

An Integrated Framework on Mining Logs Files for Computing System Management

Tao Li
School of Computer Science
Florida International University
Miami, FL 33199
taoli@cs.fiu.edu

Feng Liang
Institute of Statistics and Decision Sciences
Duke University
Durham, NC 27708
feng@stat.duke.edu

Sheng Ma
Machine Learning for Systems
IBM T.J. Watson Research Center
Hawthorne, NY 10532
shengma@us.ibm.com

Wei Peng
School of Computer Science
Florida International University
Miami, FL 33199
wpeng002@cs.fiu.edu

ABSTRACT

Traditional approaches to system management have been largely based on domain experts through a knowledge acquisition process that translates domain knowledge into operating rules and policies. This has been well known and experienced as a cumbersome, labor intensive, and error prone process. In addition, this process is difficult to keep up with the rapidly changing environments. In this paper, we will describe our research efforts on establishing an integrated framework for mining system log files for automatic management. In particular, we apply text mining techniques to categorize messages in log files into common situations, improve categorization accuracy by considering the temporal characteristics of log messages, develop temporal mining techniques to discover the relationships between different events, and utilize visualization tools to evaluate and validate the interesting temporal patterns for system management.

Categories and Subject Descriptors

I.2 [Artificial Intelligence]: Learning; I.5.4 [Pattern Recognition]: Applications

General Terms

Algorithms, Experimentation, Theory

Keywords

System Management, Log Categorization, Event Relationship, Temporal Pattern

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

KDD'05, August 21–24, 2005, Chicago, Illinois, USA.
Copyright 2005 ACM 1-59593-135-X/05/0008 ...\$5.00.

1. INTRODUCTION

When problems occur, traditional approaches for trouble shooting rely on the knowledge and experience of domain experts to figure out ways to discover the rules or look for the problem solutions laboriously. It has been estimated that, in medium and large companies, anywhere from 30% to 70% of their information technology resources are used in dealing with problems [19]. It is unrealistic and inefficient to depend on domain experts to manually deal with complex problems in ever-changing computing systems.

Modern computing systems are instrumented to generate huge amounts of system log data. The data in the log files describe the status of each component and record system operational changes, such as the starting and stopping of services, detection of network applications, software configuration modifications, and software execution errors. Analyzing log files, as an attractive approach for automatic system management and monitoring, has been enjoying a growing amount of attention. However, several new aspects of the system log data have been less emphasized in existing analysis methods from data mining and machine learning community and pose several challenges calling for more research. The aspects include disparate formats and relatively short text messages in data reporting, asynchronism in data collection, and temporal characteristics in data representation.

First, the heterogeneous nature of the system makes the data more complex and complicated [6]. As we know, a typical computing system contains different devices (e.g., routers, processors, and adapters) with different software components (e.g., operating systems, middleware, and user applications), possibly from different providers (e.g., Cisco, IBM, and Microsoft). These various components have multiple ways to report events, conditions, errors and alerts. The heterogeneity and inconsistency of log formats make it difficult to automate problem determination. For example, there are many different ways for the components to report the start up process. Some might log “the component has started”, while others might say that “the component has changed the state from starting to running”. This makes it difficult to perform automated analysis of the historical event data across multiple components when problems occur as one need to know all the messages that reflect the same status, for all the components involved in the solution [24]. To enable automated analysis of the historical event data across

multiple components, we need to categorize the text messages with disparate formats into common situations. Second, text messages in the log files are relatively short with a large vocabulary size [22]. Hence, care must be taken when applying traditional document processing techniques. Third, each text message usually contains a timestamp. The temporal characteristics provide additional context information of the messages and can be used to facilitate data analysis.

In this paper, we will describe our research efforts to address the above challenges in mining system logs. In particular, we propose to establish an integrated framework for computing system management by acquiring the needed knowledge automatically from a large amount of historical log data, possibly from different types of information sources such as system errors, resource performance metrics, and trouble ticket text records. Specifically, we will apply text mining techniques to automatically categorize the text messages into a set of common categories, incorporate temporal information to improve categorization performance, develop temporal mining techniques to discover the relationships between different events, and utilize visualization tools to evaluate and validate the interesting temporal patterns for system management.

The architecture of the proposed framework is summarized in Figure 1. The rest of the paper is organized as follows: Section 2 applies text mining techniques to categorize text messages into a set of common categories, Section 3 proposes two approaches of incorporating temporal information to improve the categorization performance, Section 4 describes our techniques for discovering the relationships between different events, Section 5 presents our experimental results, and finally Section 6 provides conclusions and discussions.

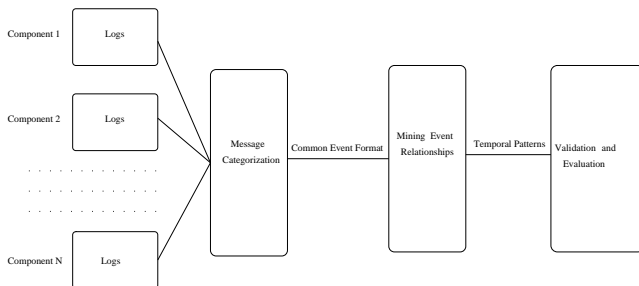


Figure 1: An Overview of the Integrated Framework.

2. SYSTEM LOG CATEGORIZATION

2.1 Common Categories

The disparate logging mechanisms impede problem investigation because of no standard approach for classifying them [24]. In order to create consistency across similar fields and improves the ability to correlate across multiple logs, it is necessary to transform the messages in the log files into a common base [2] (i.e., a set of common categories). In this paper, we first manually determine a set of categories as the basis for transformation. The set of categories is based on the CBE (Common Base Event) format established by IBM initiative [24]. CBE provides a finite set of canonical situations after analyzing thousands of log entries across multiple IBM and non-IBM products. Although we use CBE as an internal presentation here, the problem and the proposed approach are generic and can be extended for other data formats as well. Specifically, the set of categories includes *start*, *stop*, *dependency*, *create*, *connection*, *report*, *request*, *configuration*, and *other*.

2.2 Message Categorization

Given the set of common categories, we can then categorize the messages reported by different components into the prescribed categories. This can be viewed as a text categorization problem where the goal is to assign predefined category labels to unlabeled documents based on the likelihood inferred from the training set of labeled documents [21].

In our work, we use naive Bayes as our classification approach as it is among the most successful known algorithms for learning in text categorization [20]. System log files usually have a large size of vocabulary and most of the system log messages contain a free-format 1024-byte ASCII description of the event [22]. Basically Naive Bayes assumes that the text data is generated by a parametric model, and uses training data to calculate Bayes-optimal estimates of the model parameters [15]. Then, it classifies test data using the generative model by calculating the posterior probability that a class would have generated the test data in question. The most probable class is then assigned to the test data [16].

3. INCORPORATING THE TEMPORAL INFORMATION

The classification performance achieved by the Naive Bayes classifier in our practice is not satisfactory as will shown in Section 5. In many scenarios, text messages generated in the log files usually contain timestamps. The temporal characteristics provide additional context information of the messages and can be used to facilitate data analysis. As an example, knowing the current status of a network component would help forecast the possible types of messages it will generate. These structural constraints can be naturally represented using naive Bayes algorithm and hidden Markov models [7]. In this section, we describe two approaches of utilizing the temporal information to improve classification performance. Our experiments in Section 5 show that both approaches greatly improve the categorization accuracy. A preliminary report of the work has been presented in [17].

3.1 Modified Naive Bayes algorithm

To facilitate the discussion, Table 1 gives an example of the system log files extracted from a windows XP machine¹. In the example, four columns, i.e., *Time*, *EventSource*, *Msg* and *State* are listed. *Time* represents the generated timestamp of log messages, *EventSource* identifies the component which reports the message, *Msg* lists the content of text message, and *State* labels the current category.

If a sequence of log messages are considered, the accuracy of categorization for each message can be improved as the structure relationships among the messages can be used for analysis. For example, the components usually first start up a process, then stop the process at a later time. That is, the *stop* message usually occurs after the *start* message. For another example, as shown in Table 1, the component *Inventory Scanner* will report *dependency* messages if it cannot find some features or components. This happens especially just right after it generates the *create* message—“cannot create *” (* is the abbreviation of arbitrary words. It represents certain components or features here). Thus it is beneficial to take the temporal information into consideration.

In order to take the relationships between adjacent messages into account, we make some modifications to the naive Bayes algorithm. Suppose we are given a sequence of adjacent messages $D = (d_1, d_2, \dots, d_T)$. let Q_i be the category labels for message

¹The details of the data collection will be described in Section 5.

| Time | EventSource | msg | State |
|------------|-------------------|--|------------|
| 1093359583 | UPHClean | The following handles in user profile hive FIU-SCS \ lcruz01 (S-1-5-21-876885435-739206947-926709054-2036) have been closed because they were preventing the profile from unloading successfully wmioprse.exe (2432)HKCU (0x228) | report |
| 1093359583 | UPHClean | User profile hive cleanup service version 1.5.5.21 started successfully. | start |
| 1093359603 | AutoEnrollment | Automatic certificate enrollment for local system failed to contact the active directory (0x8007054b). The specified domain either does not exist or could not be contacted.Enrollment will not be performed. | dependency |
| 1093359605 | Inventory Scanner | Unable to open output file C: \invdelta.tmp. | dependency |
| 1093359606 | Inventory Scanner | LDISCN32: Can't create temporary file. | create |
| 1093359641 | Inventory Scanner | Unable to open output file C: \invdelta.tmp. | dependency |
| 1093359918 | Inventory Scanner | LDISCN32: Can't create temporary file. | create |
| 1093359583 | Userenv | | other |
| 1093359584 | MsiInstaller | Product: WebFldrs XP – Configuration completed successfully. | dependency |

Table 1: An example fraction of a log file with the manual label of each message

d_i (i.e., Q_i is one of C_1, C_2, \dots, C_L). Now we want to classify d_{i+1} . Note that

$$\begin{aligned}
 P(Q_{i+1}|d_{i+1}, Q_i) &= \frac{P(Q_{i+1}, d_{i+1}|Q_i)}{P(d_{i+1}|Q_i)} \\
 &= \frac{P(Q_{i+1}|Q_i)P(d_{i+1}|Q_{i+1}, Q_i)}{P(d_{i+1}|Q_i)}
 \end{aligned}$$

Assuming that d_{i+1} is conditionally independent of Q_i given Q_{i+1} , $P(d_{i+1}|Q_{i+1}, Q_i) = P(d_{i+1}|Q_{i+1})$. In addition, once Q_i is determined, $P(d_{i+1}|Q_i)$ is then fixed. Hence maximizing $P(Q_{i+1}|d_{i+1}, Q_i)$ is equivalent to maximizing $P(Q_{i+1}|Q_i)P(d_{i+1}|Q_{i+1})$. Observe that $P(d_{i+1}|Q_{i+1}) = \frac{P(Q_{i+1}|d_{i+1})P(d_{i+1})}{P(Q_{i+1})}$. If we assume a uniform prior on $P(Q_{i+1})$, in order to maximize $P(d_{i+1}|Q_{i+1})$, it is enough to maximize $P(Q_{i+1}|d_{i+1})$. Therefore, in summary, we assign d_{i+1} to the category Q_{i+1} which is $\text{argmax}_j(P(C_j|d_{i+1}) \times P(C_j|Q_i))$. In other words, we aim at maximizing the multiplication of text classification probability $P(C_j|d_{i+1})$ and state transition probability $P(C_j|Q_i)$. The transition probability $P(C_j|Q_i)$ can be estimated from training log files.

3.2 Hidden Markov Model

Hidden Markov Model(HMM) is another approach to incorporate the temporal information for message categorization. The temporal relations among messages are naturally captured in HMM [9]. The Hidden Markov Model is a finite set of states, each of which is associated with a (generally multidimensional) probability distribution. Transitions among the states are governed by a set of probabilities called transition probabilities [18]. In a particular state an outcome or observation can be generated, according to the associated probability distribution. The model describes a probabilistic generative process whereby a sequence of symbols is produced by starting in some state, transitioning to a new state, emitting a symbol selected by that state, transitioning again, emitting another symbol and so on until a designated final state is reached.

In our experiment, the category labels we specified in Section 2.1 are regarded as states. The emitting observation symbols are the log messages corresponding to their state labels. HMM explicitly considers the state transition sequence. It can also be used to compare all state paths from the start state to the destination state, and then choose the best state sequence that emits a certain observation sequence. So it works well in labeling continuous log messages.

When one log message has been assigned several competitive state labels by text classification, HMM selects a certain label by traversing the best state sequence. The state transition probability is calculated from the training log data sets, for example, in our experiment the probability of state *create* to state *dependency* is 0.4470. The probability of emitting messages can be estimated as the ratio of the occurrence probabilities of log messages to the occurrence of their states in the training data. Viterbi algorithm is used to find the most possible state sequence that emits the given log messages, that is, finding the most possible state labels to assign to these log messages [8].

4. MINING EVENT RELATIONSHIPS

4.1 Introduction

Once the messages in the log file are transformed into common categories, it is then possible to analyze the historical event data across multiple components to discover interesting patterns embedded in the data. Each message in log files usually has a timestamp and we will try to find the mining temporal patterns. Temporal patterns of interest appear naturally in the system management applications [4]. Specifically, a computer system problem may trigger a series of symptom events/activities. Such a sequence of symptom events provides a natural signature for identifying the root cause. As summarized in [4, 5]: a problem manifests itself as a sequence of events propagating from origin and low layer to high software layer through the dependency tree. Thus knowing the temporal patterns can help us pinpoint the root cause and take proper actions.

4.2 Notations and Problem Formulations

Once we categorize each message into a set of categories (or event types), each message can then be represented as a three-tuple $\langle s, c, t \rangle$ where s is the category label or event type, c is the source component of this log message, and t is the timestamp. In order to discover event relationships among components, we map each distinct pair of $\langle s, c \rangle$ into a unique event type e . The event sequence we obtained is then a collection of all events occurred within some time range (say, between 0 and T) $D = \{\langle e_1, t_1 \rangle, \langle e_2, t_2 \rangle, \dots, \langle e_n, t_n \rangle\}$, where $e_i \in E$, $t_i \in [0, T]$, and $t_i \leq t_{i+1}$.

As a first attempt, we focus on pairwise temporal patterns. Pairwise patterns can be easily interpreted by domain experts, and be

easily visually presented. The problem of mining event relationships can be stated as

PROBLEM 1. *Given event sequence D , find all pairwise statistically dependent patterns that can be characterized as following*

- *temporal patterns: The temporal patterns assert dependency between events and specify the timing information. Usually, they can be described as “event a happens after event b , say, about 5 minutes”.*

We refer this type of patterns as **t-patterns**.

4.3 Discovering t-Patterns

Previous work of temporal mining focuses on frequent itemsets with a predefined time window [14]. It fails to address two important aspects often required by applications. First, the fixed time-window scheme can not explore precise temporal information within a window, and misses the opportunity to mine temporal relationship longer than the window size. In our practice on system management applications, the temporal relationships discovered have time distances ranging from one second to one day. Second, as well-known for transaction data, frequent pattern framework misses significant, but infrequent patterns. In most system management applications, frequent patterns are normal operations and service disruptions are usually infrequent but significant patterns.

To address the above two problems, we develop algorithms for discovering temporal patterns without predefined time windows. The problem of discovering temporal patterns is divided into two sub-tasks: (1) using “cheap statistics” for dependence testing and candidates removal (2) identifying the temporal relationships between dependent event types. The dependence problem is formulated as the problem of comparing two probability distributions and is solved using a technique reminiscent of the distance methods used in the spatial point process, while the latter problem is solved using an approach based on Chi-Squared tests. The statistical properties provide meaningful characterizations for the patterns and are usually robust against noise. A preliminary report of the work has been presented in [10].

4.3.1 Two Interarrival Distributions

Let T_a and T_b be two point processes for event a and b respectively. We would like to check whether T_b is dependent on T_a . In this paper, we use distance methods instead of quadrat methods (e.g., time window segment). Our approach is motivated by the methods of finding spatial associations between spatial point processes developed in statistical community [1, 3]. The idea is simple: *if the occurrence of b is predictable by the occurrence of a , then the conditional distribution which models the waiting time of event type b given event type a 's presence would be different from the unconditional one.*

Define the distance from a point z to the point process T_a by $d(z, T_a) = \inf_{x \in T_a, x \geq z} (x - z)$. It is the distance from the point z to its nearest neighbor in T_a which occurs after z .

DEFINITION 1. *The unconditional distribution of the waiting time of event b is defined as $F_b(r) = \mathcal{P}\{d(x, T_b) \leq r\}$.*

The distribution can be interpreted as the probability of having event type b within time r .

DEFINITION 2. *The conditional distribution for the waiting time of event b with respect to event a is defined as:*

$$\begin{aligned} F_b^a(r) &= \mathcal{P}(d(x, T_b) \leq r | x \in T_a) \\ &= \mathcal{P}^x(d(x, T_b) \leq r) \end{aligned}$$

where \mathcal{P}^x denotes the Palm distribution [3] at an arbitrary point x with respect to the process a .

The conditional distribution can be regarded as the conditional probability distribution given there is an event of a at time x . Next we define the dependent relationship between two event types based on the two distributions.

DEFINITION 3. *Given two event types a and b , Denote their arrival times as T_a and T_b . We say that b is directly dependent on a , denoted by $a \rightarrow b$, if the two distributions $F_b(r)$ and $F_b^a(r)$ are different.*

4.3.2 Stage One: Testing the Dependences

Estimating the Two Distributions. Since the true distributions are not available, we have to first estimate them. Under the stationarity assumption, the distribution of $d(x, T_b)$ does not depend on the location of x . The assumption of stationarity is a pragmatic simplification which justifies the use of relatively simple estimation methods [3]. The two distributions can be estimated as follows: let T' be the last arrival time of event type b , i.e., $T' = b_{n_b}$. We use the information in $[0, T']$ for estimation. $F_b(r)$ can be estimated by measuring the distance $d(x, b)$ from each of several arbitrary test points to the nearest event occurrence in b and forming an empirical distribution. Similarly, $F_b^a(r)$ can be estimated by measuring the distance from each occurrence a_i of the point process a to the nearest arrival of event type b , and forming an empirical distribution function.

Let $t_i = b_i - b_{i-1}$ denote the inter-arrival time for T_b , where $i = 1, \dots, n_b$ and $b_0 = 0$. Denote $(t_{(1)}, t_{(2)}, \dots, t_{(n_b)})$ as the ordered permutation of t_i 's where $t_{(i)} \leq t_{(i+1)}$. The estimate of $F_b(r)$ is the observed proportion of distance values satisfying $d(x, T_b) \leq r$:

$$\begin{aligned} \hat{F}_b(r) &= \frac{\text{Length}\{x \in [0, T'] : d(x, T_b) \leq r\}}{\text{Length}([0, T'])} \\ &= \begin{cases} n_b r / T', & 0 \leq r \leq t_{(1)} \\ (t_{(1)} + (n_b - 1)r) / T', & t_{(1)} < r \leq t_{(2)} \\ \dots, & \dots \\ (\sum_{i=1}^{n_b-1} t_{(i)} + r) / T', & t_{(n_b-1)} \leq r \leq t_{(n_b)}. \end{cases} \quad (1) \end{aligned}$$

Let $n'_a = N_a([0, T'])$ denote the number of points in T_a within the time range $[0, T']$. For each point a_i , $1 \leq i \leq n'_a$, let d_i denote the distance from a_i to the nearest b occurring after b_i , i.e., $d_i = d(a_i, T_b)$. Let $D_{ab} = (d_1, \dots, d_{n'_a})$. Denote $(d_{(1)}, d_{(2)}, \dots, d_{(n'_a)})$ as the ordered permutation of d_i 's where $d_{(i)} \leq d_{(i+1)}$. Then $\hat{F}_b^a(r)$ is the observed proportion of distance values satisfying $d(a_i, T_b) \leq r$:

$$\begin{aligned} \hat{F}_b^a(r) &= \frac{\sum_{i=1}^{n'_a} 1\{d(a_i, b) \leq r\}}{n'_a} \\ &= \begin{cases} 0, & 0 \leq r \leq d_{(1)} \\ 1/n'_a, & d_{(1)} < r \leq d_{(2)} \\ \dots, & \dots \\ 1, & d_{(n'_a)} \leq r. \end{cases} \quad (2) \end{aligned}$$

Note that the estimate of the unconditional distribution is a piecewise linear function of r , while the estimate of the conditional one is a step function of r .

Dependency Tests On First Moments. Under the null hypothesis that there is no dependency between event type a and b , it is expected that $F_b = F_b^a$. In order to manage the computation for large scale data, we propose to test their difference in first moment. It is easy to see that if $F_b = F_b^a$, so are their first moments. But counterexamples can be constructed in which a and b are dependent but

the first moments of two distribution functions are the same. So the testing based on first moment difference is conservative and will generate false positives. However the conservative property does not affect our two-stage approach since stage one only serves as a preprocessing step to reduce the candidate space.

It is easy to check that the first moments for the two distributions $\hat{F}_b(r)$ and $\hat{F}_b^a(r)$ can be calculated as following $M_b = \frac{1}{2T'} \sum_{i=1}^{n_b} t_i^2$ and $M_b^a = \frac{1}{n_a'} \sum_{i=1}^{n_a'} d_i$. Note that M_b^a is the sample mean of d_i 's. Under the hypothesis that a and b are independent, all the d_i 's are independently distributed as $\hat{F}_b(r)$. So by Central Limiting Theorem,

$$\frac{|M_b - M_b^a|}{\sqrt{\text{var}(\hat{F}_b)/n_a'}} \sim N(0, 1)$$

where $\text{var}(\hat{F}_b)$ denotes the variance of distribution $\hat{F}_b(r)$ and is equal to

$$\text{var}(\hat{F}_b) = \frac{1}{3T'} \sum_{i=1}^{n_b} t_i^3 - \left(\frac{1}{2T'} \sum_{i=1}^{n_b} t_i^2\right)^2. \tag{3}$$

Given a confidence level α , we could compute the corresponding confidence interval $[-Z_\alpha/\sqrt{n}, Z_\alpha/\sqrt{n}]$ where z_α is the corresponding $1 - \alpha$ quantile and n is the sample size. If the value of $\frac{|M_b - M_b^a|}{\sqrt{\text{var}(\hat{F}_b)/n_a'}}$ falls outside the interval, then we conclude that the two distributions are different and hence b is dependent on a .

4.3.3 Stage Two: Identifying the Relationships

After preprocessing via “cheap statistics” in stage one, the next task is identifying the dependence between the candidate pairs and finding the waiting period (or lag) between two dependent events. Let δ be the time tolerance accounting for factors such as phase shifts and lack of clock synchronization.

DEFINITION 4. Given b is dependent on a , we say that the waiting period of b after a is p if the distance sequence $D_{ab} = (d_1, \dots, d_{n_a'})$, has a period p with time tolerance δ .

The discovery of the waiting periods is carried out using the Chi-Squared test based approach first introduced in [13]. Consider an arbitrary element τ in D_{ab} and a fixed δ . Let C_τ be the total number of elements of D_{ab} that fall into the interval $[\tau - \delta, \tau + \delta]$. Intuitively, if τ is not a period p , C_τ should be small; otherwise it should be large. The idea here is to compare C_τ with the number of elements in $[\tau - \delta, \tau + \delta]$ that would be expected from a random sequence. The procedure for identifying the relationship is essentially a one dimensional clustering process.

5. EXPERIMENTS

5.1 Log Data Generation

The log files used in our experiments are collected from several different machines with different operating systems in the School of Computer Science at Florida International University. We use logdump2td (NT data collection tool) developed by Event Mining Team at IBM T.J. Watson research center. The raw log files has 14 columns, which are sequence, time, EventType, ReportTime, Log-TypeName, LogTypeID, EventSourceName, EventSource, HostName, HostID, severity, Category, and Msg. To preprocess text messages, we remove stop words and skip html labels. We use accuracy, defined as the proportion of messages that are correctly assigned to a category, to evaluate the classification performance and use the random 30-70 data split for training and testing.

5.2 Message Categorization

The Naive Bayes classifier is built on the Bow (A Toolkit for Statistical Language Modeling, Text Retrieval, Classification and Clustering) software from Carnegie Mellon University². The implementation of HMM tool is based on the package UMDHMM version 1.02 from University of Maryland³.

The experimental results on message categorization are summarized in Table 2. In summary, our experiments show that both the modified Bayes algorithm and HMM enhance the classification accuracy and HMM achieves the best performance in our practice. However, constructing HMM requires more computational costs. Intuitively, the modified Bayes algorithm considers only adjacent state transitions while HMM considers all possible state paths. The choice between the two approaches for incorporating temporal information is problem-dependent in practice.

| Methods | Accuracy |
|----------------------|----------|
| Naive Bayes | 0.7170 |
| Modified Naive Bayes | 0.8232 |
| HMM | 0.85 |

Table 2: Performance Comparisons on Message Categorization

5.3 Discover and Visualize Event Relationships

5.3.1 Evaluation on Two Interarrival Distributions

We perform simulation studies to evaluate the effectiveness of our approach for estimating two interarrival distributions. Three sets of experiments are conducted. These simulation experiments illustrate the relationships between temporal dependence and comparisons of two distributions.

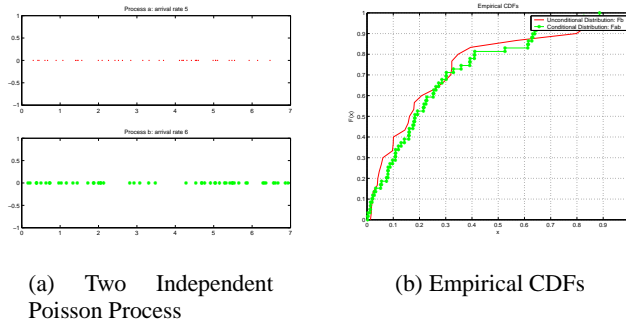


Figure 2: Two Independent Poisson Process

Independent Poisson Process : Two independent Poisson processes are generated. Figure 2(a) plots the arrivals of two generated processes and Figure 2(b) plots the two cumulative distribution functions(CDFs). As we would expect, they are very close to each other.

Branching Poisson Process : Branching Poisson process is a process model for generating dependent process types. The event sequence T_a , which is simulated by a Poisson process with arrival rate 3, is the main process. The event sequence T_b is generated as the subsidiary process with the following parameters: the number

²The tool can be downloaded at <http://www-2.cs.cmu.edu/~mccallum/bow/>.

³The package can be downloaded at <http://www.cfar.umd.edu/~kanungo/software/software.html>.

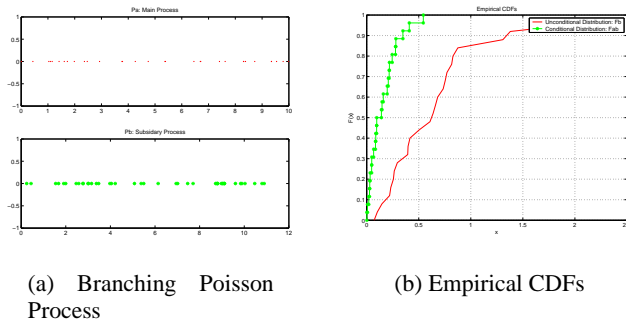


Figure 3: Branching Poisson Process

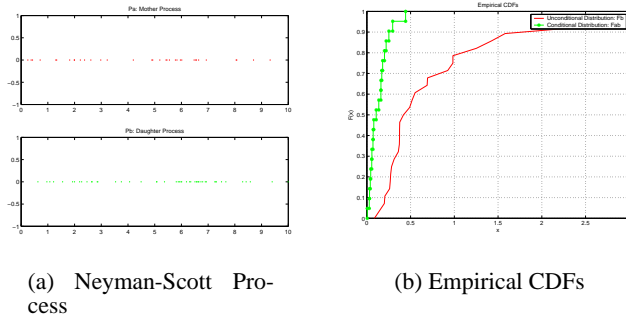


Figure 4: Neyman-Scott Process

of subsidiary events $K \sim Poisson(1)$ and the inter-arrival time between each subsidiary event $T \sim exp(0.8)$, where exp stands for exponential distribution. Figure 3(a) plots the arrivals of the two processes and Figure 3(b) plots the two CDF's. The difference between the two CDF's clearly indicates the dependence between the two processes.

Neyman-Scott Process: Neyman-scott process is a typical cluster process in which the daughter process is generated around the parent process. Like the example in branching Poisson process, we first generate the parent process T_a as a Poisson Process with arrival rate 3 and the daughter process T_b as follows: the number of daughters $\sim Poisson(1)$ and the distance between the daughter event and parent $\sim exp(0.7)$. Figure 4(a) plots the arrivals of the two processes and Figure 4(b) plots the two CDF's. The difference between the two CDF's clearly indicates the dependence between the two dependent processes.

5.3.2 Real World Study

We apply the techniques for discovering t-patterns on our log files. By reviewing the discovered patterns with operation staff, we find that many of the t-patterns we discovered are closely related to underlying problems. We discover several patterns of interest. For example, the “port up” and “port down” are temporal patterns that signify a mobile user’s login and logout; “router link down” followed by “Router down” in 10 seconds signifies a typical router down sequence: a router down results in its link down. Further, a sequence of four SNMP_requests (5, 10, 15, seconds apart in sequences) raises up a security concern of port_scans. Another pattern discovered about a performance measurement crossing critical thresholds and resetting leads to the resolution of a wrong threshold setting problem.

Visualization can help the users understand and interpret patterns in the data. For example, a scatter plot can help network operators to identify patterns of significance from a large amount of monitor-

ing data. We also use visualization tools such as EventBrowser [12] to interpret and validate interesting patterns.

6. CONCLUSION AND FUTURE WORK

In this paper, we present our research efforts on establishing an integrated framework on mining log files for computing system management by exploring the synergy of text mining, temporal data mining and visualization. There are several natural avenues for future research. First, instead of manually determining the set of common categories, we could develop techniques to automatically infer them from historical data. Second, in practice, the number of common categories for can be significantly large. To resolve this issue, it is useful to utilize the dependence relationships among different categories. Third, another natural direction is to use level-wise approach to extend the pairwise patterns to the patterns of length greater than two.

7. REFERENCES

- [1] Mark Berman. Testing for spatial association between a point process and another stochastic process. *Applied Statistics*, 35(1):54–62, 1986.
- [2] M. Chessell. Specification: Common base event, 2003. <http://www-106.ibm.com/developerworks/webservices/library/ws-cbe/>.
- [3] Noel A.C. Cressie. *Statistics for spatial data*. John Wiley & Sons, 1991.
- [4] Joseph L. Hellerstein, Sheng Ma, and Chang shing Perng. Discover actionable patterns in event data. *IBM System Journal*, 41(3):475–493, 2002.
- [5] K. Houck, S. Calo, and A. Finkel. Towards a practical alarm correlation system. *Integrated Network Management IV*, pages 226–237, 1995.
- [6] Jeffrey O. Kephart and David M. Chess. The vision of autonomic computing. *Computer*, pages 41–50, 2003.
- [7] Nicholas Kushmerick, Edward Johnston, and Stephen McGuinness. Information extraction by text classification. *Proceedings of the IJCAI-01 Workshop on Adaptive Text Extraction and Mining*, 2001.
- [8] T. R. Leek. Information extraction using hidden markov models. Master’s thesis, UC San Diego, 1997.
- [9] C. Li and G. Biswas. Temporal pattern generation using hidden Markov model based unsupervised classification. In *In Proc. of IDA-99*, pages 245–256, 1999.
- [10] Tao Li and Sheng Ma. Mining temporal patterns without predefined time windows. In *Proceedings of the Fourth IEEE International Conference on Data Mining (ICDM’04)*, pages 451–454, 2004.
- [11] Feng Liang, Sheng Ma, and Joseph L. Hellerstein. Discovering fully dependent patterns. In *SIAM DM*, 2002.
- [12] Sheng Ma and Joseph L. Hellerstein. Eventbrowser: A flexible tool for scalable analysis of event data. In *Proceedings of the 10th IFIP/IEEE International Workshop on Distributed Systems: Operations and Management*, pages 285–296. Springer-Verlag, 1999.
- [13] Sheng Ma and Joseph L. Hellerstein. Mining partially periodic event patterns with unknown periods. In *ICDE*, pages 205–214, 2001.
- [14] Heikki Mannila, Hannu Toivonen, and A. Inkeri Verkamo. Discovering frequent episodes in sequences. In *Proceedings of the First International Conference on Knowledge Discovery and Data Mining (SIGKDD’95)*, pages 210–215. AAAI Press, 1995.
- [15] A. McCallum and K. Nigam. A comparison of event models for naive Bayes text classification. In *AAAI-98 Workshop on Learning for Text Categorization*, 1998.
- [16] Tom M. Mitchell. *Machine Learning*. The McGraw-Hill Companies, Inc., 1997.
- [17] Wei Peng, Tao Li, and Sheng Ma. Mining logs files for computing system management. In *Proceedings of The 2nd IEEE International Conference on Autonomic Computing (ICAC-05)*. To appear, 2005.
- [18] L. R. Rabiner. A tutorial on hidden Markov models and selected applications in speech recognition. *Proceedings of IEEE*, 77(2):257–286, 1989.
- [19] IBM Market Research. Autonomic computing core technology study, 2003.
- [20] Irina Rish. An empirical study of the naive Bayes classifier. In *Proceedings of IJCAI-01 workshop on Empirical Methods in AI*, pages 41–46, 2001.
- [21] Fabrizio Sebastiani. Machine learning in automated text categorization. *ACM Computing Survey*, 34(1):1–47, 2002.
- [22] Jon Stearley. Towards informatic analysis of syslogs. In *Proceedings of IEEE International Conference on Cluster Computing*, Sept. 2004.
- [23] D. Stoyan, W.S. Kendall, and J. Mecke. *Stochastic Geometry and its Applications*. John Wiley and Sons, 1995.
- [24] Brad Topol, David Ogle, Donna Pierson, Jim Thoensen, John Sweitzer, Marie Chow, Mary Ann Hoffmann, Pamela Durham, Ric Telford, Sulabha Sheth, and Thomas Studwell. Automating problem determination: A first step toward self-healing computing systems. IBM White Paper, October 2003. <http://www-106.ibm.com/developerworks/autonomic/library/ac-summary/ac-prob.html>.