

Temporal Sequence Learning and Data Reduction for Anomaly Detection

TERRAN LANE and CARLA E. BRODLEY
Purdue University

The anomaly-detection problem can be formulated as one of learning to characterize the behaviors of an individual, system, or network in terms of temporal sequences of discrete data. We present an approach on the basis of instance-based learning (IBL) techniques. To cast the anomaly-detection task in an IBL framework, we employ an approach that transforms temporal sequences of discrete, unordered observations into a metric space via a similarity measure that encodes intra-attribute dependencies. Classification boundaries are selected from an *a posteriori* characterization of valid user behaviors, coupled with a domain heuristic. An empirical evaluation of the approach on user command data demonstrates that we can accurately differentiate the profiled user from alternative users when the available features encode sufficient information. Furthermore, we demonstrate that the system detects anomalous conditions *quickly* — an important quality for reducing potential damage by a malicious user. We present several techniques for reducing data storage requirements of the user profile, including instance-selection methods and clustering. An empirical evaluation shows that a new greedy clustering algorithm reduces the size of the user model by 70%, with only a small loss in accuracy.

Categories and Subject Descriptors: D.4.6 [Operating Systems]: Security and Protection

General Terms: Security, Experimentation

Additional Key Words and Phrases: Anomaly detection, clustering, data reduction, empirical evaluation, instance based learning, machine learning, user profiling

Terran Lane's work was supported, in part, by contract MDA904-97-C-0176 from the Maryland Procurement Office, and by sponsors of the Purdue University Center for Education and Research in Information Assurance and Security (CERIAS). Carla E. Brodley's work was carried out under NSF grant 9733573-IIS.

Authors' addresses: T. Lane, School of Electrical and Computer Engineering and CERIAS, Purdue University, 1285 Electrical Engineering Building, W. Lafayette, IN 47907-1287; email: terran@ecn.purdue.edu; C. E. Brodley, School of Electrical and Computer Engineering and CERIAS, Purdue University, 1285 Electrical Engineering Building, Purdue University, W. Lafayette, IN 47907-1287; email: brodley@ecn.purdue.edu.

Permission to make digital/hard copy of part or all of this work for personal or classroom use is granted without fee provided that the copies are not made or distributed for profit or commercial advantage, the copyright notice, the title of the publication, and its date appear, and notice is given that copying is by permission of the ACM, Inc. To copy otherwise, to republish, to post on servers, or to redistribute to lists, requires prior specific permission and/or a fee.

© 1999 ACM 1094-9224/99/0800-0295 \$5.00

1. INTRODUCTION

In this paper we examine the problem of anomaly detection as one of learning to characterize the behaviors of an individual, system, or network in terms of temporal sequences of discrete data. Although here we focus on user-oriented anomaly detection at the level of shell command input, the methods we present are generalizable to learning on arbitrary streams of discrete events such as GUI events, network packet traffic, or system call traces.

The anomaly-detection problem is a difficult one, especially at the level of user command traces. It encompasses a broad spectrum of possibilities, from a trusted system user turning from legitimate usage to abuse of system resources, to system penetration by sophisticated and careful hostile outsiders, to one-time use by a co-worker “borrowing” a workstation, to automated penetration launched by a relatively naive attacker via a scripted attack sequence. Time spans of interest vary from a few seconds to months. Patterns may only appear in data gathered from a number of different hosts and networks, possibly spread over thousands of miles geographically. The amount of available data can be truly staggering, as security officers may be responsible for thousands of hosts, each of which can generate megabytes of audit data per hour. Selection of data sources can also be difficult. Do the patterns of interest show themselves most clearly in command data, system call traces, network activity logs, CPU load averages, disk access patterns, or any of the hundreds of other possible sources? The patterns of interest may be corrupted by noise or interspersed with examples of normal system usage. Indeed, normal usage may vary greatly as the user changes tasks or software and learns new behaviors and command actions. Differentiating innocuous anomalies from those associated with actual abuse, misuse, or intrusion is a further difficulty. On top of all of this, a practical security system must be accurate; false alarms reduce user confidence in the system, while falsely accepting anomalous or hostile activities render the system useless.

Subsets of the general problem have been addressed by specialized techniques. Short-term (“hit and run”) attacks and attacks launched by automated scripts can often be detected by pattern matching to databases of known attack patterns (for example, Kumar [1995]; Staniford-Chen et al. [1996]). Similarly, there are numerous free and commercial programs for detecting the presence of known vulnerabilities and viruses by signatures [Farmer and Venema 1995; Gordon 1996]. In this work we address the subset of anomaly detection oriented to long-term, low profile, human-generated patterns in which known misuse signatures are insufficient to distinguish the space of possible anomalies. By long term, we mean penetration or exploitation that occurs over a period of days or weeks, rather than hit-and-run scenarios, and by low profile we mean attacks that are unlikely to flag conventional alarms (such as data theft, as opposed to data destruction or denial of service, which are more likely to be noticed through their effects). Human-generated patterns are those that occur at

the user interface level (such as the command line or GUI), rather than machine-generated events such as network packet traces, system calls, or scripted attack sequences.

The subset of anomaly detection described above encompasses not only intrusions but also hostile activities by a trusted user, and even relatively “innocuous” policy violations such as inappropriate use of system resources by an authorized user. Host-based anomaly detection at the user level can be viewed as a continuous monitoring process of the *internal* state of the system, working in conjunction with *perimeter* defenses such as passwords, firewalls, and network-based intrusion-detection systems [Heberlein et al. 1990]. Additionally, multisensor intrusion-detection systems such as AAFID [Balasubramaniyan et al. 1998] or EMERALD [Porras and Neumann 1997] allow integration of a user monitor, such as the one presented here, with other knowledge sources, to yield a comprehensive and robust view of the system state.

We take a machine-learning viewpoint of the anomaly-detection problem, in which the task is to train a classifier with known “normal” data to distinguish normal from anomalous behaviors. We focus on an *instance-based learning* (IBL) model in which query data is classified according to its relation to a set of previously encountered exemplar instances. The system we present stores historical examples of user behavior to reference when assessing the normalcy of newly encountered behavioral data.

In the rest of this paper we examine methods for representing the anomaly-detection domain as an instance-based learning task, including a temporal encoding of discrete data streams and a definition of similarity suitable for discrete temporal sequence data. We present two classes of methods for data reduction in this domain: one based on instance selection through accumulated activity statistics and one based on instance clustering. We finish with an empirical examination of performance at differentiating users under this learning scheme.

2. LEARNING FROM TEMPORAL SEQUENCE DATA

To approach anomaly detection as a machine-learning task, we must define both the learning model and representational format for the input data. One popular and highly general class of machine-learning techniques is *instance-based learning* (IBL) [Aha et al. 1991]. In this model, the concept of interest is implicitly represented by a set of instances that exemplify the concept (the instance *dictionary*). A previously unseen query instance is classified according to its relation to stored instances. A typical scheme is *k*-nearest-neighbor classification, in which a new instance is given the label of the majority of the *k* dictionary instances closest to it, where “closest” is a domain-specific measure. In continuous domains, for example, the similarity measure is often taken to be the Euclidean distance. IBL techniques may be contrasted to learning techniques that build explicit models of the data such as parametric statistical models [Fukunaga 1990], artificial neural networks [Ripley 1996], or decision trees [Quinlan 1993]. To adapt

the anomaly-detection task to the IBL learning framework, we need to choose a fixed-length vector (feature vector) representation of the data and to define the concept of “closeness,” or similarity of two vectors.

An important difference between our task and the traditional classification task of IBL is that we do not have labeled instances of multiple classes. We clearly have instances of the valid user’s normal working behaviors, but instances of hostile behaviors are a different matter—leaving aside the practical difficulty of obtaining instances of hostile activity,¹ there is an issue of coverage. Because the space of possible malicious behaviors and intruder actions is potentially infinite, it is difficult or impossible to demonstrate complete coverage of the space from a finite training corpus. Furthermore, it is often the previously unseen attack that represents the greatest threat (indeed, the very purpose of this work is to augment systems that use pattern databases to detect known threats). Finally, for reasons of privacy, it is desirable that a user-based anomaly-detection agent only employ data that originate with the profiled user or are publicly available. Releasing traces of one’s own normal behaviors, even to assist the training of someone else’s anomaly detector, runs the risk that the data will be abused to subvert the original user’s security mechanisms. Thus, we are faced with a learning situation in which only positive instances are available (where we assign positive to the class “normal behaviors” and negative to the class “anomalous behaviors”). When only positive examples are available, many standard IBL algorithms (such as the k -nearest-neighbor rule, for example) automatically classify all new examples as positive. In Section 2.2.5 we discuss a classification rule that employs only positive examples. Learning from only positive examples presents a challenge for classification, since it can easily lead to overgeneralization [Iba 1979].

A widely acknowledged difficulty with instance-based learning techniques is the overhead incurred by explicitly storing a set of class exemplars. In a dynamic environment such as anomaly detection, the size of the instance dictionary can conceivably grow without bounds, requiring data-reduction techniques to reduce the resource consumption of the IBL system. Possible solutions include removal of instances from the dictionary and rerepresentation of instances in another, less space-intensive form. In this paper we explore the use of instance selection and clustering algorithms to reduce dictionary size. Instance-selection methods use accumulated usage statistics, such as the number of times a dictionary instance is most similar to a query instance, to choose items to remove from the dictionary. In the clustering formulation, a group of similar dictionary

¹Examples (usually simulated) of machine-level attack logs (such as network packet logs or system call traces) are available, but traces of *real attacks* at the *human command* level are considerably rarer. A recent call for examples of such data by the CERIAS security research center has, to date, yielded no instances of such data. Currently such data are often not stored, not because such threats do not exist, but because adequate automated analysis tools are lacking.

instances is replaced with a single exemplar instance without regard to usage statistics.

2.1 Alternate Approaches to Sequence Learning

Many traditional approaches to learning from temporal sequence data are not applicable to the anomaly-detection domain, when the base data consists of discrete, unordered (i.e., nominal-valued) elements such as command strings. For time series of numeric values, techniques such as spectral analysis [Oppenheim and Schafer 1989], principal component analysis [Fukunaga 1990]; linear regression [Casella and Berger 1990]; linear predictive coding [Rabiner and Juang 1993]; nearest-neighbor matching, (γ, ϵ) -similarity [Bollobás et al. 1997; Das et al. 1997]; and neural networks [Chenoweth and Obradovic 1996] have proven fruitful. Such techniques typically employ a Euclidean distance, or related distance measure, defined for real-valued vectors.

There are a number of learning algorithms that are amenable to learning on spaces with nominal-valued attributes, but they typically make the assumption of the independence of attributes. For example, decision trees [Quinlan 1993] are well suited to representing decision boundaries on discrete spaces. The bias used to search for such structures generally employs a greedy search that examines each feature independently of all others. This bias ignores internal relations arising from causal structures in the data-generating process.

One method of circumventing this difficulty is to convert the data to an atemporal representation in which the causal structures are represented explicitly. Norton [1994] and Salzberg [1995] independently used such a technique to enable the learning domain to recognize coding regions in DNA fragments. DNA coding, while not temporal, does exhibit interrelations between positions that are difficult for conventional learning systems to acquire directly. The features extracted from the DNA sequences are selected by domain experts and cannot be generalized to other sequential domains. Although such an approach could be applied to the anomaly-detection domain, it would require considerable effort on the part of a domain expert, and the developed features would apply only to that data source. We are interested in developing techniques that can be applied across different data sources and tasks.

There also exist learning methods explicitly developed to model sequence data. Methods for learning the structure of deterministic finite-state automata, for example, have been widely studied [Angulin 1987; Rivest and Schapire 1989; Aslam and Rivest 1990]. DFAs, however, are not well suited to modeling highly noisy domains such as human-generated computer interface data. There exist stochastic extensions to finite-state automata such as hidden Markov models [Rabiner 1989], which are often more effective in noisy domains. We are currently investigating the application of such models to this domain.

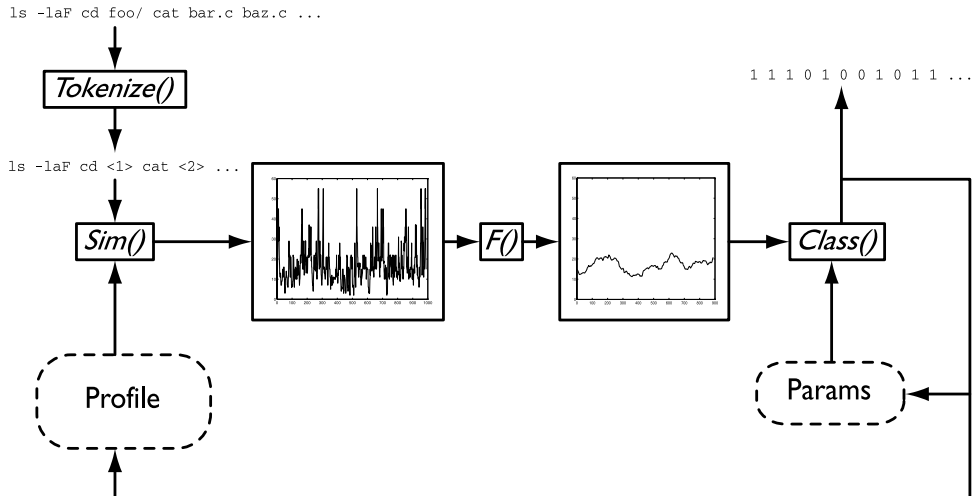


Fig. 1. Information flow in the instance-based anomaly-detection system.

2.2 The IBL Representation of Anomaly Detection

We have developed a prototype anomaly-detection system that employs an instance-based learning framework to classify observed user behaviors as normal or anomalous. We give a brief overview of the system information flow (shown schematically in Figure 1) here and describe the individual blocks in more detail below. Data enters the system, in the upper left, as an undifferentiated sequence of discrete symbols (UNIX shell command lines in this work) and is passed through a parser (`Tokenize()`) which reduces the data stream to an internal format and does preliminary feature selection. The resulting data stream is compared to the user's historical profile via a similarity measure (`Sim()`), yielding a temporal sequence of real-valued similarity measures (indicating instantaneous similarity of observed data to the profile). Because the instantaneous similarity measure stream is highly noisy, classification on the basis of this signal is difficult. To solve this problem, we introduce a noise-suppression filter (`F()`). Classification of the smoothed data stream is via a threshold decision module (`Class()`) whose decision boundaries are set through examination of an independent, parameter-selection set of the user's historical data. The final, binary class stream (upper right) is the detector's estimation of the current state of the input data (1 being "normal" and 0 "abnormal").

The components described above suffice to operate the detection system in a batch mode. In this mode the detector accumulates a single, fixed profile and employs it for all further classifications. In an operational setting we face the additional difficulty that human behaviors are dynamic and what is considered "normal" behavior is likely to change over time. This problem is known in the machine learning community as *concept drift*, and requires a system that can dynamically update the user profile and classification parameters. In such an online-learning mode, the detector employs the feedback loop shown in Figure 1 to update the profile and

classification thresholds. In this paper we are interested in examining the effects of profile data-reduction techniques in isolation. To prevent interactions with the additional complexities of the online mode feedback loop, we perform all experiments in batch mode. Elsewhere [Lane and Brodley 1998], we have examined some of the online learning issues involved with adapting user models to changing behaviors.

2.2.1 Feature Extraction. In our environment we have been examining UNIX shell command data, captured via the `(t)csH` history file mechanism. The history data are parsed by a recognizer for the `(t)csH` command language and emitted as a sequence of tokens in an internal format. Each “word” of the history data (e.g. a command name, group of command flags, or shell metacharacters) is considered to be a single token. The resulting alphabet is very large (over 35,000 distinct symbols in our complete data) and, as the frequency of some of these tokens is quite low, gathering adequate statistics over this large an alphabet is difficult. We have investigated different methods for reducing the alphabet size and found that omitting file names in favor of a filename count (e.g. `cat foo.c bar.c gux.c` is converted to `cat <3>`) greatly constrains the alphabet size (to just over 2,500 distinct tokens) and in an empirical evaluation, yields improved anomaly-detection performance.

2.2.2 The Similarity Measure. We examined several measures for computing the similarity between two discrete-valued temporal instances [Lane and Brodley 1997b]. Here we describe the measure that we found performs the best on average in empirical evaluations.

The similarity measure operates on token sequences of equal, fixed length. Although we examine only UNIX shell command data in this work, tokens may in general be any symbols drawn from a discrete, finite, unordered alphabet (e.g., GUI events, keystrokes, system calls). For a length l , the similarity between sequences $X = (x_0, x_1, \dots, x_{l-1})$ and $Y = (y_0, y_1, \dots, y_{l-1})$ is defined by the pair of functions:

$$w(X, Y, i) = \begin{cases} 0 & \text{if } i < 0 \text{ or } x_i \neq y_i \\ 1 + w(X, Y, i - 1) & \text{if } x_i = y_i \end{cases}$$

(where $w(X, Y, i) = 0$ for $i < 0$, so that $w(X, Y, 0)$ is well defined when $x_0 = y_0$) and

$$\text{Sim}(X, Y) = \sum_{i=0}^{l-1} w(X, Y, i).$$

The converse measure, *distance*, is defined to be

$$\text{Dist}(X, Y) = \text{Sim}_{\max} - \text{Sim}(X, Y)$$

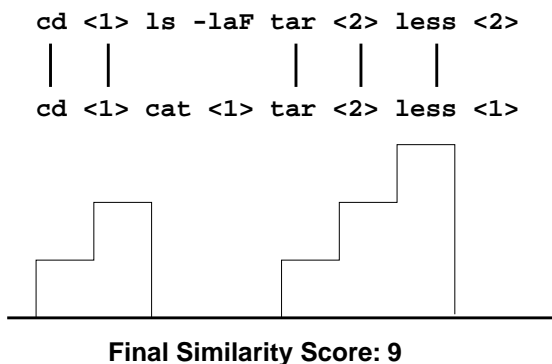


Fig. 2. Example of sequence similarity calculation. Two sequences are compared, element by element. The bottom curve represents the weight contributed by each match and the final similarity is the area under this curve.

where Sim_{\max} is the maximum attainable similarity value for a given sequence length: $\text{Sim}_{\max} = \text{Sim}(X, X)$.

An example similarity calculation is depicted in Figure 2. The function $w(X, Y, i)$ accumulates weight linearly along matching subsequences (bottom curve) and $\text{Sim}(X, Y)$ is the integral of total weight over time (area under the weight curve). In the limiting case of identical sequences, this measure reduces to $\text{Sim}_{\max} = \sum_{i=1}^l i = l(l+1)/2$. Thus, a run of contiguous matching tokens will accumulate a large similarity, while changing a single token, especially in the middle of the run, can greatly reduce the overall similarity. This measure depends strongly on the interactions between adjacent tokens as well as comparisons between corresponding tokens in the two sequences (i.e., tokens at the same offset, i , within each sequence). The sequence length l is a user-dependent parameter, explored in Lane and Brodley [1997a], where the best value was found to be dependent on the profile and on the opponent being detected. Plausible values for l are small integers in the range 8 . . . 15, and the setting $l = 10$ was found to be an acceptable compromise across users.

A user profile is a collection of sequences, \mathbf{D} , selected from a user's observed actions.² The similarity between the profile and a newly observed sequence, X , is defined to be

$$\text{Sim}_{\mathbf{D}}(X) = \max_{Y \in \mathbf{D}} \{\text{Sim}(Y, X)\}.$$

This rule is related to the 1-nearest-neighbor classification rule [Fukunaga 1990], although we are not actually performing classification at this stage, we are defining similarity to known patterns. We examined the possibility

²The problem of guaranteeing that the observed history used to profile a user actually originates with that user is critical. We do not examine this problem here, but assume that the known data are accurate.

of using an average similarity to the entire profile, but found that such a measure has much lower accuracy than the measure given here. An average across the entire profile decreases the classifier’s ability to resolve fine-structure patterns in the classification space.

2.2.3 Segmenting the Event Stream. Because the similarity measure is defined for fixed-length sequences only, it is necessary to partition the raw event stream into component subsequences. This raises the question of optimal sequence alignments: where should each sequence be defined to start? Our approach is based on the data-reduction techniques in Section 4. Initially the system segments the data stream into all possible overlapping sequences of length l (thereby replicating each token l times). Thus, every position i of the event stream is considered the starting point for a sequence of length l , referred to as the i^{th} sequence, or the sequence at time step i . For example, under sequence length $l = 6$, the tokenized data stream “cat <3> > <1> 1s -1 | more” is converted to three instance sequences “cat <3> > <1> 1s -1”, “<3> > <1> 1s -1 |”, and “> <1> 1s -1 | more”. After data reduction via instance selection or clustering (see below), the sequences remaining in the profile are considered as defining the desired alignments.

2.2.4 Noise Suppression. In practice, we found that the instantaneous similarity stream, produced by comparing an input data stream to a user profile, is far too noisy for effective classification. We attribute the high degree of noise to natural variations in the user’s actions and patterns. For example, the user may temporarily suspend writing a paper to deal with urgent incoming email, thus disrupting his or her standard paper-writing routine. Such a disruption will appear as a spuriously low similarity spike within an overall high similarity period. We therefore employ a noise-reduction filter before selecting decision thresholds or performing classification. We employ a trailing window mean-value filter, defined as

$$v_{\mathbf{D}}(j) = \frac{1}{W} \sum_{i=j-W+1}^j \text{Sim}_{\mathbf{D}}(i)$$

where $\text{Sim}_{\mathbf{D}}(i)$ is the similarity of the i^{th} token sequence to the user profile \mathbf{D} ; W is the window length; and $v_{\mathbf{D}}(j)$ is the final value of sequence j with respect to \mathbf{D} . In a comparison of the mean-value filter with a median-value filter, we found that, while the median filter is generally more effective at short window lengths ($W < 80$), performance for the two methods is approximately equivalent at longer window lengths [Lane and Brodley 1997b]. We use the mean-value filter here because we employ $W = 100$, and the mean filter can be made to run quickly more easily. We note that while a great deal of damage can be inflicted in less than the window length, such short-term attacks can be handled more readily by matching known attack signatures [Kumar and Spafford 1994]. We are primarily

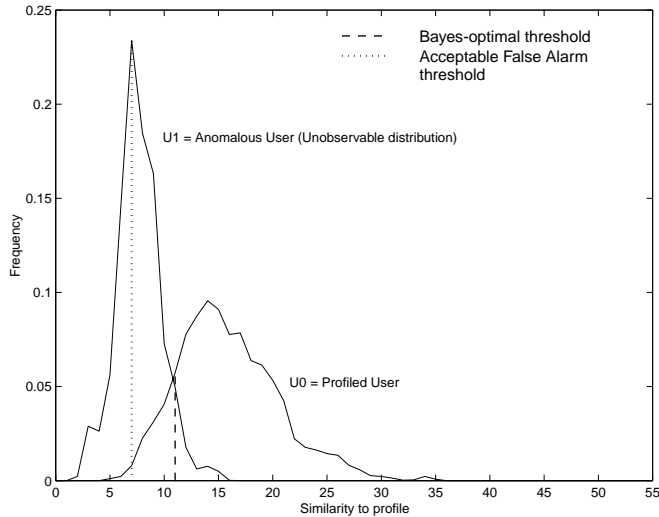


Fig. 3. Comparison of unweighted Bayes-optimal decision boundary and acceptable false alarm rate boundary. The rightmost curve (user U0) represents the profiled user.

concerned here with the class of long-term, low-profile attacks such as resource theft or industrial data theft.

2.2.5 Classification and Threshold Selection. The similarity-to-profile measure transforms the nominal vector (vectors of l discrete, unordered elements) representation into a real-valued time series on which numerical classification can be performed. Assessing the similarities of all points in a new command trace to the user profile yields a probability distribution $P_{\{T\}}$ over similarity values. Due to the structure of the similarity measure (Section 2.2.2), the range of similarity values (and thus $P_{\{T\}}$) is $0 \dots l(l + 1)/2$. When multiple classes are observable, their probability distributions can be used to construct Bayes-optimal decision boundaries [Fukunaga 1990]. In the anomaly-detection domain, however, we possess data from the profiled user only, so the Bayes-optimal boundary is unobservable to us. Furthermore, for most of the data sets we examined, the *unweighted* Bayes-optimal threshold is overly critical of the profiled user. In Figure 3, normal (U0) and anomalous (U1) similarity distributions are displayed together with the Bayes-optimal classification threshold and an alternative possible classification threshold (the acceptable false alarm threshold, described below). Sequences whose similarity to the profile falls to the right of the classification threshold are labeled normal, while points falling to the left are labeled abnormal. The area under distribution U0 and to the left of the threshold is then the false alarm probability (the probability of the valid user being falsely accused of being anomalous), while the area under distribution U1 and to the right of the threshold is the probability of falsely accepting an anomalous user. In this example, employing the unweighted

Bayes-optimal threshold for classification yields an unacceptably high false alarm rate.

In light of the above, we must seek another method for selecting a decision boundary. Conveniently, the constraints of our domain provide us with a practical heuristic: constrain the false alarm rate. This yields a Neyman-Pearson hypothesis test [Casella and Berger 1990], for a nonparametric distribution of the form:

$$class(v) = \begin{cases} 1 & \text{if } P_{\{T\}}(v) \geq r \\ 0 & \text{if } P_{\{T\}}(v) < r \end{cases}$$

where v is the similarity of a sequence to be classified to a particular profile D , 1 denotes “normal,” 0 denotes “anomalous,” and r is the specified false alarm rate. The test, as given above, does not uniquely determine decision thresholds. We implement this test by selecting two decision thresholds, t_{\max} and t_{\min} , based on the upper and lower $r / 2$ quantiles of the similarity distribution observed on a “parameter selection” data set that is independent of both training and testing data. Similarity values falling between the decision thresholds are labeled normal and those falling outside are labeled abnormal. The lower threshold, t_{\min} , detects sequences that are too different from known behaviors, while the upper threshold, t_{\max} , detects sequences that are improbably *similar* to historical behaviors, perhaps indicating a replay attack.

The parameter r selects the width of the acceptance region on the similarity-to-profile axis. Thus, it encodes the trade-off between false alarm and false accept error rates. A smaller value of r yields a wider acceptance region and corresponds to a lower false alarm rate. Simultaneously, however, anomalous values have a greater chance of falling into the broader acceptance region, so false accept rates *rise*. The choice of a particular point on the error rate trade-off curve depends on site-specific factors such as security policy and estimated cost of errors.

3. EMPIRICAL ANALYSIS: BASE SYSTEM

Before presenting methods for reducing the amount of data that must be stored in the user profile, we give an empirical performance analysis of the base classification system as described in Section 2.2. In this mode the profile contains all available training data.

3.1 Performance Criteria

We employ two methods for evaluating the performance of anomaly-detection systems. In addition to the traditional *accuracy* measurements, we argue that the *mean time to generation of an alarm* is a useful quantity to consider.

The goal in the anomaly-detection task is to identify potentially malicious occurrences, while falsely flagging innocuous actions as rarely as

possible. We denote the rate of incorrectly flagging normal behaviors as the *false alarm* rate and the rate of failing to identify abnormal or malicious behaviors as the *false acceptance* rate. Under the null hypothesis that all behavior is normal, these correspond to Type I and Type II errors, respectively. The converse accuracy rates are referred to as the true accept (ability to correctly accept the profiled user as normal) rate and the true detect (ability to correctly detect an anomalous user) rate. For the detector to be practical, it is important that the false alarm rate be low. Users and security officers will quickly learn to ignore the “security system that cries wolf” if it flags innocuous behavior too often. Finally, a practical security system must be resource-conservative in both space and time.

Detection accuracy does not, however, tell the complete story. A second important issue is *time to alarm* (TTA), which is a measure of how quickly an anomalous or hostile situation can be detected. In the case of false alarms, the time to alarm represents the expected time until a false alarm occurs. We wish the time to alarm to be short for hostile users so that they can be dealt with quickly and before doing much harm, but long for the valid user so that normal work is interrupted by false alarms as seldom as possible. We define TTA to be the mean run length of “normal” (i.e., nonalarm) classifications. This value represents the mean number of commands over which no alarm will be generated for a given user. For the valid user, this gives a sense of how long work can progress before being interrupted by a false alarm, while for an opponent this measures how long malicious use can continue before being spotted.

3.2 Data Sources and Collection

Of the literally thousands of possible data sources and features that might characterize a system or user, we chose to examine UNIX shell command data. The UNIX operating system is widely used and extensively studied in both the security and operating systems communities. The user environment is highly configurable with a rich command language, and permits a large range of possible behaviors. In the UNIX model, most user interactions take place through a command line environment (a shell), so command data is strongly reflective of user activities. Finally, there are available mechanisms to make collection of shell command data convenient in the UNIX environment.

Lacking shell traces of actual misuse or intrusive behaviors, we demonstrate the behavior of the detection system on the task of differentiating different authorized users of the UNIX hosts in the Purdue MILLENNIUM machine learning lab. In this framework, an anomalous situation is simulated by testing one legitimate user’s command data against another legitimate user’s profile. This framework simulates only a subset of the possible misuse scenarios — that of a naive intruder gaining access to an unauthorized account — but it allows us to evaluate the approach. It is to be hoped that the “naive intruder” scenario comprises a large enough fraction of all attacks to make progress in this domain of practical benefit.

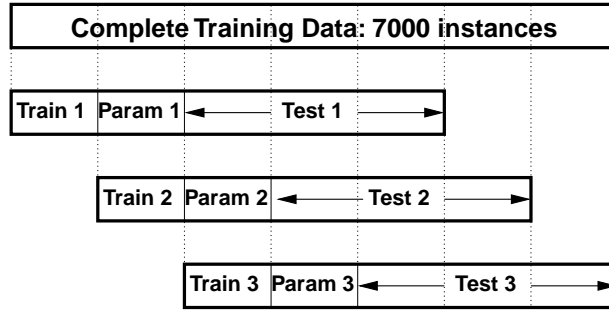


Fig. 4. Division of training data into train, parameter selection, and test data sets.

Nonetheless, we acknowledge our inability to generalize these result to broader definitions of abuses until we are able to test these techniques against real misuse data.

We acquired shell command data from eight different UNIX users over the course of more than two years. The amount of data available varies among the users from just over 15,000 tokens to well over 100,000 tokens, depending on their work rates and when each user entered and left the study. Because of computational constraints and for testing uniformity, we employ a subset of 7,000 tokens from each user, representing approximately three months of computer usage.

3.3 Experiment Structure

Because user behaviors change over time, the effective lifetime of a *static* user profile, as employed in the work described here, is limited. Thus we constructed experiments to evaluate the detector’s performance over a limited range of future activities. The separation of the 7,000 token training data into three groups (or *folds*): train, parameter selection, and test data is shown in Figure 4. The initial 1,000 tokens of each user’s data are taken as training (profile construction) data, the following 1,000 tokens for parameter-selection data (used to set the decision thresholds t_{\max} and t_{\min}), and the 3,000 following tokens to test performance for that profile. To guard against isolated data anomalies,³ three folds of train, parameter selection, and test data are produced for each user. All tests are repeated for each fold.

From each test set, a profile is constructed with $l = 10$ (the fixed sequence length for the similarity measure) and $W = 100$ (the window length for the noise-suppression filter). The resulting profile was tested against the corresponding test set for each user (a total of 8^2 test pairings). A “self” test pairing — testing the profiled user’s data against his or her

³For example, we found that our users tend to experience large behavioral changes at the beginning of academic semesters. The batch mode detection system presented here is highly sensitive to such changes.

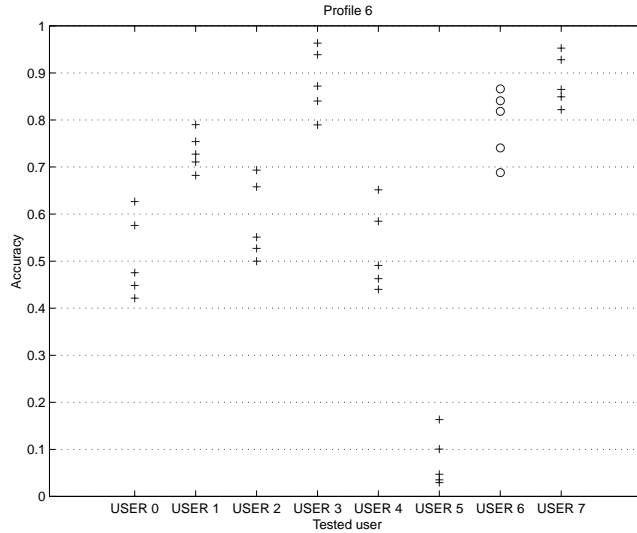


Fig. 5. Accuracy results for base system, single fold, profile for USER6. Each column shows accuracy for one test set against Profile 6. The plus symbols denote adversaries, while the circle symbols denote USER6.

own profile — allows us to examine false alarm rates, while a “nonself” pairing allows us to examine false accept rates.

The acceptable false alarm rate, r , determines how the classification thresholds, t_{\max} and t_{\min} , are set and has a substantial impact on the trade-off between false alarm and false accept errors. Because the notion of an “acceptable” false alarm rate is a site-dependent parameter, we wish to characterize the performance of the system across a spectrum of rates. We took $r \in \{0.5, 1, 2, 5, 10\}\%$, which yields a performance curve for each profile/test set pair. This curve, which expresses the trade-off between false alarm and false accept errors, with respect to r , is known as a Receiver Operating Characteristic (ROC) curve [Provost and Fawcett 1998]. A ROC curve allows the user to evaluate the performance of a system under different operating conditions or to select the optimal operating point for a given cost trade-off in classification errors.

It should be noted that this experimental structure *assumes* that the training data is “pure,” i.e., that no instances of anomalous or hostile behaviors appear in the training data. In practice, the issue of initializing the system with known behaviors is critical and the assumption of data purity may not hold. Although we do not attempt to address this issue, it has been examined by other researchers. HAYSTACK, for example, initializes users as members of hand-crafted behavioral classes [Smaha 1988].

3.4 Accuracy Results

An example of accuracy results for a test fold of a single profile (that of USER6) is shown in Figure 5. Each column in this plot displays the

accuracy results for a single test set when tested against the profile. When the test set originates with the profiled user (i.e., USER6 tested against Profile 6), the results indicate the ability to correctly identify the valid user (true accept rate). This condition is denoted with an “o” symbol on the plot. When the test set originates with a different user (e.g., USER3 tested against Profile 6), the results indicate the ability to correctly flag an anomalous condition (true detect rate). This condition is denoted with a “+” symbol on the plot. In both cases, accuracy increases in the positive direction on the Y axis. The spectrum of results in each column is generated by testing at different values of r , the acceptable false alarm rate, as described in Sections 2.2.5 and 3.3. Because r encodes the size of the acceptance region, it yields a trade-off in detect versus accept accuracies. The smallest value of r tested ($r = 0.5\%$) yields the widest acceptance region and corresponds to the highest (most accurate) point on the true accept column (USER6). But because the acceptance region is wide, more anomalous points fall into it and are accepted falsely. Thus, this value of r corresponds to the *lowest* accuracy in each of the true detect columns (USER{0-5,7}).

Profile 6 was selected to highlight a number of points. First is that accuracy is highly sensitive to the particular opponent. USER1 and USER3, for example, display quite different detection accuracies. Because of this variance, simple statistics such as mean detection accuracy are insufficient to evaluate the system’s performance. Second, although the acceptable false alarm rate parameter r was tested across the range 0.5%–10%, all of the observed false alarm rates are greater than this (13.4%–21.2%). This is a result of the training and parameterization data failing to fully reflect the behavioral distribution present in the testing data. Because the user has changed behaviors or tasks over the interval between the generation of training and testing data, the profile does not include all of the behaviors present in the test data. This phenomenon is actually exacerbated by the batch-mode experimental setup used here. In tests of the online version of this system [Lane and Brodley 1998], we found that continuously adapting to the user’s behaviors (thus shortening the delay between training and testing) improves true accept accuracy.

A third source of false accept error is demonstrated in Figure 6, where the profiled user (USER6) and USER5 have many behaviors in common—mostly “generic” account maintenance such as directory creation and file copy and remove operations. This high degree of similarity is reflected in the substantial overlaps in the similarity distributions, making differentiation impossible within this space. By contrast, USER3 is mainly engaged in programming and writing during this time. There are two possible sources for the degree of overlap between USER5 and USER6. First, the underlying observations do not encode sufficient information to distinguish the two users. Many other data sources are available for user profiling and could be used in conjunction with the techniques presented here in an

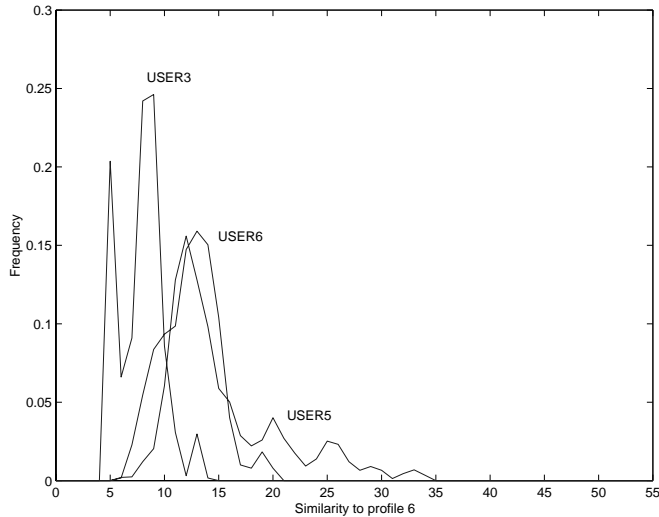


Fig. 6. False accept errors: USER5's data bears high resemblance to the profiled user's (USER6).

operational security system.⁴ The second, and more fundamental, source of error is in the similarity measure itself. The measure presented in this paper is fairly coarse (with only $O(l^2)$ possible values for a sequence length of l) and models only a single type of temporal interaction. We are currently investigating more sophisticated similarity measures such as edit distance [Cormen et al. 1992] and hidden Markov models [Rabiner 1989].

The complete set of all accuracies for all folds and all profiles is displayed in Figure 7. This plot is not intended as a reference for individual accuracy values, but to convey a sense of the general performance of the system under different operating conditions and to highlight some behavioral characteristics of the detection system. In this figure each column displays the same data as in Figure 5 — all tests against a single profiled user. Now, however, all three folds are given for each profile.

The overall impression in this plot is that many of the accuracy points are clustered toward the top of the plot, indicating that accuracy performance is generally high. The notable exception is Profile 4, which has high accuracy only for true accept (USER4 tested against Profile 4). This is an example of the decision thresholds being set to artificially extreme values, resulting in a spuriously large acceptance region. Thus the system has effectively decided that “everything is USER4,” and no real differentiation is being done — it simply accepts most behaviors as normal. Examination of USER4's training data reveals that this user appears to devote entire shell sessions to single tasks (such as the compile-debug cycle), which appear as rather repetitious and monotonous patterns. Because this user is

⁴A number of such data sources are described in Denning [1987]; Lunt and Jagannathan [1988]; Smaha [1988]; and Heberlein et al. [1990].

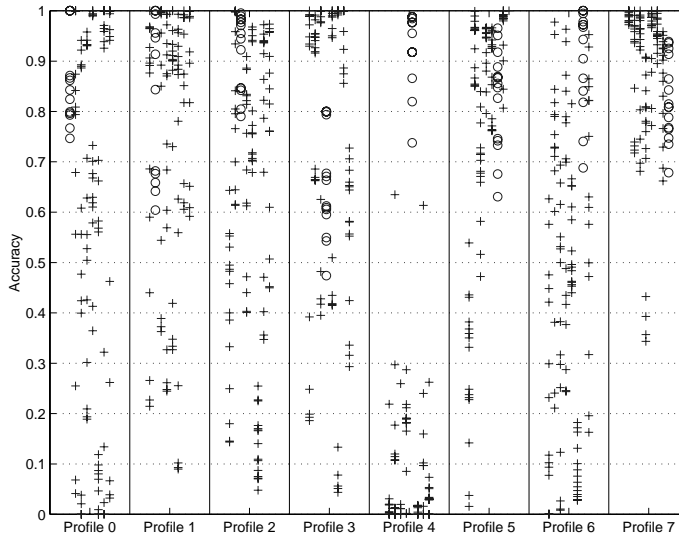


Fig. 7. Base system accuracies for all users and folds. Each column now displays a single profile tested against all test sets (the equivalent of Figure 5). Again, plus symbols denote adversary accuracies, while the circle symbols denote valid user accuracies.

working in the X-Windows environment, tasks can be assigned to single shell sessions, and those shell sessions may be long-lived (some are over 2,000 commands). So the training data may display only one or two sessions and a very small number of behaviors, while the parameter-selection data displays a different (but also small) set of behaviors. Because there may be little overlap between training and parameter selection data, the observed similarity-to-profile frequency distribution is distorted and the selected decision thresholds are poorly chosen.

A converse behavior occurs with Profile 3. Although not as dramatic as the Profile 4 case, this profile displays relatively low true accept rates (denoted by the “o” symbol) in comparison to other profiles. This is an example of the system deciding that “nothing is USER3” because the acceptance region was set too narrowly. As with USER4, this arises because different behaviors are displayed in the training and testing data. In this case the parameter-selection data reflects the training data well, but the test data is different from both of them. As a result, the acceptance range is narrowly focused to high-similarity behaviors, but the behaviors encountered in the testing data are of lower similarity.

Both of these cases can be ameliorated by online training. When the system is constantly updating the profile, the chances of missing important behaviors are smaller. In general, however, changing user behaviors present a serious problem, because sudden large changes may appear very similar to the very anomalies that the system is designed to detect. We explore these issues in Lane and Brodley [1998].

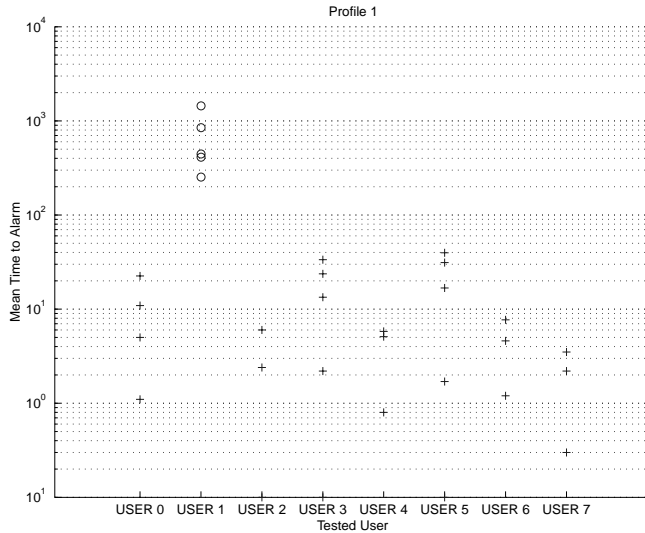


Fig. 8. Time-to-alarm results for the base system, single fold, profile for USER1. Each column shows mean time-to-alarm for one test set against Profile 1. The plus symbols denote adversaries, while the circle symbols denote USER1. Times are in token counts and the time axis is logarithmic.

3.5 Time to Alarm Results

We measure all alarm times in token counts rather than wall-clock time. Token count is more nearly correlated with the quantity of interest—how much damage can a hostile user cause before detection — than is wall-clock time.

An example of time-to-alarm results is shown in Figure 8. Analogously to Figure 5, this plot displays a single fold of tests against a single profile (USER1’s profile in this case). Each point denotes a mean time to generation of an alarm. The circle symbols represent the generation of false alarms, which we wish to be rare occurrences, while the plus symbols are true alarms, which we wish to receive as quickly as possible. Each column shows a spectrum of results, corresponding to the different settings of the acceptable false alarm rate, r (Section 2.2.5). Note that the time axis here is logarithmic and that the false alarms occur nearly an order of magnitude more slowly than do the true alarms. This difference is disproportionately greater than the accuracy difference between true accept and true detect, indicating that *false alarms occur in clusters separated by long stretches, while true alarms occur more sporadically but more often*. This is a desirable behavior, as it leads to rapid detection of adversaries, while false alarms occur rarely and in blocks, so that many false alarms can be verified at once.

The complete time-to-alarm results for the base system are shown in Figure 9. Here each column shows the results of all tests and all folds for a single profiled user. The important result in this plot is that the times to generation of false alarms (represented by the “o” symbols) are generally

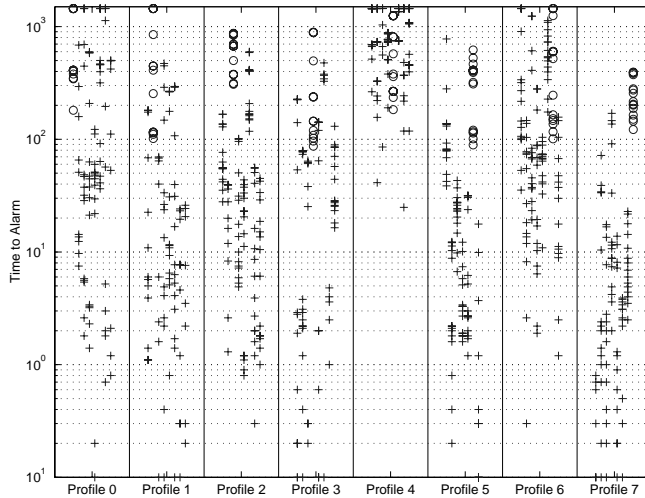


Fig. 9. Base system times-to-alarm for all users and folds. Each column displays a single profile tested against all tests sets (the equivalent of Figure 8). The plus symbols denote adversary times, while the circle symbols denote valid user times.

much longer than the times to true detections (denoted by the “+” symbols) for a given profile. The difference is often better than an order of magnitude, and is sometimes better than two. When the underlying accuracy results are poor, as in the case of USER4’s profile, the time-to-alarm results are similarly poor. In this case, since the system is marking most behaviors as normal, true alarms are generated nearly as rarely as false alarms.

4. STORAGE REDUCTION

A widely acknowledged weakness of instance-based learning algorithms is the large data storage requirement for accurate classification. A number of techniques have been examined for reducing this memory overhead, many of which are reviewed by Wilson and Martinez [1999]. In an operational setting, data reduction is critical, since the size of the profile directly impacts the time required for classification. In this section we describe two classes of techniques that we investigated for reducing the storage requirements of the user profile and give empirical results demonstrating the performance of these data-reduction methods.

4.1 Instance Selection

We note, first, that the chosen similarity measure selects only a single historical sequence as most similar to a given input sequence. If we assume that the characteristics of a user’s behavior change relatively slowly, we can invoke locality of reference to predict that recently matched dictionary sequences will be used again for detection in the near future. This suggests an analogy to tasks in operating systems, such as page replacement, in which some resources must be released in favor of others.

To examine the utility of this analogy, we implemented the least-recently-used (LRU) pruning strategy. As parameter-selection instances are examined and classified, the profile instance selected as most similar to the query instance is time-stamped. The profile is constrained to the desired size by removing the least-recently-used sequences. By analogy, we also constructed and tested the pruning heuristics FIFO (equivalent to preserving the most recently stored n sequences) and LFU (remove least frequently used sequences). Finally, we employ uniformly random instance selection as a base-line instance-selection method.

We employed the experimental structure described in Section 3.3, with the addition of a final profile size parameter S , which sets the number of sequences to remain in the user profile after instance selection. S describes the severity of the pruning to be done, and can have a substantial impact on the descriptiveness of the profile. Again, the trade-off between available system resources and desired classification performance is a site-dependent issue, so we examined the performance of each pruning method for $S \in \{100, 250, 500, 750\}$ instances.

4.1.1 Accuracy Results. Comparisons for accuracies in the performance of instance-selection methods to the base system (described in Section 3.4) are given in Figure 10. Each of these figures plots the accuracy results for one instance-selection technique on the vertical axis against those for the base system on the horizontal axis. The diagonal lines are equal-performance lines. In the region to the right of these lines, the base system has superior performance, while in the region above the diagonal the instance-selection technique has superior performance. The points shown represent all folds and tests for all profiles for a single setting of the profile size parameter S . The results displayed here are for $S = 500$ (i.e., preserve only 500 of the original thousand instances in the profile after instance selection). Other settings of S yield more-or-less extreme versions of the results displayed here, but do not change the fundamental nature of the results. The “o” symbols denote true accept accuracies (rate of correctly recognizing the profiled user) while “+” symbols denote true detection accuracies (rate of correctly distinguishing an impostor).

The first notable feature in these graphs is that overall the base system has superior true accept accuracy but inferior true detect accuracy.⁵ Accuracy impacts are a natural result of reducing the profile size because as information is removed from the profile, it becomes more likely that new query instances will fail to find a good match in the profile. This can potentially affect true accept accuracy only, while leaving true detect accuracy unchanged. The results from the instance-selection techniques, however, demonstrate a different class of behavior. Here true detect rates actually increase. In these cases the range of similarity values accepted as

⁵Although visually it may appear that the distinction is complete, the base system actually has superior true detection accuracy in approximately 16% of the cases. These points all occur at the extreme ends of the accuracy scale; the margin of victory is slight.

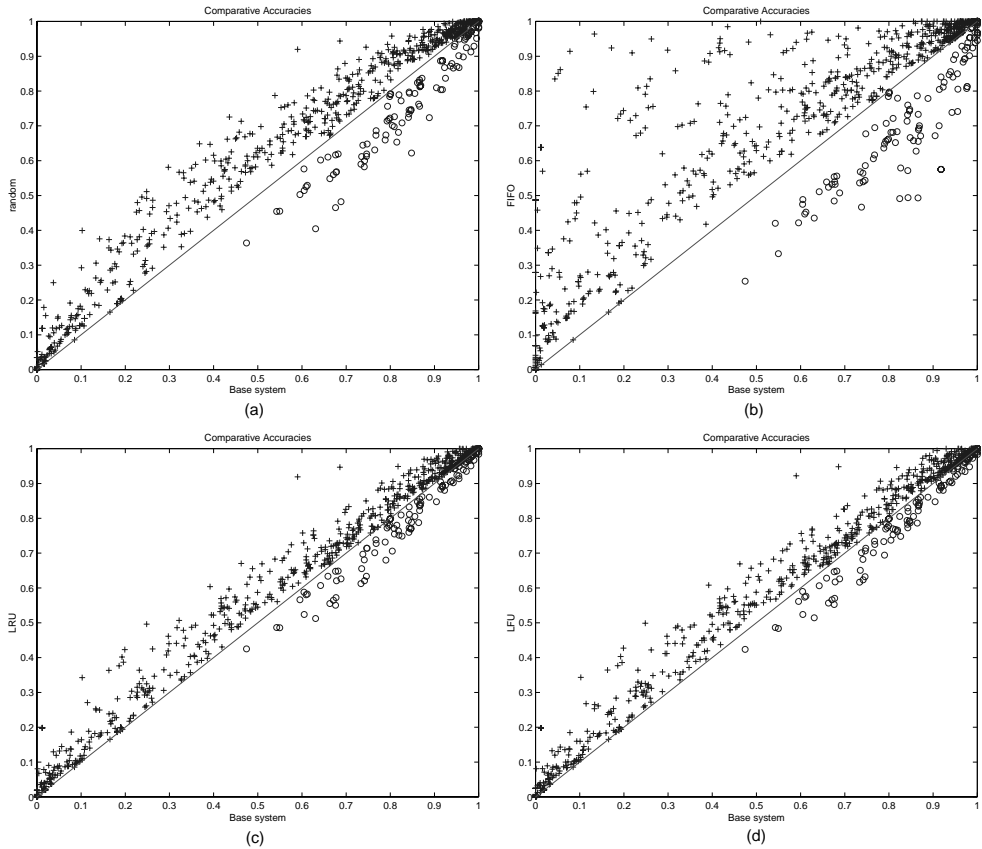


Fig. 10. Comparison of accuracy results for base system versus random (a), FIFO (b), LRU (c), and LFU (d) instance-selection strategies. Base system accuracies are plotted on the horizontal scale, while accuracies of the instance-selection techniques are plotted on the vertical scale. Points falling to the right of the diagonal indicate higher accuracies for the base system, while points above the diagonal indicate higher accuracies for the instance-selection strategy. The circle symbols denote true accept accuracies and the plus symbols are true detect accuracies.

normal is set too narrowly. This improves true detections because an adversary's behavior has a smaller chance of being similar enough to be accepted, but decreases true acceptance. The instance-selection systems biased their concepts toward rejecting more classes of behaviors. They have effectively decided that “nobody is the profiled user.”

The basic process by which accuracy is impacted is the same for all of the instance-selection techniques. The important difference between them is the ability to control informational loss. By carefully choosing the instances to preserve in the profile, the accuracy impacts can be minimized. Smaller accuracy impacts appear as clustering of the points near the iso-performance line at the diagonal. The worst technique in terms of accuracy is the FIFO method (b), which keeps only the S newest instances in the profile. Clearly, a great deal of important information is being lost by discarding

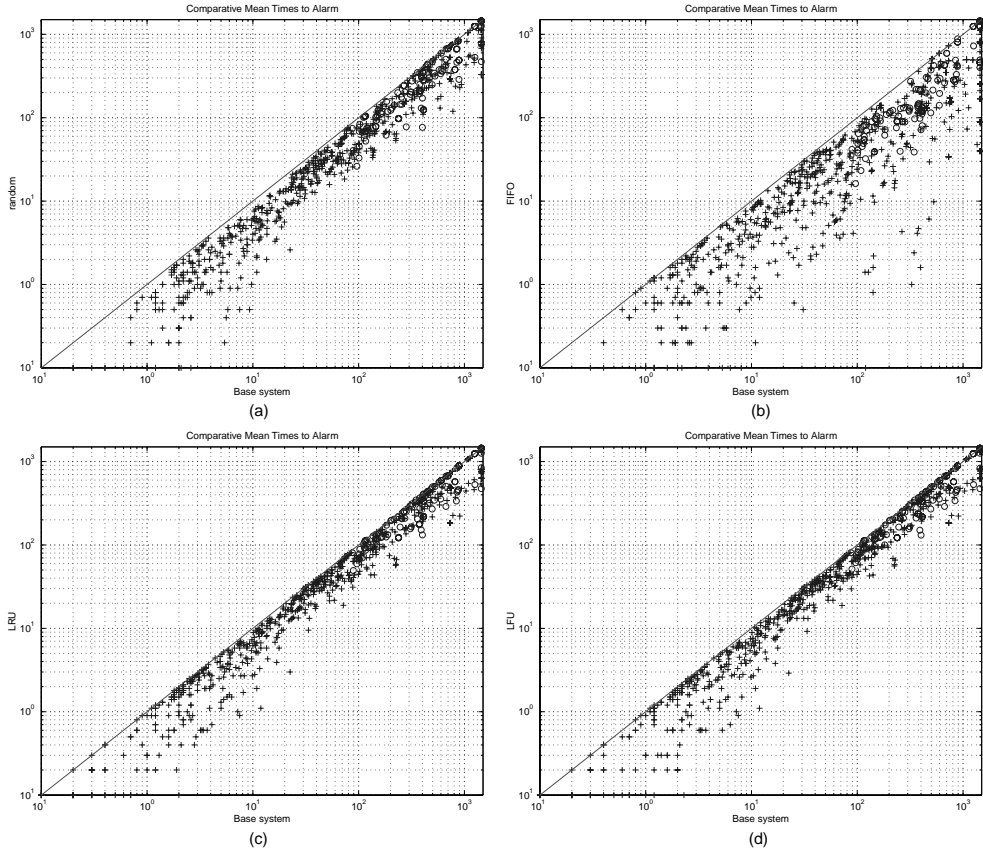


Fig. 11. Comparison of time-to-alarm results for base system versus random (a), FIFO (b), LRU (c), and LFU (d) instance-selection strategies. Base system TTAs are plotted on the horizontal scale, while TTAs of the instance-selection techniques are plotted on the vertical scale. Points falling to the right of the diagonal indicate higher TTAs for the base system, while points above the diagonal indicate higher TTAs for the instance-selection strategy. The circle symbols denote time to generation of a false alarm and the plus symbols are times to generation of true alarms. All times are in token counts, and both time axes are logarithmic.

older instances (a result that casts some doubt on the assumption of locality of reference employed by the instance-selection techniques discussed here). Simply randomly selecting S instances to keep, (a), is superior to FIFO. Both LRU (c) and LFU (d) beat random. The difference between the two is slight, even by nonvisual measures. LRU sacrifices an average of 3.6% on true accept rate to gain 3.5% on true detect rate, while LFU loses 3.7% to gain 3.4%.⁶ At smaller values of S (final profile size), the

⁶To measure the relative accuracy performance between two systems, we employ a mean of accuracy value differences. Thus the difference in true detect rates between method 1 and method 2 is $1/N \sum_{t \in \text{opponent_test_sets}} (\text{accuracy}_{\text{method1}}(t) - \text{accuracy}_{\text{method2}}(t))$ where N is the number of opponent test sets.

effect is more pronounced, but the margins are still not dramatic. At $S = 100$, for example, (reducing the profile to a tenth of its original size), the corresponding figures for LRU are 28.6% and 17.8%, compared to 23.7% and 16.1% for LFU.

4.1.2 Time-to-Alarm Results. Comparative results for the time-to-alarm measure are given in Figure 11. The plots are analogous to those in Figure 10, but give times-to-alarms for the same conditions. Times are in token counts and both time axes are logarithmic. In the region to the right of the diagonal, the base system has longer time-to-alarm, while above the diagonal the instance selection technique has longer TTA. It is apparent that instance selection decreases (or leaves unchanged) time-to-alarm in all cases. This is a consequence of the accuracy trade-off observed in the previous section. Improved ability to detect an adversary corresponds to a shorter TTA, while poorer ability to validate the profiled user also yields shorter TTA (i.e., more false alarms per unit time). Again, points close to the iso-performance diagonal correspond to minimal performance impact due to instance selection. Not surprisingly, we observe the same relations among the instance-selection techniques that we found in the last section. In particular, FIFO incurs the largest performance changes, while LRU and LFU are approximately equivalent, both beating random instance selection.

4.2 Clustering

An alternate method for reducing data storage is to modify the representation of sets of points within the data space. For example, Salzberg [1991] represented sets of points as hyper-rectangles, while Domingos [1995] induced rules that cover subsets of the instance base. We have examined techniques that attempt to locate *clusters* of similar instances (with respect to the similarity measure defined in Section 2.2.2) within the data. A cluster can be represented by a single exemplar instance, or center point, which is the instance having the smallest distance to all other instances in the cluster. By discarding all other elements of the cluster, substantial space and time savings can be realized. Although the practical effect of this process is the same as that of the instance-selection methods described above, the clustering process employs knowledge about the relationships among elements within the profile, while the pruning methods employ knowledge about the relationship between individual profile elements and the “external” parameter-selection data.

One popular class of clustering algorithms is based on the Expectation-Maximization (EM) procedure [Moon 1996]. These methods attempt to simultaneously maximize an optimality criterion across all clusters through gradient descent on the cluster likelihood space. The basic process is to assign all points to clusters and evaluate the optimality criterion under that labeling. The evaluation yields a parameter set used to reassign points to clusters. This basic loop is repeated until the optimality criterion converges to a stable point. A common implementation of this process is the

K -means algorithm, or its discrete analog, K -centers. In these algorithms, clusters are parameterized by their centroid (continuous) or center (discrete) points and radii, and the search locates K clusters.

Clustering methods based on EM are popular because they are general and often highly effective. However, when many local optima are present in the likelihood space, the quality of the solution produced can be very sensitive to the initial assignment of points to clusters. A larger difficulty for the anomaly-detection domain is that K , the the number of clusters to be sought, must be known *a priori*, yet it is not clear how to determine the number of “natural” clusters in the user behavioral data. Furthermore, the convergence rate of these methods is not guaranteed and is often related to the number of clusters sought. For a large K , the search time for a stable solution can be prohibitive.

As a response to these difficulties, we propose a greedy clustering algorithm that builds individual clusters consecutively, attempting to minimize the criterion

$$val(\mathbf{C}) = \frac{\sum_{x \in \mathbf{C}} \sum_{y \in \mathbf{C}} \text{Dist}(x, y)}{|\mathbf{C}|^2}$$

for each cluster \mathbf{C} . This measure is a generalization of the mean intercluster distance employed for clustering [Fukunaga 1990]. From an initial seed point, the cluster is grown incrementally by including the point that increases $val(\mathbf{C})$ the least, until the halting criterion is reached. Growth is halted when the value of the cluster’s criterion function reaches a local minimum. Because, in some cases, the cluster value monotonically approaches Sim_{\max} , the halting criterion we actually use is that the first derivative of $val(\mathbf{C})$ be within ϵ of 0 for some (empirically selected) value of ϵ . As each sequence is added to a cluster, it is removed from the set of available sequences. When the cluster is complete, we define the center of the cluster \mathbf{C}_{cent} to be the point possessing the minimum total distance to all other points in \mathbf{C} . The similarity between a sequence X and a cluster is then $\text{Sim}(X, \mathbf{C}_{\text{cent}})$.

In practice, we have found that this cluster-selection algorithm is somewhat too lenient — it accepts points that decrease the cluster’s effectiveness in classification. We solve this in a manner analogous to the pruning process employed in decision-tree learning [Quinlan 1993]. After growing a single cluster to completion according to the halting criterion, the clustering algorithm removes outlying points and returns them to the pool of available sequences (to make it possible for them to contribute to different clusters). Our pruning function removes points from the cluster that fall outside the cluster mean radius — i.e., points whose distance to the center is greater than the mean distance to the center of all points in the cluster. Points falling within the mean radius are discarded and the final cluster is

represented only by its center and mean radius. We realize substantial space savings by discarding all cluster elements other than the center.

The complete clustering algorithm is structurally similar to the single cluster construction algorithm. We sequentially select individual clusters by their ability to maximize the analog of mean intracluster distance:

$$\text{val}\{\mathbf{C}_1, \mathbf{C}_2, \dots, \mathbf{C}_n\} = \frac{\sum_{i=1}^n \sum_{j=1}^n \text{Dist}(\mathbf{C}_{i, \text{cent}}, \mathbf{C}_{j, \text{cent}})}{n^2}.$$

In this case we found the single cluster halting criterion to be ineffective because, typically, all of a data set's points are exhausted before the derivative of the intracluster distance approaches 0. When we allow the clustering process to absorb all available points, many of the clusters were found to either not contribute to classification accuracy or to be actively harmful. Instead, we halt the clustering process when the the minimum intercluster value of all current clusters falls below a threshold C . C determines when the clustering process will be halted and how many clusters will be constructed. A small value of C produces many clusters, while a large one allows many of the original profile points to remain. We examine the effects of the parameter C below.

4.2.1 Comparison of Clustering Techniques. We examined the performance of the K -centers and greedy clustering methods for storage reduction under the experimental conditions described in Section 3.3. The additional parameters examined were K , the number of clusters to locate for the K -centers method, and C , the global halting criterion for the greedy clustering method. For K -centers, we ran the search to 10,000 cycles or convergence, and examined $K \in \{50, 75, 100, 125, 150\}$ clusters. We note that this corresponds to a fairly extreme degree of data compression — retaining only K sequences total in the profile — but we were unable to examine the system behavior for larger K , as the search failed to converge or reach an acceptable solution in the allotted number of search cycles. For the greedy clustering method, we examined $C \in \{0.25, 0.5, 0.75\}$ and used a fixed, empirically selected $\epsilon = 0.005$.

A comparison of the accuracies of the two clustering techniques is given in Figure 12 (a). In this figure the test accuracies of the greedy clustering system with $C = 0.25$ are plotted vertically against the accuracies for the K -centers clustering system with $K = 150$ on the horizontal. The diagonal line is the iso-performance surface. In the region to the right of this line the K -centers system has superior performance, while in the region above the diagonal the greedy clustering technique has superior performance. The “o” symbols represent true acceptance rates, while the “+” symbols represent true detection rates. The corresponding times-to-alarm are given in (b).

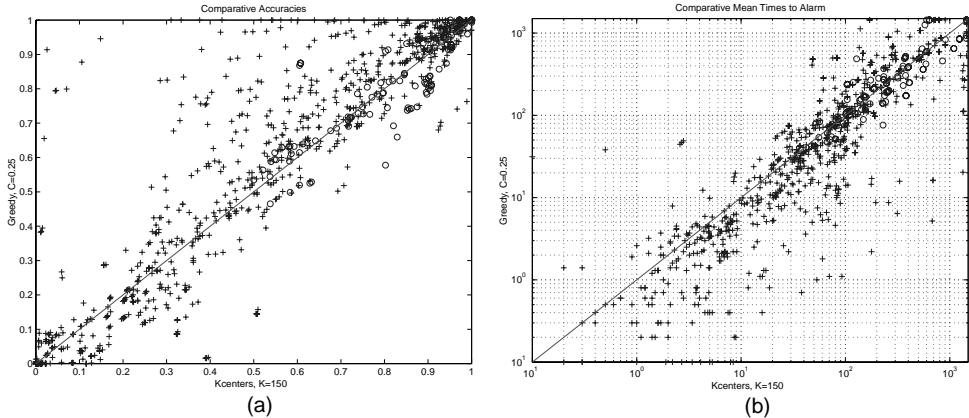


Fig. 12. Comparisons of accuracy (a), and time-to-alarm (b) for greedy clustering method with $C = 0.25$ (vertical axis) versus K -centers clustering with $K = 150$ (horizontal axis). The diagonal is the equal performance line. Circles denote true accept rate/time to false alarm and pluses denote true detect rate/time to true detection.

The interpretation of these plots is not as clear-cut as those of the plots in Section 4.1. Because both methods employ a reduced profile, neither has uniformly superior performance in the way that the base system is uniformly superior to the instance-selection methods. Nonetheless, some trends are evident. The greedy clustering method has generally stronger true detection accuracy (at a mean accuracy value difference of 4.5% over K -centers), but slightly poorer true acceptance accuracy (0.5% lower). The important distinction is revealed in the time-to-alarm results. In this measure, the greedy system has superior results in both time to false alarm (an average of 5.3 longer) and time to true detection (an average of 16.3 shorter). Thus, in terms of practical impact, although the greedy clustering system has a marginally lower true acceptance rate, it is generating false alarms more rarely.

In part, the disparity between the two systems can be attributed to the final profile size. The settings $K = 150$ and $C = 0.25$ were chosen for display here because the profiles they produce are most nearly comparable in size. After reduction by clustering, the K -centers profiles contain only 150 instances, while the greedy profiles contain approximately 300 instances, of which about 180 are cluster centers and the other 120 are original instances not assigned to any cluster. Reaching complete parity in profile sizes is difficult; increasing K beyond 150 could, in theory, produce profiles of equivalent size to those located by the greedy clustering algorithm. But it was found that the K -centers algorithm did not converge to an acceptable solution in the allowed time frame (10,000 search iterations) for $K \gg 150$. The search also did not converge for $K = 150$, but the solutions located were generally of much higher performance than those for larger values of K . This occurs because the K -centers algorithm attempts to optimize all K clusters *simultaneously*, and the combinatorics of assigning

```

**EOF** **SOF** elm findall finger <1> elm cd <1> ls
cd <1> cat <1> finger <1> cd cd <1> ls
<1> elm nn findall finger <1> elm cd <1> ls
**EOF** **SOF** elm findall finger <1> finger <1> chfn finger
**EOF** **SOF** elm findall finger <1> exit **EOF** **SOF** **EOF**
fg cd elm findall finger <1> talk <2> findall **EOF**

```

Fig. 13. An example cluster produced by the greedy clustering algorithm from User 7's data. Each line is a single subsequence of commands and flags. The first line shown is the center point of the cluster. The symbols ****SOF**** and ****EOF**** denote the start and end of shell sessions, respectively.

points to K clusters makes the search difficult for large K . Decreasing K improves the convergence of the algorithm, but yields a smaller final profile and poorer performance. For the data sets tested here, $K = 150$ is found to be most effective.

Decreasing C (the parameter controlling the termination of the greedy search for new clusters) allows the greedy clustering algorithm to search further and locate more clusters, thereby decreasing total profile size. But as $C = 0.25$ already yields about 180 clusters, it is not possible to reduce the profile to 150 instances solely by decreasing C .

Increasing C , on the other hand, causes premature halting of the clustering process and yields fewer clusters (e.g., typically fewer than 5 clusters at $C = 0.75$) and a large number of retained original points (more than 900, or greater than 90% of the original profile for $C = 0.75$). Doing so does improve accuracy with respect to the K -centers system, but this is to be expected given the much larger profile of the greedy system in this case.

Thus, the best performance of the K -centers system is beaten in this domain by a rather conservative setting of the greedy clustering system. We attribute this to the dimensionality of the space (attempting to optimize $K = 150$ clusters simultaneously is combinatorially difficult). While the K -centers algorithm is guaranteed to converge to a locally optimal cluster assignment, there is no guarantee of time bounds on the convergence. In this domain, we have found that the performance of the K -centers based system is bounded largely by the convergence rate. The greedy clustering method does not suffer from the same scaling difficulties because it is guaranteed to terminate.

We note, in passing, that the clusters constructed by the greedy clustering algorithm make intuitive sense, in terms of the actions being performed by the underlying sequences. For example, we have identified clusters that correspond to “programming,” “paper writing,” “reading email,” and “navigating directories.” An example of such an “intuitive” cluster is shown in Figure 13.

4.2.2 Comparison of Clustering to Base System. Because the greedy clustering algorithm is shown to have superior performance to the

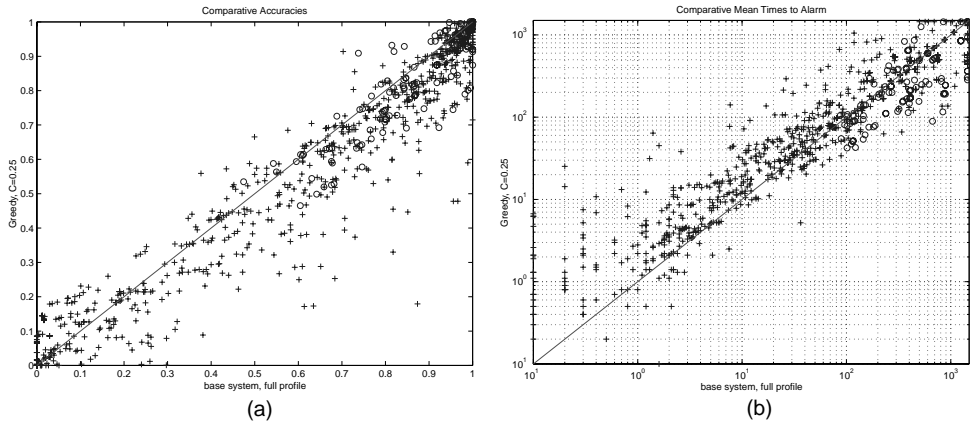


Fig. 14. Comparisons of accuracy (a) and time-to-alarm (b) for greedy clustering method with $C = 0.25$ (vertical axis) versus the base system with the full user profile (horizontal axis). The diagonal is the equal performance line. Circles denote true accept rate/time to false alarm and pluses denote true detect rate/time to true detection.

K -centers algorithm for this data, we present only its results here. The relative performance of the greedy clustering and base systems is displayed in Figure 14. In this figure, the test accuracies of the greedy clustering system with $C = 0.25$ are plotted vertically against the accuracies for the base system (full user profile) on the horizontal. The diagonal line is the iso-performance surface. In the region to the right of this line the base system has superior performance, while in the region above the diagonal the greedy clustering technique has superior performance. The circle symbols represent true acceptance rates, while the plus symbols represent true detection rates. The corresponding times-to-alarm are given in (b).

The situation in this comparison is analogous to those in Section 4.1; a reduced user profile (approximately 300 instances for the greedy system with $C = 0.25$) compared to the full user profile (1,000 instances) of the base system. Not surprisingly, the base system has better overall performance in both accuracy and time-to-alarm (note that the time to false alarms is longer for the base system, while the time to true detection is shorter for the base system). What *is* interesting is that the greedy clustering technique does not display the same behavior that the instance-selection techniques do. Recall that the instance-selection systems produced true accept and true detect rates that were well separated by the iso-performance diagonal line. That result is characteristic of a “nobody is the profiled user” problem — enough data was removed from the profile that the systems were unable to characterize the profiled user well, so they marked too many situations as anomalous. The greedy clustering system does not seem to suffer from this fault. It has lower accuracy on both true detect and true accept, indicating that it has lost ground against the full profile, but has done so gracefully. Rather than learning an incorrect generalization (“nobody is the profiled user”), the greedy clustering method

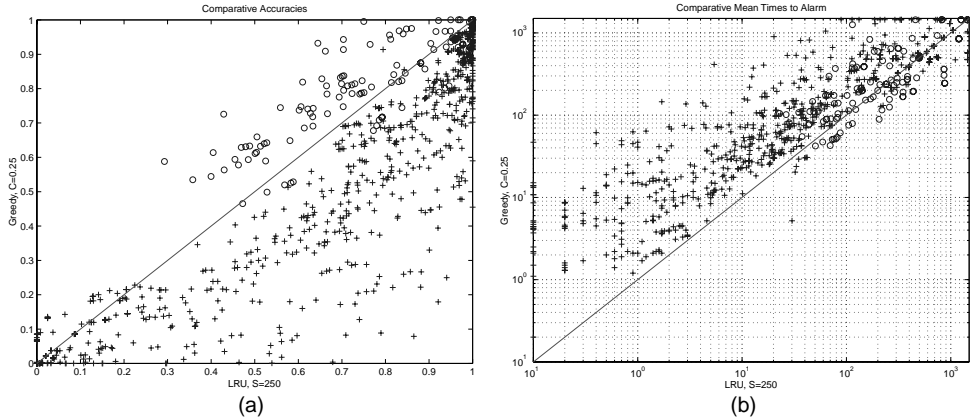


Fig. 15. Comparisons of accuracy (a) and time-to-alarm (b) for greedy clustering method with $C = 0.25$ (vertical axis) versus LRU instance selection with $S = 250$ (horizontal axis). The diagonal is the equal performance line. Circles denote true accept rate/time to false alarm and pluses denote true detect rate/time to true detection.

has preserved the correct characterization of the user’s behavior and has degraded gracefully under the pressure of reduced profile size. The K -centers clustering method displays a similar form of degradation.

4.2.3 Comparison of Clustering to Instance Selection. Since the LRU and LFU instance-selection techniques have comparable performance and are the most effective of the instance-selection techniques we examined, we compare only LRU to the greedy clustering method. The most nearly comparable settings, in terms of profile size, for the two systems are $C = 0.25$ (yielding approximately 300 instances in the final user profile) and $S = 250$ (or 250 instances retained in the profile after LRU instance selection). Comparisons of the two system are displayed in Figure 15. In this figure the test accuracies of the greedy clustering system with $C = 0.25$ are plotted vertically against the accuracies for the LRU instance-selection system with $S = 250$ on the horizontal. The diagonal line is the iso-performance surface. In the region to the right of this line the LRU system has superior performance, while in the region above the diagonal the greedy clustering technique has superior performance. The “o” symbols represent true acceptance rates, while the “+” symbols represent true detection rates. The corresponding times-to-alarm are given in (b).

The difference between the systems is substantial. The LRU system generally has much higher true detection rates, but much lower true acceptance rates. This difference is reflected in the TTA plot, where we find that LRU has much shorter time to detection of an anomaly (by an average of 61.3 tokens) at the cost of a shorter time to generation of false alarms (38.9 tokens on average). This type of result is reminiscent of those in Section 4.1 and, indeed, the cause is similar. The profiles produced by the LRU algorithm often do not preserve important behavior exemplars that

are needed to provide high similarity matches to the test data. The result is improperly set decision thresholds that reject a large proportion of test data, yielding high detection rates for opponents, but low acceptance rates for the profiled user. The characteristic separation between accept rate and detect rate points is not as dramatic as those in Section 4.1.1 because the comparison is not against the full profile, but against a profile selected by the greedy clustering algorithm. The existence of a clear separation is an additional indication that the greedy clustering system is performing in the same fashion (albeit at lower overall accuracies) as the base system. Overall the greedy algorithm appears in most cases to produce superior profiles to those generated by the LRU selection process.

4.3 Summary of Data-Reduction Techniques

Because user behavioral profiles can, potentially, grow without bounds, it is necessary to examine methods for reducing the size of the profile. Profile reduction not only saves space but time, since the running time of the classification operation in an instance-based learning system is related to the number of instances in the instance dictionary. We examined two classes of techniques for reducing the profile size: instance-selection techniques based on examining the interactions of profile instances with *external* parameter selection data and clustering techniques based on examining the relations of instances *within* the profile.

Within the class of instance-selection techniques, we found that the LRU and LFU methods were comparable, and both soundly beat other contenders. FIFO, the other “intelligent” instance-selection technique, performed poorly, even losing to random-element selection. All of the instance-selection techniques, however, display the same class of performance degradation. They all generalize to an incorrect concept of user behavior, in which they tend to classify novel data as anomalous, resulting in an improved ability to recognize impostors at the cost of a decreased ability to accept the profiled user. The major difference between these methods is in the degree to which they succumb to this incorrect generalization.

The clustering methods exhibit a different class of behavior. Rather than learning to incorrectly label too many query instances as anomalous, they degrade in a more uniform manner, losing both true accept and true detect accuracy as the profile is reduced. This is a preferable form of degradation, as it preserves the concept of user behavior rather than replacing it with a “nobody is the profiled user” concept. In the limit, the correct concept is likely to scale better, while the “nobody is the profiled user” concept could reduce to marking all behaviors as anomalous.

The greedy clustering algorithm outperforms the K -centers method in this domain. The difference appears to result from K -centers’s slow convergence for large K . As K (the number of clusters to locate) grows, K -centers takes longer to converge to a solution. By the time K is as large as the number of clusters located by the greedy method (approximately 180),

Table I. Impersonation Incidence for the Greedy Clustered System at an acceptable false alarm level of 0.005. The valid users are listed on the left and impostors across the top. The character “I” represents an instance of a successful impersonation

	USER0	USER1	USER2	USER3	USER4	USER5	USER6	USER7
USER0	I	-	-	-	-	-	-	-
USER1	-	I	-	-	-	-	-	-
USER2	-	-	I	-	-	I	-	-
USER3	-	-	-	I	-	-	-	-
USER4	I	I	I	I	I	I	I	I
USER5	I	-	-	-	-	I	-	-
USER6	-	-	-	-	-	-	I	-
USER7	-	-	-	-	-	-	-	I

K -centers does not locate an acceptable cluster solution in an acceptable time (10,000 search iterations).

5. PRACTICAL IMPLEMENTATION ISSUES

Although it is likely that the system presented here does not possess the level of performance necessary for a fielded anomaly-detection system, it is nonetheless educational to examine the results given above in terms of practical utility and to consider some ways of rectifying weaknesses in the present system.

5.1 Impersonation Rates

One measure of the practical effectiveness of an anomaly-detection system is the chance that an impostor can find a valid user to impersonate successfully. “Successful impersonation” can be interpreted in a number of ways, and we examine two of them. First, we may consider that an impersonation is successful if an impostor appears indistinguishable from a valid user. For example, if the impostor displays times-to-alarm (TTAs) that are at least as long as those of a valid user, the impostor may be considered successful. A summary of this measure appears in Table I.

Table I displays incidences of successful impersonations for the greedy clustering system at a halting threshold of $C = 0.25$ (see Section 4.2) and an acceptable false alarm rate threshold of 0.005 (see Section 2.2.5). Whenever an impostor’s mean time-to-alarm (averaged across all test folds) meets or exceeds that of the valid user, the impersonation is considered successful and is denoted “I.” The valid user is always considered successful in impersonating him or herself. We see that under this measure, successful impersonations are relatively difficult to achieve, the notable exception being impersonations of USER4. Overall, there are only nine impersonations out of 56 possible, or about a one in six chance. USER4 is particularly bad, as observed in Section 3.4, because of the way in which his behaviors change over time. As we discuss in the next section, online learning methods can probably compensate for most or all of the difficulty. If we disregard USER4, there are only two impersonations out of 49

Table II. Impersonation Incidents for fixed time frames under the greedy clustering system with a acceptable false alarm rate of 0.005. The valid user appears on the left and the impostors across the top. “-” characters denote successful detections and “I” characters denote successful impersonations.

	USER0	USER1	USER2	USER3	USER4	USER5	USER6	USER7
10 tokens								
USER0	I	I	I	I	I	I	I	I
USER1	I	I	I	I	I	I	I	I
USER2	I	I	I	I	I	I	I	I
USER3	I	I	I	I	I	I	-	I
USER4	I	I	I	I	I	I	I	I
USER5	I	-	I	I	I	I	-	-
USER6	I	I	I	I	I	I	I	I
USER7	-	I	I	I	I	-	I	I
100 tokens								
USER0	I	I	I	I	-	I	I	I
USER1	-	I	-	I	I	I	-	-
USER2	I	-	I	-	-	I	-	-
USER3	-	-	-	I	-	I	-	-
USER4	I	I	I	I	I	I	I	I
USER5	I	-	-	-	-	I	-	-
USER6	I	I	I	I	-	I	I	I
USER7	-	-	-	-	-	-	-	I
500 tokens								
USER0	I	-	I	-	-	I	-	-
USER1	-	I	-	-	-	-	-	-
USER2	-	-	-	-	-	I	-	-
USER3	-	-	-	-	-	-	-	-
USER4	I	I	I	I	-	I	I	I
USER5	-	-	-	-	-	-	-	-
USER6	I	-	I	-	-	-	I	-
USER7	-	-	-	-	-	-	-	-

possible, or about a one in 25 chance. If we relax the acceptable false alarm threshold to 0.1 (allowing more false alarms in exchange for fewer false accepts), the figures improve to $6/56 = 10.7\%$ with USER4 and no successful impersonations without USER4.

A more conservative way to estimate the impersonation rate, however, is to examine the chance of an impostor going undetected for some fixed period of time. The results of such an analysis appear in Table II. The three subtables display the incidence of successful impersonations at TTA cutoffs of 10, 100, and 500 tokens (recall, however, that all TTAs are measured *after* the smoothing window, Section 2.2.4, was applied so an additional $W = 100$ tokens were observed in addition to the TTA cutoff value). If an alarm has not, on average, been issued within the specified cutoff, the impersonation is considered successful and is denoted with “I.” Note that under this definition, if a false alarm occurs within the cutoff period, then the valid user is considered to have been unsuccessful at impersonating him or herself.

We see from Table II that impersonation is relatively easy at short time frames; at a cutoff of 10 tokens, impersonations are successful 50 out of 56 possible times (or 89%). This is, however, not terribly surprising, since the explicit goal of the techniques employed here is to detect long-term attacks rather than hit-and-run scenarios (Section 1). When we extend the TTA cutoff to 100 tokens (or just twice the minimum possible detection time of the smoothing window length), the situation improves. At this level, there are 30 out of 56 successful impersonations, or 61%, and by 500 tokens it has dropped to 15 incidents, or about 27%. When we relax the acceptable false alarm threshold to 0.1, the corresponding rates are 29% and 11%. As above, the presence of USER4 is an important factor. When we omit USER4 from consideration, the rates are 55% at a cutoff of 100 tokens and an acceptable false alarm rate of 0.005, and only 16% at 500 tokens and a 0.005 acceptable false alarm rate.

5.2 Performance Improvements

Although the system described in this paper has a number of demonstrated weaknesses, there are many possible modifications that would likely improve performance substantially. Here we discuss some approaches that could be taken to improve performance, with the goal of reaching a fieldable system. We are currently investigating a number of these issues.

Online training. As we have seen in the empirical results, the detection system displays dramatically poor detection ability for some profiles (either ability to detect impostors or to correctly identify the valid user, Section 3.4). To some extent this is a result of actual behavioral similarities between the users (as illustrated, e.g., in Figure 6). But a second source of such errors stems from an interaction between the normal behavioral changes of the users and the batch mode training employed in these experiments. USER4, for example, displays a very narrow set of behaviors across the training data, but a different set of behaviors in the parameter-selection data. Thus the profile does not reflect observed behaviors well and the classification thresholds are set very permissively, resulting in poor ability to detect impostors. The batch mode training framework is simply not responsive enough to adapt to such behavioral changes. In other work [Lane and Brodley 1998], we examined an online training extension to this system, in which the user profile and classification thresholds are modified after every observed token. We found that such a framework dramatically improves performance and greatly reduces the impact of behavioral change. We did not employ the online training framework in this work because we wished to examine the effects of the various data reduction strategies in isolation, without the possibility of interaction with another factor; but online training techniques should be employed in a fielded application.

Feature set engineering. Although here we employ only command line data (even suppressing file names), many other data sources are available and could provide valuable information for discrimination. Features such as file names, extensions, or types; login times and time of day; network

activity; working directory; and even GUI events are all potential sources of important user characteristics. Careful analysis and engineering of such data sources is likely to provide greater discrimination between users and reduce the degree of overlap between users illustrated in Figure 6.

Alternative similarity measures. In addition to adding more information sources, it may be possible to take better advantage of the information we already employ. The similarity measure we employ here (see Section 2.2.2) is relatively coarse and captures only one notion of temporal interactions. Using a different measure, or multiple together, may allow finer discrimination within the data framework already available. In other work, [Lane 1999], we investigated the use of hidden Markov models as user profiles and similarity measures. Other possible measures that we are or will investigate include edit distance [Cormen et al. 1992]; stochastic context-free grammars [Charniak 1997]; or frequent episode analysis [Srikant and Agrawal 1996].

Metalearning and hierarchical classifiers. It is likely that no single sensor will possess the breadth of capability to handle the full anomaly-detection problem. But multiple sensors can be combined in a metalearning framework such as voting, boosting [Freund and Schapire 1997] or stacking [Schaffer 1994] to yield an overall more accurate classifier (assuming that the base-level classifiers are at least minimally competent and make statistically independent errors). It is notable that recent intrusion-detection frameworks such as AAFID [Balasubramaniyan et al. 1998] or EMERALD [Porras and Neumann 1997] take advantage of such possibilities by employing hierarchical classification schemes.

Domain knowledge. The system presented here employs relatively little domain knowledge, assuming only that the data is represented as a stream of nominal events. By including *a priori* domain knowledge available from a security expert (for example, “Pay special attention to activity on critical files such as `/etc/passwd` or the Windows NT registry”) performance could be improved at the cost of decreasing generality of the techniques (a trade-off worth making for a fielded system).

Automatic parameter selection. The system presented here possesses a number of parameters that must be set: sequence length l , noise suppression window length W , and greedy cluster halting criteria ϵ and C . The results presented here apply single parameter settings to all users and profiles simultaneously. We found that there is a significant impact of parameter settings on both detection accuracy and time to detection. If the parameters can be properly set on a per-profile basis, then performance can likely be improved. We are currently investigating methods for automatically adapting system parameters to the profiled user.

6. CONCLUSIONS

This work has demonstrated a technique for mapping the temporal data in the anomaly-detection task onto a space in which IBL learning can be formulated. The results demonstrate that this technique can be useful in

such tasks when the underlying data supports sufficient class separation. Furthermore, the system is biased toward detecting anomalous conditions *quickly*, but rarely generating false alarms. We demonstrate two classes of data-reduction methods for limiting profile size. We found that instance-selection techniques, based on relations between instances in the user profile and external parameter-selection data, were susceptible to improper generalizations of the normalcy concept. Clustering-based methods, stipulated on examining relations solely among instances within the profile, did not suffer from this misgeneralization. We showed a new clustering technique based on sequential, greedy selection of clusters. The greedy clustering technique was able to produce a large saving in storage requirements, with an overall small loss in accuracy. K -centers clustering was unable to match the performance of greedy clustering in this domain in either convergence rate or detection accuracy.

In summary, we have presented a machine-learning-oriented approach to anomaly detection, and demonstrated that, in this framework, it is possible to learn to distinguish anomalous behavior patterns from normal ones. We believe that, in general, both the computer security and machine-learning communities can benefit from further interactions. The machine-learning community has studied many pattern-recognition techniques that could be valuable for a variety of security problems, while computer security tasks present a number of challenging issues that can motivate new research directions for the machine-learning community.

ACKNOWLEDGMENTS

We thank Gene Spafford and the members of the Purdue University Center for Education and Research in Information Assurance and Security, the members of the Purdue machine-learning lab, and our reviewers for their helpful comments. We also thank our data donors, and especially USER4 whose data forced us to examine this domain more closely than we might otherwise have done.

REFERENCES

- AHA, D. W., KIBLER, D., AND ALBERT, M. K. 1991. Instance-based learning algorithms. *Mach. Learn.* 6, 1 (Jan. 1991), 37–66.
- ANGLUIN, D. 1987. Learning regular sets from queries and counterexamples. *Inf. Comput.* 75, 2 (Nov. 1, 1987), 87–106.
- ASLAM, J. A. AND RIVEST, R. L. 1990. Inferring graphs from walks. In *Proceedings of the 3rd Annual Workshop on Computational Learning Theory (COLT '90, Rochester, NY, Aug. 6–8)*, M. Fulk and J. Case, Eds. Morgan Kaufmann Publishers Inc., San Francisco, CA, 359–370.
- BALASUBRAMANIYAN, J. S., GARCIA-FERNANDEZ, J. O., ISACOFF, D., SPAFFORD, E., AND ZAMBONI, D. 1998. An architecture for intrusion detection using autonomous agents. Tech. Rep. COAST TR 98/05. Purdue University, West Lafayette, IN.
- BOLLOBÁS, B., DAS, G., GUNOPULOS, D., AND MANNILA, H. 1997. Time-series similarity problems and well-separated geometric sets. In *Proceedings of the 13th Annual Symposium on Computational Geometry (Nice, France, June 4–6, 1997)*, J.-D. Boissonnat, Ed. ACM Press, New York, NY, 454–456.
- CASELLA, G. AND BERGER, R. L. 1990. *Statistical Inference*. Brooks-Cole, CA.

- CHARNIAK, E. 1997. Statistical techniques for natural language parsing. *AI Mag.* 18, 4, 33–43.
- CHENOWETH, T. AND OBRADOVIC, Z. 1996. A multi-component nonlinear prediction system for the S&P 500 index. *Neurocomputing* 10, 3, 275–290.
- CORMEN, T. H., LEISERSON, C. E., AND RIVEST, R. L. 1990. *Introduction to Algorithms*. MIT Press, Cambridge, MA.
- DAS, G., GUNOPULOS, D., AND MANNILA, H. 1997. Finding similar time series. In *Proceedings of the ACM SIGMOD International Workshop on Data Mining and Knowledge Discovery* (Aug.), R. Ng, Ed. ACM Press, New York, NY.
- DENNING, D. E. 1987. An intrusion-detection model. *IEEE Trans. Softw. Eng.* 13, 2, 222–232.
- DOMINGOS, P. 1995. Rule induction and instance-based learning: A unified approach. In *Proceedings of the 14th International Joint Conference on Artificial Intelligence (AAAI-95)*, Montreal, Que., Canada). Morgan Kaufmann, San Mateo, CA, 1226–1232.
- FARMER, D. AND VENEMA, W. 1995. SATAN overview (Security Administrator Tool for Analyzing Networks).
- FREUND, Y. AND SCHAPIRE, R. E. 1997. A decision-theoretic generalization of on-line learning and an application to boosting. *J. Comput. Syst. Sci.* 55, 1, 119–139.
- FUKUNAGA, K. 1990. *Introduction to Statistical Pattern Recognition*. 2ND Academic Press Prof., Inc., San Diego, CA.
- GORDON, S. 1996. Current computer virus threats, countermeasures, and strategic solutions, White paper. McAfee Associates.
- HEBERLEIN, L. T., DIAS, G. V., LEVITT, K. N., MUKHERJEE, B., WOOD, J., AND WOLBER, D. 1990. A network security monitor. In *Proceedings of the IEEE Symposium on Research in Security and Privacy* (Oakland, CA). IEEE Computer Society Press, Los Alamitos, CA, 296–30304.
- IBA, G. A. 1979. Learning disjunctive concepts from examples. Master's Thesis. MIT Press, Cambridge, MA.
- KUMAR, S. 1995. Classification and detection of computer intrusions. Ph.D. Dissertation. Purdue University, West Lafayette, IN.
- KUMAR, S. AND SPAFFORD, E. 1994. An application of pattern matching in intrusion detection. Tech. Rep. CSD-TR-94-013. Purdue University, West Lafayette, IN.
- LANE, T. 1999. Hidden Markov models for human/computer interface modeling. In *Proceedings of the IJCAI-99 Workshop on Learning About Users*. 35–44.
- LANE, T. AND BRODLEY, C. E. 1997a. An application of machine learning to anomaly detection. In *Proceedings of the 20th National Conference on National Information Systems Security. Vol.1* (Baltimore, MD). National Institute of Standards and Technology, Gaithersburg, MD, 366–380.
- LANE, T. AND BRODLEY, C. E. 1997b. Sequence matching and learning in anomaly detection for computer security. In *Proceedings of the AAAI-97 Workshop on AI Approaches to Fraud Detection and Risk Management (AAAI-97)*. 43–49.
- LANE, T. AND BRODLEY, C. E. 1998. Approaches to online learning and concept drift for user identification in computer security. In *Proceedings of the Fourth International Conference on Knowledge Discovery and Data Mining*. 259–263.
- LUNT, T. F. AND JAGANNATHAN, R. 1988. A prototype real-time intrusion-detection expert system. In *Proceedings of the IEEE Symposium on Research in Security and Privacy*. IEEE Computer Society Press, Los Alamitos, CA, 59–66.
- MOON, T. K. 1996. The expectation-maximization algorithm. *IEEE Trans. Signal Process.* 44, 1, 47–59.
- NORTON, S. W. 1994. Learning to recognize promoter sequences in *E. coli* by modeling uncertainty in the training data. In *Proceedings of the 12th National Conference on Artificial Intelligence (vol. 1)* (AAAI '94, Seattle, WA, July 31-Aug. 4, 1994), B. Hayes-Roth and R. E. Korf, Eds. Amer. Assn. for Artificial Intelligence, Menlo Park, CA, 657–663.
- OPPENHEIM, A. V. AND SCHAFER, R. W. 1989. *Discrete-Time Signal Processing*. Prentice-Hall, Inc., Upper Saddle River, NJ.
- PORRAS, P. AND NEUMANN, P. 1997. EMERALD: Event monitoring enabling responses to anomalous live disturbances. In *Proceedings of the 20th National Conference on National*

- Information Systems Security. Vol.1* (Baltimore, MD). National Institute of Standards and Technology, Gaithersburg, MD, 353–365.
- PROVOST, F. AND FAWCETT, T. 1998. Robust classification systems for imprecise environments. In *Proceedings of the 15th National Conference on Artificial Intelligence: Innovative Applications of Artificial Intelligence (AAAI '98/IAAI '98, July 26–30, 1998, Madison, WI)*, J. Mostow, C. Rich, and B. Buchanan, Eds. Amer. Assn. for Artificial Intelligence, Menlo Park, CA, 706–713.
- QUINLAN, J. R. 1993. *C4.5: Programs for Machine Learning*. Morgan Kaufmann Publishers Inc., San Francisco, CA.
- RABINER, L. AND JUANG, B.-H. 1993. *Fundamentals of Speech Recognition*. Prentice-Hall signal processing series. Prentice-Hall, Inc., Upper Saddle River, NJ.
- RABINER, L. R. 1989. A tutorial on hidden Markov models and selected applications in speech recognition. *Proc. IEEE* 77, 2 (Feb.).
- RIPLEY, B. D. 1996. *Pattern Recognition and Neural Networks*. Cambridge University Press, New York, NY.
- RIVEST, R. L. AND SCHAPIRE, R. E. 1989. Inference of finite automata using homing sequences. In *Proceedings of the 21st Annual ACM Symposium on Theoretical Computing*. 411–420.
- SALZBERG, S. 1991. A nearest hyperrectangle learning method. *Mach. Learn.* 6, 3 (May 1991), 251–276.
- SALZBERG, S. 1995. Locating protein coding regions in human DNA using a decision tree algorithm. *J. Comput. Biology* 2, 3, 473–485.
- SCHAFFER, C. 1994. Cross-validation, stacking, and bi-level methods for stacking: Meta-methods for classification learning. In *Selecting Models from Data: Artificial Intelligence and Statistics*, P. Cheeseman and W. Oldford, Eds. Springer-Verlag, Vienna, Austria.
- SMAHA, S. E. 1988. Haystack: An intrusion detection system. In *Proceedings of the Fourth Conference on Aerospace Computer Security Applications*. 37–44.
- SRIKANT, R. AND AGRAWAL, R. 1996. Mining sequential patterns: Generalizations and performance improvements. In *Proceedings of the Fifth International Conference on Extending Database Technology* (Avignon, France).
- STANIFORD-CHEN, S., CHEUNG, S., CRAWFORD, R., DILGER, M., FRANK, J., HOAGLAND, J., LEVITT, K., WEE, C., YIP, R., AND ZERKLE, D. 1996. GrIDS--a graph-based intrusion detection system for large networks. In *Proceedings of the 19th Conference on National Information Systems Security* (Oct.). National Institute of Standards and Technology, Gaithersburg, MD.
- WILSON, D. R. AND MARTINEZ, T. R. 1999. Reduction techniques for exemplar-based learning algorithms. *Mach. Learn.* 34.

Received: February 1999; accepted: July 1999