike traditional ICs, which are time-invariant, PACs are adaptive and govern temporal localities. We can specify associated time windows over which to measure and apply each constraint.

For example, a constraint could enforce a regular rate that a network device is polled every 5 minutes. However, interfaces are polled via UDP, which is prone to small and unpredictable delays. Hence, if this constraint were to be enforced using traditional ICs then, in practice, a very small portion of the data would be allowable and would thus result in false dismissals. One can attempt to get around this by relaxing this IC, for example, by permitting polling intervals of 4-6 minutes. However, this may not be deemed acceptable as it would allow consistent 4-minute polls as well as alternating polling intervals of 4 and 6 minutes. Hence, we need a way to express *degrees* of constraint satisfaction, both in terms of magnitude and rate. We give the PAC for this constraint in the following section. Below we describe three general classes of PACs (domain constraints, functional dependencies, and unique key constraints) and give examples of each. These classes are inspired by standard integrity constraints from a DBMS.

#### 3.1 Domain PACs

A *domain constraint* requires that an attribute value must be drawn from a specific set or range. The PAC analog is a straightforward generalization of this.

**Definition 1 (Domain PAC)** *Given a legal ordered domain  $D$  on an attribute, a Domain PAC specifies that all attribute values  $x$  fall within  $\epsilon$  of  $D$  with at least probability  $\delta$ , that is,*

$$Pr(x \in [D \pm \epsilon]) \geq \delta.$$

We will give a few examples of the constraints in network database and show how they are analogs of traditional ICs. The advantages of PACs will become obvious.

A simple domain constraint of interest is that traffic should be bounded by the capacity of an individual link. Thus, if  $V_i$  is the traffic *rate* between  $[t_{i-1}, t_i)$  (that is,  $V_i$  is the volume per second) and  $C$  is the total capacity of bytes that can be transferred during the unit time interval,<sup>2</sup> the constraint is  $\forall i, V_i \leq C$ . This is an exact constraint. However, in network operations, one often overengineers the link and, hence, there is a designed-for capacity of the link  $DC < C$ . One typically wants the traffic level to be less than  $DC$ . This is an approximate constraint. The rule can be expressed as follows.

**Rule 1 (Capacity Constraint)** *The traffic rate is mostly less than a designed-for capacity  $DC$ .*

Here is an effort to express this in traditional IC:

**Traditional IC 1**

$$V_i \leq DC \quad \forall i$$

<sup>2</sup> $C$  is usually constant except when there is a change in the network topology (e.g., the link speed is upgraded).

The problem of the traditional IC for expressing the constraint is that “mostly” in the rule is replaced by “always”. To overcome the problem, we will propose PAC as follows. Two ways to formalize this are based on *instantaneous* or *windowed average* terms. For example, the instantaneous PACs is as follows:

**PAC 1**

$$Pr\left(\frac{V_i}{DC} \leq \epsilon\right) \geq \delta \quad (1)$$

Another way to formalize this is in terms of the average over a window of size  $\omega$ . The PAC for this is as follows:

**PAC 1'**

$$Pr\left(\frac{avg(V_i, \omega)}{DC} \leq \epsilon\right) \geq \delta \quad \forall i$$

where  $avg(V_j, \omega)$  is the average traffic rate in  $\omega$ .

Next we explore the *smoothness* property of the data. SNMP data aggregates flows from numerous sources, so it tends to have smooth properties in general. Data quality problems often lead to violations of such smoothness. So we explore this using constraints that are designed to capture the posited smoothness of SNMP traffic. One approach might be to bound the change in traffic between consecutive polls:

**Rule 2 (Smoothness Constraint)** *The difference between consecutive traffic rates should be small most of the time.*

Let  $\bar{V}$  be the current daily average traffic rate, this rule in IC and its PAC analog are expressed as follows:

**Traditional IC 2**

$$\frac{|V_{i+1} - V_i|}{\bar{V}} \leq \mu \quad \forall i$$

where  $\mu$  is a pre-defined constant.

**PAC 2**

$$Pr\left(\frac{|V_{i+1} - V_i|}{\bar{V}} \leq \epsilon\right) \geq \delta \quad (2)$$

Note that the IC version is very difficult to express in SQL due to the sequential nature of it. Also, the IC version can not capture the probabilistic nature of the constrains as does the PAC version.

We can also measure the regularity of the polling rate as we had discussed. This involves setting up a constraint over the time field.

**Rule 3 (Polling Regularity Constraint)** *The polling rate should be at a constant rate  $\Delta$ .*

The IC and PAC for this rule could be expressed as follows.

**Traditional IC 3**

$$(t_i - t_{i-1}) = \Delta$$

**PAC 3**

$$Pr(|(t_i - t_{i-1}) - \Delta| \leq \epsilon) \geq \delta \quad (3)$$

## 3.2 Functional Dependency PACs

A *functional dependency*  $X \rightarrow Y$  enforces that two tuples must agree on the values in the set of attributes  $Y$  if they agree in attributes  $X$ .

**Definition 2 (Functional Dependency PAC)** *A Functional Dependency PAC  $X \rightarrow Y$  specifies that, if*

$$|T_i.A_\ell - T_j.A_\ell| \leq \Delta_\ell \quad \forall A_\ell \in X,$$

*then*

$$Pr(|T_i.B_\ell - T_j.B_\ell| \leq \epsilon_\ell) \geq \delta \quad \forall B_\ell \in Y.$$

We generalize the concept of functional dependency further by allowing for an *aggregate*  $f$  over a set of values in attributes  $X$  to functionally determine an *aggregate*  $g$  over a set of values in attributes  $Y$ . Notationally, we say  $f(X) \rightarrow g(Y)$ . Below we give two examples of Functional Dependency PACs.

Let  $e_1$  and  $e_2$  be the endpoints of a link. Let  $V_{out}^{e_1}$  be the outbound traffic polled at interface  $e_1$  and  $V_{in}^{e_2}$  be the inbound traffic polled at the other endpoint  $e_2$ . If the link is functioning properly, we should have the following rule.

**Rule 4 (Left-right Balance Constraint)** *The outbound and inbound traffics at the two end points of a link should be roughly the same.*

In this case, a trigger could be set to ensure that the traffic volumes at both endpoints at the same timestamps are identical.

**Traditional IC 4**

$$V_{out}^{e_1}(t_i) = V_{in}^{e_2}(t_i)$$

However, the endpoints may be polled at different times and hence this constraint is inflexible. Let  $t_i$  be a timestamp at which  $e_1$  is polled and let  $t_j$  be the closest timestamp to  $t_i$  at which  $e_2$  is polled. We give the FD PAC as follows:

**PAC 4**

$$Pr(|V_{out}^{e_1}(t_i) - V_{in}^{e_2}(t_j)| \leq \epsilon) \geq \delta \quad (4)$$

For another example, consider a network traffic version of “Kirchoff’s Law”.

**Rule 5 (In-out Balance Constraint)** *The total traffic flowing into a router should approximately equal that flowing out.*

**Traditional IC 5**

$$\sum_{in-links} V_{in}(t_i) = \sum_{out-links} V_{out}(t_i)$$

The above IC is too rigid for a network database. This constraint should be expressed as a PAC:

## PAC 5

$$Pr(|\sum_{in-links} V_{in} - \sum_{out-links} V_{out}| \leq \epsilon) \geq \delta \quad (5)$$

The sums are over traffic volumes close together in time from the same router. We can also specify a time window  $\omega$  in which these sums are computed.

### 3.3 Unique Key PACs

A *unique key constraint* says that no two tuples may agree in their values for all of the attributes that constitute a key. The PAC analog is a straightforward generalization of this.

**Definition 3 (Unique Key PAC)** *Let  $T$  be a table and  $A_\ell$  be an attribute in  $T$ . A Unique Key PAC specifies it is unlikely that more than one tuple exists with approximately the same key:*

$$Pr(|T_i.A_\ell - T_j.A_\ell| \leq \epsilon_\ell) \leq \delta_\ell$$

for each attribute  $A_\ell$  in the key.<sup>3</sup>

We can set up a unique key PAC to detect spurious traffic values whose values occur at a suspiciously high rate.

**Rule 6 (Nonrepeating Values Constraint)** *The same traffic value should not repeat too many times.*

Let  $t_i$  and  $t_j$  be consecutive polling times for an interface. A unique key IC is as follows.

#### Traditional IC 6

$$V(t_i) \neq V(t_{i+1})$$

This is infeasible because some instances of repeats may be valid (if unlikely). Hence, we use the following PAC to measure this likelihood:

#### PAC 6

$$Pr(V(t_i) = V(t_{i+1})) \leq \delta \quad (6)$$

## 4 Managing PACs: PACMAN

In this section, we describe the PAC Manager, called *PAC-Man*, that manages PACs in a network database. It has a visual interface for the DBA or user to specify PACs. Each PAC has parameters  $\epsilon$  and  $\delta$ , as we saw in the previous section, with  $\epsilon$  being the approximation to a constraint we allow and  $\delta$  being the probability for each data item to satisfy the constraint with the given approximation. PAC-Man monitors the data in the database and performs two main tasks.

- It maintains a collection of PACs and picks thresholds  $\delta^*$  and  $\epsilon^*$  suitably.

<sup>3</sup>The values of  $\epsilon_\ell$  (as well as  $\delta_\ell$ ) are independent of one another.

- It monitors new data in the database over time and *alarms* when data deviates from the PACs.

It performs all the bookkeeping activities associated with the two tasks and manages all the PACs for the given database over time. We will now discuss how PAC-Man performs its two main tasks.

Initially, the list of PACs in PAC-Man is empty. The user enters a set of rule-templates, using a visual form, whose parameters are to be instantiated as PACs over the existing training data. New PACs can be added at any time. Which PACs does PAC-Man choose to maintain and how are thresholds chosen? One may rely on the user to specify both but with thousands of elements and PACs in the database, this may be infeasible. Hence, we rely on a variety of ways to automate this process.

PAC-Man will fit the rule-template on different grouping partitions of the training data: on a per-interface basis, aggregated by interface class (e.g., OC-48s), by day-of-week (e.g., Mondays), by time-of-day (e.g., evenings), etc. In total, we consider the following group-by lattices:

- *topological*: interface  $\subset$  device pair  $\subset$  device  $\subset$  ALL;
- *categorical*: interface  $\subset$  link speed class  $\subset$  ALL;
- *temporal*: instant  $\subset$  hour  $\subset$  day  $\subset$  week  $\subset$  ALL; or  
instant  $\subset$  time-of-day  $\subset$  ALL; or  
instant  $\subset$  day-of-week  $\subset$  ALL.

Notice that the product of these hierarchies defines a lattice, and there is a well-defined notion of going up and down this lattice to refine or coarsen the granularity in one or more of these attributes. Following factors are relevant in evaluating a PAC.

- **Usefulness:** The distribution must be skewed so that it is possible to tell the good from the bad. For example, the approximation distribution from a PAC may reveal that practically any  $\epsilon$ -value will make  $\delta$  small and always be violated, or that  $\delta$  will be large and never be violated. (More formally,  $\delta$  may be 0 or 1 except for trivial values of  $\epsilon$ .) In these cases, the PAC is effectively useless.
- **Stability:** A distribution that is not stable cannot be used for a PAC rule. For example, the approximation distributions may differ substantially on each interface and on each day, which might render the PACs unstable for a given rule-template.
- **Granularity:** The distribution of a statistic depends on the sampling space. The PACs can be derived on different grouping partitions and different granularity of the data. Suppose we have

the current distribution of the data at certain granularity. If any distribution of a subpartition of the data is significantly different from the distribution, the conclusion drawn from the current distribution will not apply to that subgroup. One can have the grouping of the data in the finest granularity, but this will result in numerous distributions that are similar to each other. A better strategy is to merge these fine partitions of the data with similar distribution into coarser partitions recursively.

Currently, PAC-Man works as follows. We use a month of data (so ALL in the hierarchy above in time attribute is a month) as training data, and learn thresholds and pick a set of PACs at suitable granularity in the lattice. This is used to monitor data henceforth. That is, learning and tuning may be thought of as being *offline* with sufficient resources to perform all the computations. In particular, we have the space to store a month of data and the computations on a month of data will take less than a month.<sup>4</sup>We will only describe algorithms for implementing this framework. However a close look at the algorithms will quickly reveal that the algorithms we describe below can be implemented using recently discovered streaming algorithms in small space to a guaranteed approximation if needed, so this task can be performed by PAC-Man in a dynamic, online manner. However, we do not explicitly explore the online learning of the PACs in this work. This is because of two reasons: first, online learning of thresholds often leads to overfitting of data to local effects and further research is needed to finely tune the learning process, and second, in our preliminary study, we have not yet found this to be very useful in our SNMP database. Henceforth, we will only discuss the offline learning framework.

#### 4.1 Choosing PAC Thresholds

Based on the above criteria, PAC-Man will use some heuristic algorithms to learn and manage PACs. We will now describe how PAC-Man chooses PACs and learns thresholds in some detail. The algorithms we use are simple; in fact, our focus has been on choosing straightforward algorithms that will work effectively in our setting, rather than design sophisticated algorithms that are more general than needed.

A specific PAC rule in a particular granularity is instantiated as an *Empirical Cumulative Distribution Function* (ECDF). We maintain the percentiles (or histograms) over the approximation errors  $\epsilon$ 's to estimate their probabilities. That is, for a PAC rule, we maintain  $\epsilon_1, \epsilon_2, \dots, \epsilon_{99}$ , where  $Pr(\epsilon < \epsilon_i) \geq \frac{i}{100}$ . In other words, the empirical probability that the approximation errors of the PAC rule are less than  $\epsilon_i$  is more

than  $i\%$ . We have observed that, for the PACs used in our application, only the tails in the ECDFs are useful. We care about only those points where the approximation errors are unacceptable. PAC-Man can therefore maintain only the higher percentiles in the tails, say,  $\epsilon_{90}, \epsilon_{91}, \dots, \epsilon_{99}$ . In general, the threshold (here, 90) can be set based on the  $\delta$  one considers the absolute lower bound on any reasonable threshold. PAC-Man picks thresholds from ECDFs. Typical ECDFs have a natural “knee” or “sweet spot” where choosing  $(\epsilon^*, \delta^*)$  in various PACs is preferable. This is because increasing tolerance  $\epsilon^*$  does not change  $\delta^*$  significantly. The heuristic algorithm PAC-Man uses to decide the sweet spot  $\epsilon_s$  is based on finding the largest difference in the list  $\epsilon_{90}, \epsilon_{91}, \dots, \epsilon_{99}$ . That is,

$$s = \operatorname{argmax}_{90 < i \leq 99} (\epsilon_i - \epsilon_{i-1}).$$

This simple rule turns out to be very effective in practice, as is confirmed in our experimental study in section 6. Notice that we could have actually taken the entire training data, sorted it and used the precise ECDF to make the threshold selection, but in all our experience, the simpler heuristic above of maintaining “equiwidth histogram” on the tail of the ECDF sufficed.

**Stability of PACs.** We also need algorithm to decide the stability of the PACs over time. This can be done in a number of ways dependent on comparing the ECDFs over different windows of a given granularity. Here, we will present a simple heuristic we use that proves adequate; more sophisticated comparison techniques can be used if necessary. PAC-Man achieves this stability by monitoring the changes of the ECDFs. First we define the differences between two ECDFs  $\xi_1 = (\epsilon_{90}, \epsilon_{91}, \dots, \epsilon_{99})$  and  $\xi_2 = (\hat{\epsilon}_{90}, \hat{\epsilon}_{91}, \dots, \hat{\epsilon}_{99})$  to be their  $L_\infty$ -norm distance:

$$d(\xi_1, \xi_2) = \max_{90 \leq i \leq 99} |\epsilon_i - \hat{\epsilon}_i|.$$

The historical incremental differences between ECDFs for a particular PAC are stored as a vector  $\vec{d}$ . When a new ECDF is updated based on the new data, PAC-Man checks the change in ECDFs,  $d_{new}$ . If  $d_{new}$  exceeds  $\operatorname{avg}(\vec{d}) + k \cdot \operatorname{std}(\vec{d})$ , then PAC-Man knows that the ECDF is not stable, where  $\operatorname{avg}(\vec{d})$  and  $\operatorname{std}(\vec{d})$  are the average and standard deviation of  $\vec{d}$  respectively, and  $k$  is a predefined parameter. Given the definition of the distance between ECDFs above, we can also check if a group of ECDFs are close to each. That is, given a predefined parameter  $\eta$ , we say that a group of ECDFs  $\xi_1, \xi_2, \dots, \xi_k$  are close to each if the following holds:

$$\forall 1 \leq i, j \leq k, d(\xi_i, \xi_j) < \eta.$$

#### 4.2 Overview of PAC Selection

Combining the above algorithms, PAC-Man can manage a large number of PAC rules automatically. PAC-Man starts by constructing ECDF for each PAC rule

<sup>4</sup>This is trivial since we can now do this on a standard laptop for the entire SNMP database from a large ISP.

in the finest granularity. Then, we perform a traversal up (and down) the granularity lattice. PAC-Man will merge the ECDFs of the same PAC rule in the same granularity to a coarser granularity if they are close to each other. Also those ECDFs that are not stable will be separated into finer granularity to check whether stability can be achieved. Those PAC rules with unstable ECDFs will not be used. For PAC rules with stable ECDFs, the sweet spots will be computed and be set as their thresholds. In this manner, PAC-Man automates the process of learning set of PACs and their thresholds for data quality monitoring.

### 4.3 Alarming

The second major task that PAC-Man performs is “alarming”. The rationale is that the thresholds associated with PACs were chosen to be indicative of the data distribution, so any significant deviation from it is an occasion for an alarm. Alarming is most naturally based on the thresholds learnt. In the most general scenario, users may choose a pair  $(\alpha, \beta)$  of *alarming thresholds* where  $\alpha$  and  $\beta$  are fractions. A PAC with thresholds  $(\epsilon^*, \delta^*)$  leads to an alarm if the number of tuples in the window at specified granularity of the PAC that violate the constraint with approximation  $\alpha\epsilon^*$  is at least  $(1 - \delta)/\beta$ . Intuitively  $\alpha$  and  $\beta$  are fractions that are allowances for alarming over the thresholds that represent the data distribution. In practice, one does not want to choose these fractions for each of the PACs monitored by the PAC-Man. In our experience, it suffices to set  $\alpha = 1$  by default because the approximation thresholds are fairly robust. In our experiments we studied the effect of gradually changing  $\beta$ , but our experience is that setting  $\beta$  to be close to 0.5 is a good choice since often when violations occur, they occur in significant fraction of tuples. The online monitoring algorithm is trivial, since when a new item arrives, we can determine in  $O(1)$  time if it contributes to a violation or not with simple arithmetic to check the approximation thresholds.

## 5 Integrating PACMAN into a Database

As described so far, PAC-Man is a database application. It interacts with the database in a straightforward manner. The underlined database (in our case, Daytona) handles the transactions and manages the SNMP data. PAC-Man manages the satellite data (configuration files, routing details etc.) and perform the analysis. So, as such, PAC-Man works like a data quality browser and monitor. In this section, we discuss how to integrate PAC-Man into a database. There are several issues to consider that we describe below.

An important issue concerns specification of PACs. PAC-Man currently supports a few standard templates for specifying the checks and balances via a graphi-

cal form to be filled in. This will do as an application. However, if it is integrated into a database, one needs a systematically designed language to specify the checks. Integrity constraints which are our underlying inspiration are either incorporated into SQL or as special purpose TRIGGERS. To exploit the full power of PACs, we need to carefully design an extension of TRIGGERS. We leave this issue open.

Another issue concerns violations. In classical context of ICs, a violation automatically leads to an exception, so the data item is not inserted into the database. However, automatic datafeeds such as network traffic data come at a great rate, and contain many data quality problems both on a regular basis as well as in transitory basis. It is unrealistic to fjord such a stream so that all updates pass data quality checks. Instead, a practical approach is to populate the database with the logs and postprocess the database appropriately. There are multiple ways to postprocess the database for addressing data quality problems concerns.

One approach is *data cleaning*, which is reasonable in many traditional settings [4]. But this is not always appropriate in our setting. For example, one suggestion is to extrapolate the missing values and populate the database with these premised values. However missing polls may be an indication of network anomaly and cleaning this out of the database greatly affects the accuracy of the data analysis for network management.

A different approach is to abstract the burden away from the user and to rewrite queries automatically. The database uses the data quality monitor such as PAC-Man to understand the data, and automatically rewrites users’ queries into “equivalent” ones. This has an intuitive appeal. But it hides major conceptual and technical problems and is highly nontrivial.

The final approach is to involve the user in the loop. We use PAC-Man to expose the data quality issues to the user through a visual interface. An informed user can pose database queries more carefully. This may seem like punting the problem, and putting the entire burden on the user. But we have routinely found that this works. Let us give an example. Say a user wants to determine the correlation between the traffic at two links on a given day. The user who is made aware of the fact that there are several missing polls may modify the query as follows: *compute the correlation only restricted to observations which were made in the system*. This approach is typically not a favorite in academic research, but it is successful in real systems. Engaging a knowledgeable and thoughtful user is facilitated by data quality monitors such as PAC-Man. This is what inspired us to develop PAC-Man as a database application, merely monitoring data quality and reporting alarms.



## 6 Experiments

We tested our PAC-based approach on a real network traffic database of a large ISP. Our system is fully operational on over more than 6 months of data, with a rate of approximately 128 Mb of data arriving per day. In all of our experiments, we used a (landmark) windows of one month’s worth of data for training and one day’s worth of data for testing; we did not consider the effects of varying these window sizes.

### 6.1 Granularity Selection

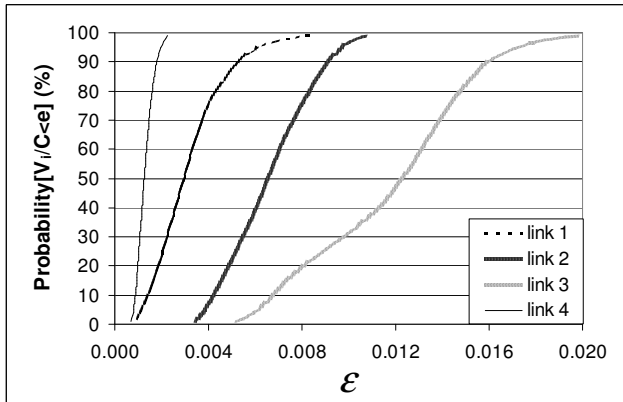


Figure 2: Distributions of PAC 1 (Capacity)

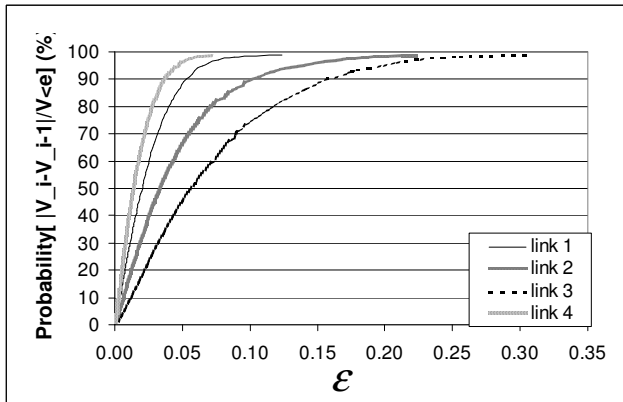


Figure 3: Distributions of PAC 2 (Smoothness)

In the first set of experiments, we examined the ECDFs of several PACs to understand their behavior at different granularities. Figure 2 shows the distributions of PAC 1 (Capacity) over all days in Month 2 for four randomly selected links. We observed that different links have different distributions. This is because each link has its own designed-for capacity; hence, a rule trained over (the aggregate of) all links would not be very meaningful. The same observation holds for PAC 2 (Smoothness) as well, as shown in Figure 3

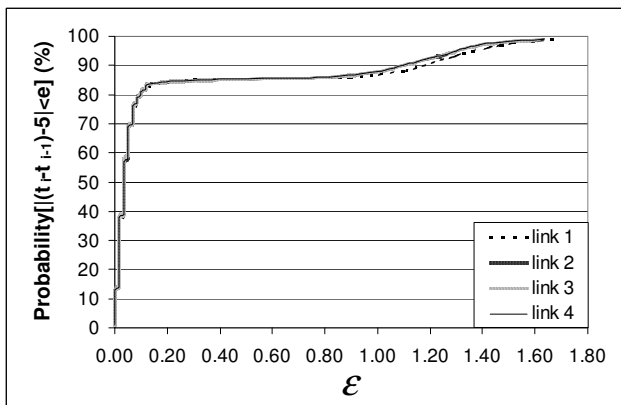


Figure 4: Distributions of PAC 3 (Polling Regularity)

for the same sample of links. On the other hand, we observed the ECDFs of different links to be roughly the same for PAC 3 (Polling Regularity), as shown in Figure 4. This makes it unnecessary to maintain a separate PAC instantiation for each link when using this rule. These examples illustrate how sensitive PACs are to topological granularity selection. In all these cases, PAC-Man automatically determined the appropriate granularity levels based on the heuristics discussed in Section 4.

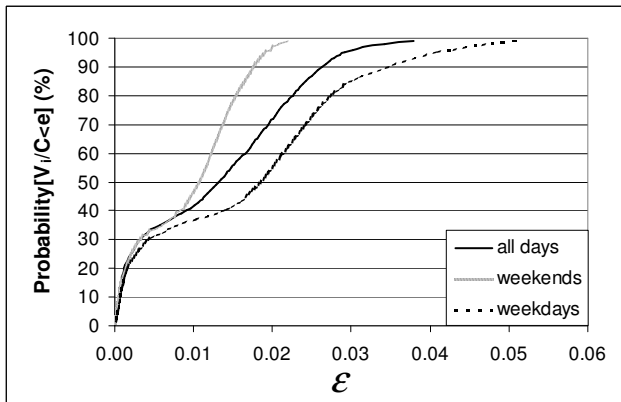


Figure 5: PAC 1 at different temporal granularities, for a single link

In addition to topological granularity, the *temporal* granularity can also affect PACs. For example, Figure 5 shows ECDFs of PAC 1 at two different temporal granularities (within Month 2) on the training data from a single link: keeping all days together in one group, and partitioning the days into two groups corresponding to weekdays and weekends. It can be seen that, at the finer granularity level, the distribution for weekdays is quite different from that for weekends. The distribution based on the (coarser) all-days granularity averages the weekend and weekday distributions together. As a result, it may not be effective

enough to detect data quality problems on a single day. Indeed, this is confirmed by considering alarms on a weekend day based on these PACs. Given the values of  $(\epsilon^*, \delta^*)$  in the respective ECDFs, the percentage of tuples that are above  $\epsilon^*$  is 47% based on the weekend granularity, but only 3% based on the all-days granularity. Hence, data quality problems may not be detected if the granularity of a PAC is not partitioned into weekdays and weekends.<sup>5</sup> Based on the heuristics from Section 4, PAC-Man appropriately chose the finer temporal granularity of partitioning into weekday/weekend for PAC 1.

## 6.2 Parameter Selection

Here we illustrate by example that, for the PACs we considered, there is a clear “sweet spot”, or “knee”, in the ECDF curves. This characteristic, which was evident throughout the data, is what enables useful PAC parameters for  $(\epsilon^*, \delta^*)$  to be chosen by PAC-Man. In Figure 4, the curve bends abruptly at roughly  $(\epsilon^* = 0.1, \delta^* = 0.85)$ . Note that  $\epsilon$  is a fraction of 5-minute units (300 seconds), so  $\epsilon^*$  is actually  $0.1 * 300 = 30$  seconds. This means that the polling interval is in the range [4:30, 5:30] for 85% of the polls. The sweet spots are less definitive in Figures 2 and 3, but are still evident. All of these occur in the tails of the curves for  $\delta$ -values in the range [90%,100%].

## 6.3 Alarms Based on PACs

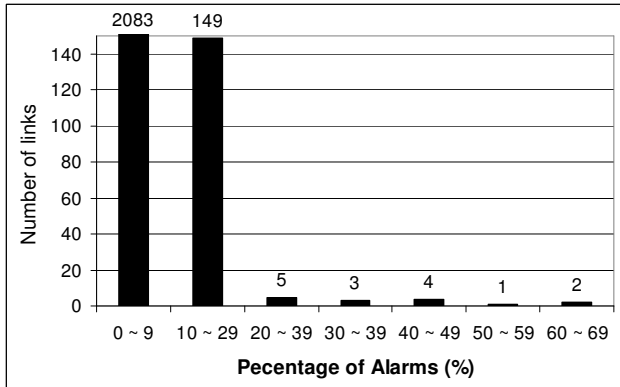


Figure 6: Percentages of violations of PAC 2 for all links in a certain day

Note that the  $\epsilon$ -values above the sweet spots spread over a wide range in Figures 2-4. This separation (of outliers) is what enables effective alarming without many false-positives. Figure 6 illustrates this further by plotting the distribution of alarm thresholds from PAC 2 after training over all links in Month 1 and

<sup>5</sup>In this case, the problem resulted from the link capacity being upgraded without a prompt change to the configuration tables.

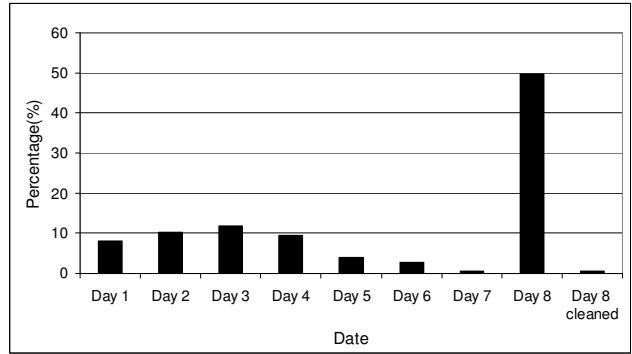


Figure 7: Percentages of violations of PAC 3 for a certain link in a certain week

applying the PACs to monitor traffic on a day from Month 2. (Recall that PAC-Man maintains a separate PAC instantiation for each link.) For different threshold ranges, it shows the number of additional alarms that would be raised if  $\delta^*$  were increased to within each range. The rapid decrease in the number of marginal alarms allows for outliers to be identified from non-outliers (due to a fairly crisp separation), thus enabling useful PACs.

## 6.4 Case Study

Here we describe an example of a data quality problem discovered using our approach. This case involves the detection of double polling. On a particular day during Month 2 (namely, Day 8), an alarm was raised by PAC 3 on an individual link. Figure 7 shows the percentage of tuples that were above the threshold  $\epsilon^*$  for that day. We can see that there is a clear spike in the percentage on Day 8. A possible reason for double polling is that two different links are polled at the same time (perhaps due to mislabeling). As a reasonable heuristic to “clean” the data, we regarded two consecutive polls occurring close together in time as coming from separate pollers; that is, the polls alternate between two different pollers. In this figure, we can see that the percentage looks normal after separating out the even and odd polls.

## 7 Related Work

Several tools exist and are widely used for monitoring network traffic such as RRDTool [8]; these tools provide graphical plots of network traffic, but put the burden on the user to detect data quality problems. Recently there has been some work in the networking literature on automatically identifying anomalies in aggregated traffic measurements using wavelets [1]; but they address specific quality problems in the network database context, and there is no underlying mechanism for detecting data quality problems in general.

There is a large body of work in the statistics literature dedicated to defining and evaluating data quality

(e.g., see [15]): outlier detection, imputation of missing values, etc. Mostly, these too tend to be focused on developing efficient, statistically grounded methods for detecting particular data quality problems.

A more principled approach for enforcing data quality in general in large scale data has been taken in the database literature. (There are a number of results in the database literature too that focus on detecting specific data quality problems; we do not survey them here because they are not directly relevant to our PAC approach here.) Traditionally, integrity constraints (ICs) and triggers are used to define deterministic constraints *a priori* as opposed to mining existing, non-deterministic rules from the data. Ad-hoc approaches to approximate functional dependencies have been defined in [7, 10, 9, 12]. Our work here is more general, abstracting from a variety of ICs, and employing a probabilistic, approximate framework in a systematic manner, to detect user-directed constraints and violations.

Recently, there has been work in providing various *browsers* for quality mining and cleaning. For example, Bellman [5] is a data quality browser that helps identify fields with similar values, estimate join sizes, tracing join paths, etc. This is useful in a database with many tables, but it does not address most of the crucial data quality problems that arise in network database context summarized in Section 2.3. Also, no principled approach akin to PACs is evident in Bellman.

The related problem of data cleaning has been studied extensively in the literature for several specific data quality problems (mapping schemas, finding approximate matches, etc.). For excellent surveys, see [3, 4, 13]; for additional references, see [11, 2]. Potter's Wheel [14], like our work here, makes the observation that data discrepancies are best dealt with using extensible, domain-specific techniques, but addresses the problem of cleaning data rather than discovering data quality violations. Ajax [6] uses a language to express data cleaning specifications declaratively. The important problems of how to express general integrity constraints and how to automatically discover such constraints were missing in these works. Our flexible framework of probabilistic ICs can provide a solid base for other data cleaning tools.

Finally, our underlying approach to use PACs (that is, to introduce tolerance and confidence parameters into ICs) is a generic approach in data mining and statistical data analysis. However, we are not aware of any application of this approach to address data quality problems.

## 8 Concluding Remarks

We have focused on network database of aggregated IP traffic data, gathered using the SNMP protocol, which is an integral part of daily operations of today's major ISPs. A fundamental problem in such databases

is one of data quality which is poor due to many network-specific issues. We have proposed a novel, principled approach to detecting and monitoring such data quality problems. Our proposal is to use probabilistic, approximate rules inspired by integrity constraints (PACs). These are specified as user-defined templates with tolerances and probability thresholds. We use the statistical properties of the data to learn the parameters. PAC-Man manages the PACs.

We have built the PAC-Man based system for a large ISP. We use this system to perform extensive experiments and show that our PAC-based approach is highly effective for detecting and monitoring data quality problems. PACs are general, powerful, logical constructs to verbalize users' insight into structural properties of the data (much as ICs are intended to be, but PACs are more flexible and powerful); hence, we believe they will prove fundamental within the general, database context beyond their effectiveness in addressing data quality problems within network databases. Future work will involve adapting PACs to study "bursty" data quality problems in a systematic manner.

## References

- [1] P. Barford, P. Kline, D. Plnka, and A. Ron. A signal analysis of network traffic anomalies. In *Proceedings of ACM SIGCOMM Internet Measurement Workshop*, 2002.
- [2] Data Cleaning Bibliography. <http://epoch.cs.berkeley.edu:8000/rshankar/pwheel/bibliography.html>.
- [3] T. Dasu and T. Johnson. *Exploratory Data Mining and Data Cleaning*. John Wiley & Sons, New York.
- [4] T. Dasu and T. Johnson. Problems, solutions and research in data quality. In *SIAM International Conference on Data Mining*, 2002.
- [5] T. Dasu, T. Johnson, S. Muthukrishnan, and V. Shkapenyuk. Mining database structure; or, how to build a data quality browser. In *SIGMOD 2002, Proceedings ACM SIGMOD International Conference on Management of Data*, 2002.
- [6] H. Galhardas, D. Florescu, D. Shasha, E. Simon, and C.-A. Saita. Declarative data cleaning: Language, model, and algorithms. In *VLDB 2001*, 2001.
- [7] R. Haux and U. Eckert. Nondeterministic dependencies in relations: An extension of the concept of functional dependency. *Information Systems*, 10(2):139–148, 1985.
- [8] RRDtool homepage. <http://www.caida.org/tools/utilities/rrdtool/>.

- [9] Wen-Chi Hou. Extraction and applications of statistical relationships in relational databases. *TKDE*, 8(6):939–945, 1996.
- [10] Ykä Huhtala, Juha Kärkkäinen, Pasi Porkka, and Hannu Toivonen. Efficient discovery of functional and approximate dependencies using partitions. In *Proceedings of the Fourteenth International Conference on Data Engineering, February 23-27, 1998, Orlando, Florida, USA*, pages 392–401. IEEE Computer Society, 1998.
- [11] Research in Data Quality Web Site. <http://www.dataquality-research.com>.
- [12] H. Mannila and K. Raiha. Algorithms for inferring functional dependencies from relations. *Data and Knowledge Engineering*, 12(1):83–99, 1994.
- [13] E. Rahm and H.H. Do. Data cleaning: Problems and current approaches. In *Data Engineering Bulletin 23(4)*, pages 3–13, 2000.
- [14] Vijayshankar Raman and Joseph M. Hellerstein. Potter’s wheel: An interactive data cleaning system. In *The VLDB Journal*, pages 381–390, 2001.
- [15] T Redman. *Data Quality: Managemetn and Technology*. Bantam Books, 1992.