

Cost-based Modeling and Evaluation
for Data Mining
With Application to Fraud and Intrusion
Detection:
Results from the JAM Project*

Salvatore J. Stolfo, Wei Fan, Wenke Lee, Andreas Prodromidis
Computer Science Department
Columbia University
500 West 120th Street, New York, NY 10027
{sal,wfan,wenke,andreas}@cs.columbia.edu

Philip K. Chan
Computer Science Department
Florida Institute of Technology
150 W. University Blvd., Melbourne, FL 32901
pkc@cs.fit.edu

August 27, 1999

Abstract

In this paper we describe the results achieved using the JAM distributed data mining system for the real world problem of fraud detection in financial information systems. For this domain we provide clear evidence that state-of-the-art commercial

*This research is supported in part by grants from DARPA (F30602-96-1-0311) and NSF (IRI-96-32225 and CDA-96-25374).

fraud detection systems can be substantially improved in stopping losses due to fraud by combining multiple models of fraudulent transaction shared among banks. While studying this problem we have discovered the same important principle reported by Provost and Fawcett [?] for the related domain of cellular phone fraud, namely, that the traditional statistical metrics used to train and evaluate the performance of learning systems, (i.e. statistical accuracy or ROC analysis) are misleading and perhaps inappropriate for some important applications. *Cost-based* metrics are more relevant in certain domains, and defining such metrics poses significant and interesting research questions both in evaluating systems and alternative models, and in formalizing the problems to which one may wish to apply data mining technologies. This view is shared by Kleinberg et al[?] in which a microeconomic view of data mining is espoused, especially with respect to mining association rules and market segmentations.

This paper also demonstrates how the techniques developed for fraud detection can be generalized and applied to the important area of Intrusion Detection in networked information systems. We report the outcome of recent evaluations of our system applied to *tcpdump* network intrusion data specifically with respect to statistical accuracy. This work involved building additional components of JAM that we have come to call, MADAM (Mining Audit Data for Automated Models for Intrusion Detection). However, taking the next step to define cost-based models for intrusion detection poses interesting new research questions. We describe our initial ideas about how to evaluate intrusion detection systems using cost models learned during our work on fraud detection.

1 Introduction

In this paper we discuss the results achieved over the past several years on the JAM Project¹. JAM was initiated as a DARPA and NSF sponsored research project studying algorithms, techniques and systems for distributed data mining. JAM

¹Browse the URL <http://www.cs.columbia.edu/~sal/JAM/PROJECT> for the complete set of progress reports, publications and JAM software available for download. This paper also serves as a guided tour to the project web site.

stands for Java Agents for Meta-Learning. The initial conceptions that we proposed involves the use of agent-based technologies to dispatch machine learning and data analysis programs to remote database sites, resulting in a distributed collection of derived models [?, ?, ?, ?]. Subsequently, the same agent-based technology provides the means for derived “base” models to migrate in a network information system and be collected at any of the participating sites where they may be combined. A single “aggregate” model would then be computed by combining the remotely computed base models, by a technique we have called “meta-learning”, that outputs a “meta-classifier”. Each parallel site may compute its own meta-classifier, trained on locally stored data. These meta-classifiers likewise can migrate around the network to be combined with each other in a hierarchical fashion. The architecture of JAM has been reported extensively in prior papers [?, ?].

A variety of algorithms and techniques have been studied by many researchers that aim to combine multiple computed models. (See [?, ?, ?, ?], for example.) The JAM system provides generic features to easily implement and “snap-in” any of these combining techniques (as well as a large collection of base learning algorithms) and it has been broadly available for use, and continues to be developed. These concepts were studied in order to solve several significant generic problems dealing with scalability of data mining over large and inherently distributed data sources. However interesting these concepts may be (and, by some researchers, these are now regarded as “standard techniques”), JAM’s utility has been demonstrated in the context of real-world problems that in and of themselves are very important. Fraud and intrusion detection are key elements of a new national challenge to protect our nation’s critical infrastructures.

2 The Fraud Detection Problem

We consider the problem of detecting fraudulent transactions after they have been subject to *fraud prevention* methods and processes. There is a vast literature on various security methods to protect transactors from unauthorized use or disclosure of their private information and valuable assets. Even so, thieves find a way through many clever means of circumventing a myriad of prevention techniques. Fraud detection is the last wall of defense.

Many transaction media, such as ATM bank cards or debit cards, require the use of pins, passwords, and in some cases “biometrics” to authenticate the legitimate owner. Credit cards pose interesting special problems since generally no pin

is required for their use; only the name, expiration date and account number is required². A common means of illicitly transacting with credit cards is by stealing someone's identity, and in some cases, creating a new bogus identity. Fraudulent electronic transactions by credit card are a significant problem especially in the context of e-commerce on the internet³. Credit cards need not be physically presented to transact (i.e., they may be used to purchase by phone), and over the internet they can be more rapidly used to fraudulently transact netting better returns for thieves, and heavier losses for banks and their customers.

Financial institutions today typically develop custom fraud detection systems targeted to their own asset bases. The key concept in fraud detection is that fraud may be detected by noticing significant deviation from the "normal behavior" of a customer's account. The behavior of an account can thus be used to protect that account. Notice, it is considerably easier to steal someone's identity information than it is to steal their behavior⁴.

Recently though, banks have come to realize that a unified, global approach is required to detect fraud, involving the periodic sharing with each other of information about attacks. We have proposed a new wall of protection consisting of pattern-directed inference systems using models of fraudulent transaction behaviors to detect attacks. This approach requires analysis of large and inherently distributed databases of information about transaction behaviors to produce models of "probably fraudulent" transactions. We have used JAM to compute these models.

JAM is used to compute local fraud detection agents that learn how to detect fraud and provide intrusion detection services within a single information system; JAM provides an integrated meta-learning system that combines the collective knowledge acquired by individual local agents from among participating bank sites. Once derived local classifier agents or models are produced at some data-site(s), two or more such agents may be composed into a new classifier agent by JAM's meta-learning agents. JAM allows financial institutions to share their models of fraudulent transactions by exchanging classifier agents in a secured agent infrastructure. But they will not need to disclose their proprietary data. In this way their competitive and legal restrictions can be met, but they can still share

²which is private information readily available from trash cans.

³Recent studies issued by the US Department of Commerce indicate that of the \$13 Billion in credit card sales on the web in 1998, the percentage of fraud is a large and growing problem, prompting some to dub the WWW as the Wild Wild West.

⁴This may change as vast databases of "click data" revealing personal behavior information continues to be acquired over the internet typically without the user's knowledge.

information. The meta-classifiers then act as sentries forewarning of possibly fraudulent transactions and threats by inspecting, classifying and labeling each incoming transaction.

2.1 The Fraud Learning Task

The learning task here is quite straightforward. Given a set of “labeled transactions”, $T = \{t | t = \langle f_1, \dots, f_n \rangle\}$, compute a model or classifier, C , by some learning algorithm L , that predicts from the features $\langle f_1, \dots, f_{n-1} \rangle$ of some unknown the target class label f_n , “fraud”. Hence, $C = L(T)$, where L is a learning algorithm. Each element $t \in T$ is a vector of features, where we denote f_1 as the “transaction amount” (*tranamt*), and f_n as the target class label, denoted $fraud(t) = 0$ (legitimate transaction) or 1 (a fraudulent transaction). Given a “new unseen” transaction, x , with unknown class label, we compute $f_n(x) = C(x)$. C serves as our fraud detector.

Much of our work on the JAM project has focussed on developing, implementing and evaluating a range of learning strategies and combining techniques. Our work on “meta-learning” strategies has focussed on computing sets of “base classifiers” over various partitions or samplings of the training data, T , and various performance metrics to evaluate base classifiers over test data. (Nearly all of our reported results are based upon k-fold cross validation.)

In meta-learning, we first seek to compute a set of base classifiers, $\{C_i, i = 1, \dots, m\}$, where $C_i = L_j(T_k)$, $\bigcup_k T_k = T$, varying the distributions of training data (T_k) and using a variety of different machine learning algorithms (L_j) in order to determine the “best” strategies for building good fraud detectors. The “best” base classifiers are then combined by a variety of techniques in order to boost performance. One of the simplest combining algorithms proposed in[?], and independently by Wolpert[?], is called “class-combiner” or “stacking”. A separate hold out training dataset, V , is used to generate a meta-level training data to learn a new “meta-classifier”, M . M is computed by learning a classifier from training data composed of the predictions of a set of base classifiers generated over a set of validation data (V) along with the true class label. Hence, $M = L(\langle C_1(v), \dots, C_m(v), f_n(v) \rangle), v \in V$. The resultant meta-classifier works by inputting the predictions for some unknown into its constituent base classifiers, and then generating its own final class prediction from these base classifier predictions. Thus, for unknown x , $f_n(x) = M(C_1(x), \dots, C_m(x))$.

Notice, M is as well a classifier, or fraud detector. In the following sections when we make reference to a classifier C , it may either be a “base classifier”,

or an ensemble “meta-classifier”, learned over some training distribution, unless explicitly stated otherwise.

For notational convenience, we define $J = F[L, M, T, D_T]$ as the set of all classifiers that may be computed where $L = \{L_1, \dots, L_p\}$ is a set of learning algorithms, $M = \{M_1, \dots, M_q\}$ is a set of meta-learning algorithms, T is a set of labeled ground truth data used to train classifiers, and D_T is some training distribution of interest. JAM is a computational environment that is designed to assist data miners in generating desired classifiers.

2.2 Credit Card Datasets

Chase and First Union Banks, members of the FSTC (Financial Services Technology Consortium) provided us with real credit card data for our studies. The two data sets contain credit card transactions labeled as fraudulent or legitimate. Each bank supplied .5 million records spanning one year with 20% fraud and 80% non-fraud distribution for Chase bank and 15% versus 85% for First Union bank. The schemata (or feature sets) of the databases were developed over years of experience and continuous analysis by bank personnel to capture important information for fraud detection. We cannot reveal the details of the schema beyond what is described in [?]. The records have a fixed length of 137 bytes each and about 30 attributes including the binary class label (f_n). Some of the fields are numeric and the rest categorical, i.e. numbers were used to represent a few discrete categories.

The features in this data defined by the banks essentially describe the “usage behavior” of an individual credit card account. The data is rich enough to allow alternative modeling strategies. For example, it is possible to segment the data into classes of accounts based upon “payment history”⁵. It is well known that there are at least two classes of card holders, “transactors” who pay their bills in full each month, and “revolvers” who pay their minimum charge and roll over balances and interest charges to the next month. Models can thus be built for each of these market segments to determine finer distinctions between account transaction behaviors. In this work, we do no such segmentation and partitioning, but rather we compute models of “fraudulent transaction” for the entire asset base.

⁵The reader is encouraged to do a personal study of their own checking account behavior. You will likely find the same type of payments from month to month in your own checking account. Variations or distinct payments in any particular month are likely to reoccur in the same month in the prior year. Such repetitive behaviors can be regarded as a normal profile for the checking account. Significant variations from this normal profile possibly indicate fraud.

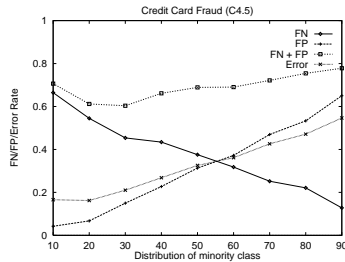


Figure 1: Percentage of minority class in training vs. \mathcal{FN} , \mathcal{FP} , $\mathcal{FN} + \mathcal{FP}$ & error rate

We believe the results achieved in this study can be improved had we segmented the data⁶.

Many experiments were conducted using JAM to evaluate the performance of different learning algorithms, and different meta-level training sets using this data. Prior publications report on these experiments and indicate that the meta-classifiers consistently outperform the best base classifiers. Unfortunately, the studies that considered different training distributions, different learning algorithms and a variety of combining techniques demonstrated that the best strategies for this target domain are not immediately discernible, but rather requires extensive experimentation to find the best models, and the best meta-classifiers [?].

2.3 Skewed Distributions and Real-world constraints

In the credit card domain, the number of legitimate transactions is much larger than the fraudulent ones. Hence, the class (legitimate or fraudulent) distributed is highly skewed. Other domains like cellular phone fraud and network intrusion have similar characteristics. Using the “natural” (highly skewed) class distribution for training, some learning algorithms might treat the minority class as noise or simply produce a classifier that always predicts the majority class. That is, the generated classifier has a very low error rate (high accuracy), which is the most common measure for evaluating performance. For a bimodal problem, an alternative performance criterion is to calculate the false positive and false negative rates.

We conducted some experiments to study how the underlying training distribution affects the error rate, the false-negative rate, the false-positive rate, and the

⁶However, such a strategy would delve dangerously close to industry trade secrets.

sum of FN and FP rates. (Later, we report results in terms of the $TP - FP$ spread, a related measure.) The results are plotted in Figure 1. As expected, our empirical results demonstrate that the error rates vary with the class distributions in the training set. In the fraud domain the error rate generally increases with the amount of fraud (minority class) in the training set, but minimizes at around 20% fraud (natural percentage). Also, the false-negative rate decreases when the negative (minority) percentage increases. That is, increasing the number of minority instances in training produces fewer errors on the minority instances. Conversely, the false-positive rate has the opposite trend. An alternative metric to error rate is the sum of \mathcal{FN} and \mathcal{FP} . This sum considers both types of errors. From the figure, the sum minimizes at around 30% (which is not the natural percentage of 20%). In other words, if the sum of \mathcal{FN} and \mathcal{FP} is used as the performance metric, using the natural distribution might not yield the most effective classifier.

Experiments were also performed by varying the distribution of training data over time. For example, we wondered whether classifiers trained on recent data would perform better than classifiers trained on older data. We noticed a significant methodological danger when experimenting with different training strategies.

In some experiments, we noticed that classifiers trained in a particular month generally performed better in predicting fraud in the next consecutive month (i.e., classifiers trained on May data, were best at catching fraud in June). Unfortunately, this observation has no value. One cannot train a classifier in one month and have it available for use in the next month simply because the data is unavailable. Credit cards have a two month business cycle before the status of a transaction can be determined. There is first a 30-day billing period, followed by a 30-day payment and investigation period before a transaction’s status as either legitimate or fraudulent can be definitely determined.

Hence, data miners must be cautious in understanding the real world context and constraints in developing experiments that are meaningful and appropriate to the target domain.

2.4 Performance Criteria

The “performance criterion” has also been the subject of our research. Here, the “best” detectors should maximize a “predictive performance” criterion as well as a “computational performance” criterion. That is, we seek the best detectors that are

- good at catching fraudulent transactions,

- are efficient to compute over massive datasets of transactions, and
- whose evaluation maintains real-time throughput rates in some target environment.

We address the “predictive performance” criterion, what we mean by “good”, after describing the cost model for this domain. We next briefly touch on the computational performance issues we have studied.

2.5 Computational Performance and Scaling

With respect to computational performance, we have shown how a distributed or parallel implementation of meta-learning in the JAM system is scalable[?, ?] (by computing base classifiers in parallel over distributed sites with reduced or samples of data) and can produce meta-classifiers that can be “sped up” by pruning base classifiers[?, ?] to maximize transaction throughputs. The techniques developed here are general and broadly applicable to a wide range of data mining applications.

Here we report some of our results on scaling distributed learning systems. Results achieved on pruning are presented later in the context of the “cost performance” of the meta-classifiers and the throughput increases achieved.

Our meta-learning approach is highly scalable because training from multiple data subsets at the base level is computationally independent and can be conducted in parallel. Once all the base classifiers are trained (in parallel), we train the meta-classifier. One can easily observe that this approach, for large amounts of data, will provide substantial improvement in speed for super-linear-time learning algorithms.

We investigated both the theoretical and empirical aspects of scalability in meta-learning [?]. *Scaled speedup* [?, ?] provides a metric for *scalability*. It measures the speedup of a parallel system when the problem size increases linearly with the number of processors. Since one level of meta-learning might not recover the loss of predictive performance due to subsampling at the base level [?], we investigated hierarchical meta-learning, in which meta-learning is performed at multiple levels in a tree fashion [?]. Similarly, hierarchical meta-learning is easily parallelizable because sibling classifiers can be generated independently in parallel. In many cases hierarchical meta-learning maintains the predictive performance with the additional benefit of speed via coarse-grained parallel processing [?]. (Note that this approach does not require the time-consuming task of parallelizing a learning algorithm for a particular parallel architecture.)

2.6 Cost-based Models for Fraud Detection

Most of the machine learning literature concentrates on model accuracy (either training error or generalization error on hold out test data computed as overall accuracy, True Positive/False Positive rates, or ROC analysis). This domain provides a considerably different metric to evaluate performance of learned models; models are evaluated and rated by a “cost model.” Within the context of financial transactions, cost is naturally measured in dollars. However, any unit of *measure of utility* applies here. The credit card domain provides an excellent motivating example domain familiar to most people.

Due to the different dollar amounts of each credit card transaction and other factors, the cost of failing to detect a fraud varies with each transaction. Hence, the cost model for this domain is based on the sum and average of loss caused by fraud. We define for a set of transactions S , a fixed *overhead* amount, and a fraud detector (or classifier) C :

$$CumulativeCost(S, C, overhead) = \sum_{t \in S}^n Cost(C(t, overhead)) \quad (1)$$

and

$$AverageCost(S, C, overhead) = \frac{CumulativeCost(S, C, overhead)}{n} \quad (2)$$

where $Cost(t, overhead)$ is the cost associated with transaction t and n is the total number of transactions in a test set S . The cost of a transaction is not simply its “transaction amount”, but is also a function of an *overhead* amount.

After consulting with a bank representative, we jointly settled on a simplified cost model that closely reflects reality. Since it takes time and personnel to investigate a potential fraudulent transaction, a fixed *overhead* value is incurred for each investigation. That is, if the amount of a transaction is smaller than the overhead, it is not worthwhile to investigate the transaction even if it is suspicious. For example, if it takes ten dollars to investigate a potential loss of one dollar, it is more economical not to investigate. Therefore, assuming a fixed *overhead*, we devised the following cost model for each transaction, t and classifier C . (Recall, $f_1(t) = tranamt(t)$.)

Outcome	Cost(t,overhead)	
Miss (False Negative, <i>FN</i>)	$tranamt(t)$	
False Alarm (False Positive, <i>FP</i>)	$overhead$	if $tranamt(t) > overhead$
	0	if $tranamt(t) \leq overhead$
Hit (True Positive, <i>TP</i>)	$overhead$	if $tranamt(t) > overhead$
	$tranamt(t)$	if $tranamt(t) \leq overhead$
Normal (True Negative, <i>TN</i>)	0	

The *overhead* threshold, for good reason, is a closely guarded secret for important reasons discussed later in section 5.5, and may vary over time. The range of values used in our studies is probably reasonable as appropriate bounds for the data set provided by the banks. All the empirical studies we conducted are evaluated using this cost model.

It is important to note that this *overhead* is not a “score threshold” for classifiers that may output continuous values (i.e., density estimators). Rather, the threshold is a simple “decision boundary”. Transactions whose amounts are under this threshold are immediately authorized (subject to account credit availability of course). Decisions are therefore made by detectors only for transactions above this threshold⁷.

The target application is described as detecting frauds to minimize cumulative cost, or maximize cost savings in this model. In describing our results, we report the *maximum savings* (or *stop loss* in bank parlance), as the total dollar amount saved from detection under this cost model. The total potential dollar loss for a (test) set of transactions (S) is defined as the total dollar amount of all fraudulent transactions:

$$TotalPotentialLoss(S) = \sum_{t \in S_{\text{fraud}(t)=\text{true}}} tranamt(t) \quad (3)$$

A complete comparative evaluation between purely statistical error rates versus cost-model savings can be found in the cited papers. The important lesson here is that the data mining problem is actually a straightforward cost optimization problem, namely to capture or recover the *TotalPotentialLoss* due to fraud. That

⁷When training classifiers, one may think that simply ignoring all transactions in the training data under the threshold will produce better detectors. This may not be true. The fraudulent transaction behavior we seek to learn may not vary with the transaction amount (thieves will do the same things whether stealing ten dollars, or a hundred dollars), and the learning of fraudulent behavior may be more accurate when studying low cost transactions. Even so, some experiments were performed varying the underlying distributions, and we describe these later.

is, given a test set of transactions, S , a fraud model $C_{overhead}$, and the *overhead*, the *TotalCostSavings* is defined as

$$TotalCostSavings(S, C_{overhead}, overhead) = TotalPotentialLoss(S) - CumulativeCost(S, C_{overhead}, overhead) \quad (4)$$

We can now state concretely what the learning task is for fraud detection. We seek to compute the

$$\underset{C_{overhead} \in J}{argmax} \{TotalCostSavings(S, C, overhead)\} \quad (5)$$

where individual classifiers, $C_{overhead}$, can be models computed (by JAM) over training data under the cost model with a fixed *overhead*, T , by many different learning and meta-learning strategies. (Note, the worst possible outcome is that the detector is so bad we actually lose money; the total cost savings may be negative.) We may of course add additional constraints on $C_{overhead}$ so that it, for example, meets real-time, or memory constraints, or is computed over inherently distributed data (partitions of T) as rapidly as possible. Explorations of these issues have been much of the focus of the JAM project.

It is interesting to note here another reason why pure accuracy measures are inappropriate in some domains. Observe in this and related domains (like cellular phone fraud) the distribution of frauds may constitute a very tiny percentage of the total number of transactions. If, for example, the percentage of frauds is 1% or less, than the null detector will be 99% accurate or better! Naturally, we must depend upon at least the TP/FP rates of a detector to measure alternatives. But, more to the point, training detectors in such domains begs an important question: what is the appropriate distribution of data used in training to produce “good detectors”? In our prior work many experiments were conducted to evaluate classifier performance over different samplings and distributions. For example, in one set of tests, 50:50 distributions produced the best detectors⁸. Of particular note is that the best training distributions varied according to the *overhead* and thus the cost-model. This is why we subscript the classifier $C_{overhead}$ to reveal that alternative classifiers can be computed simply by varying this quantity. Here again, determining the best training distribution under realistic environments, with possibly highly skewed distributions is a matter of considerable empirical investigation.

⁸In this context, data is plentiful, so we could afford to construct many alternative distributions without much fear of generating “knowledge poor” training samples.

3 Results using JAM

A large number of experiments were performed to determine whether various alternative training strategies would produce improved or superior models. The strategies included different temporal distributions (different months of training data), different partitions or samples (random, or different distributions of fraud versus non-fraud) and different cost distributions (training over large frauds versus lower cost frauds). The results we report here are only one broad view of what we achieved using the best strategies determined empirically as compared to the best possible outcome for a commercial off the shelf system (COTS).

CHASE provided us with data that had embedded within each transaction record a field recording a score (in the range of 0-999) generated by a commercial off-the-shelf (COTS) fraud detection system for that transaction. From this information we were able to easily compute the best possible detection *cost savings* of this system as well as its statistical accuracy. We do not know what the actual score thresholds may be in practice. However, for the same data provided by CHASE we can compute what the optimal setting for the COTS should be for optimal performance.⁹ Using this, we were able to compare JAM’s models to see if we were indeed doing better.

3.1 Baseline Performance Under the Cost Model

Tables 1 and 2 summarize our results for the Chase and First Union banks respectively. Both tables display the accuracy, the $TP - FP$ spread and savings for each of the fraud predictors examined. Recall, overall accuracy is simply the percentage of correct predictions of a classifier on a test set of “ground truth”. TP means the rate of predicting “true positives” (the ratio of correctly predicted frauds over all of the true frauds), FP means the rate of predicting “false positives” (the ratio of incorrectly predicted frauds over those test examples that were not frauds, otherwise known as the “false alarm rate”.) We use “ $TP - FP$ spread” to indicate how well the system finds true frauds versus false alarms. A “1.00 $TP - FP$ spread” is optimal performance.¹⁰ The maximum loss potential of these test sets is approximately \$1,470,000 for the Chase data and \$1,085,000 for the First Union

⁹We computed the performance of the COTS for a whole range of score thresholds: 5, 10, 20, 30, 40, 50, 100, 150, 200, 250, ... , 850, 800, 950.

¹⁰These are standard terms from the statistical “confusion” matrix that also includes FN and TN rates as well.

Table 1: Performance results for the Chase credit card data set.

Type of Classification Model	Size	Accuracy	TP - FP	Savings
COTS scoring system from Chase	-	85.7%	0.523	\$ 682K
Best base classifier over one subset	1	88.5%	0.551	\$ 812K
Best base classifier over entire set	1	88.8%	0.568	\$ 840K
Meta-classifier, Chase	50	89.6%	0.621	\$ 818K

Table 2: Performance results for the First Union credit card data set.

Type of Classification Model	Size	Accuracy	TP - FP	Savings
Best base classifier over one subset	1	95.2%	0.749	\$ 806K
Best base classifier over entire set	1	95.3%	0.787	\$ 828K
Meta-classifier, First Union	50	96.5%	0.831	\$ 944K

data. The column denoted as “size” indicates the number of base-classifiers used in the meta-classifier.

3.2 JAM versus COTS

The first row of Table 1 shows the best possible performance of Chase’s own COTS authorization/detection system on this data set. The next two rows present the performance of the best base classifiers over the entire set and over a single month’s data, while the last rows detail the performance of the unpruned (size of 50) meta-classifiers. Similar data is recorded in Table 2 for the First Union set, with the exception of First Union’s COTS authorization/detection performance (it was not made available to us).

The outcome was clearly in favor of JAM for this dataset. According to these results, the COTS system achieves 85.7% overall accuracy, 0.523 ‘ $TP - FP$ spread’ and saves \$682K when set to its optimal “score threshold”.

A comparison of the results of Tables 1 and 2 indicates that in almost all instances, meta-classifiers outperform all base classifiers, and in some cases by a significant margin. The most notable exception is found in the “savings” column of Chase bank where the meta-classifier exhibits reduced effectiveness compared to that of the best base classifier.

This shortcoming can be attributed to the fact that the learning task is ill-defined. Training classifiers to distinguish fraudulent transactions is not a direct

Table 3: Pruned meta-classifiers on Chase data.

Type of Classification Model	Size	Accuracy	TP - FP	Savings
Meta-classifier, Chase	46	89.7%	0.595	\$ 667K
Meta-classifier, Chase	32	89.2%	0.632	\$ 870K
Meta-classifier, Chase	5	88.9%	0.589	\$ 903K

Table 4: Pruned meta-classifiers on First Union data.

Type of Classification Model	Size	Accuracy	TP - FP	Savings
Meta-classifier, First Union	30	96.6%	0.816	\$ 898K
Meta-classifier, First Union	12	96.5%	0.848	\$ 942K
Meta-classifier, First Union	29	96.5%	0.844	\$ 950K

approach to maximizing savings (or the $TP - FP$ spread). Traditional learning algorithms are not biased towards the cost model and the actual value (in dollars) of the fraud/legitimate label; instead they are designed to minimize statistical misclassification error. Hence, the most accurate classifiers are not necessarily the most cost effective. Similarly, the meta-classifiers are trained to maximize the overall accuracy not by examining the savings in dollars but by relying on the predictions of the base-classifiers. Naturally, the meta-classifiers are trained to trust the wrong base-classifiers for the wrong reasons, i.e. they trust the base classifiers that are most accurate instead of the classifiers that accrue highest savings.

3.3 Boosting Cost and Computation Performance by Pruning

To improve the performance of our meta-classifiers, we pruned [?] the meta-classifier, i.e. we discarded the base classifiers that do not exhibit the desired property. Although this approach addresses the cost-model problem at a late stage, after base classifiers are generated, it has the advantage of fitting to the requirements of this problem (financial institutions import pre-computed classification models, and hence treat them as black boxes) and also reduces the size of the meta-classifier, thus allowing for faster predictions and better use of system resources.

Tables 3 and 4 present a summary result of the best pruned meta-classifiers and their size (number of constituent base classifiers). As can be seen, depending on the evaluation criterion, the method computes different meta-classifiers of differ-

ent sizes. An additional result, evident from these tables, is the invaluable contribution of pruning. In all cases, pruning succeeded in computing meta-classifiers with similar or better fraud detection capabilities, while reducing their size and thus improving their efficiency. A detailed description on the pruning methods and a comparative study between predictive performance and meta-classifier throughput can be found in [?].

3.4 Bridging Classifiers for Knowledge Sharing

The final stage of our experiments on the credit card data involved the exchange of base classifiers between the two banks. To meta-learn over a set of classifier agents, however, we had to overcome additional obstacles in order to share their knowledge of fraud. The two databases had differences in their schema definition of the transactions, and hence learning over these different sites produced *incompatible* classifiers:

1. Chase and First Union defined an attribute with different semantics (i.e. one bank recorded the number of times an event occurs within a specific time period while the second bank recorded the number of times the same event occurs within a different time period).
2. Chase includes two (continuous) attributes not present in the First Union data.

To address these problems we followed the approaches described in [?, ?, ?]. For the first incompatibility, we had the values of the First Union data set mapped via a linear approximation to the semantics of the Chase data. For the second incompatibility, we deployed special *bridging* agents that were trained to compute the missing values of First Union data set. The training involved the construction of regression models [?] of the missing attributes over the Chase data set using only the attributes that were common to both banks. When predicting, the First Union classifier agents simply disregarded the real values provided at the Chase data sites, while the Chase classifier agents relied on both the common attributes and the predictions of the bridging agents to deliver a prediction at the First Union data sites.

Tables 5 and 6 display the accuracy, $TP - FP$ spread and cost savings of each Chase and First Union meta-classifier. These results demonstrate that both Chase and First Union fraud detectors can be exchanged and applied to their respective data sets. The most apparent outcome of these experiments is the superior

Table 5: Combining Chase and First Union classifiers on Chase data.

Composition of Meta-Classifier	Size	Accuracy	TP - FP	Savings
Chase + First Union	110	89.7%	0.621	\$ 800K
Chase + First Union	63	89.7%	0.633	\$ 877K

Table 6: Combining Chase and First Union classifiers on First Union data.

Composition of Meta-Classifier	Size	Accuracy	TP - FP	Savings
First Union + Chase - bridge	110	96.6%	0.843	\$ 942K
First Union + Chase + bridge	110	97.8%	0.901	\$ 963K
First Union + Chase + bridge	56	97.8%	0.899	\$ 962K

performance of the First Union meta-classifiers and the lack of improvement on the performance of the Chase meta-classifiers. This phenomenon can be easily explained from the fact that the attributes missing from the First Union data set were significant in modeling the Chase data set. Hence, the First Union classifiers are not as effective as the Chase classifiers on the Chase data, and the Chase classifiers cannot perform at their best at the First Union sites without the bridging agents. The latter was verified by a separate experiment, similar to the above, with the exception that no bridging agents were used, i.e. Chase classifiers produced predictions without using any information on the missing values.

The bottom line is that our hypothesis was correct: better performance resulted from combining multiple fraud models by distributed data mining over different transaction record sources (including multiple banks) even when bridging the differences among their schema.

3.5 Cost-sensitive Learning: AdaCost

An interesting observation to be made here is that the various “off the shelf” learning algorithms employed in JAM and used in this study are likely inappropriate for this type of application domain. Since much of the focus in machine learning has been on algorithms and techniques to learn models with high statistical accuracy, it is fair to say that the algorithms primarily employ heuristics and underlying strategies based upon statistical- and information-theoretic metrics like entropy, *not* “cost-based” strategies.

Much of our experimental work has been to “bias” the outcome of the learned

classifiers towards improved cost performance by varying training distributions, or pruning poor cost performing classifiers. This approach is somewhat akin to hammering a square peg into a round hole.

An alternative strategy is called *cost sensitive learning*. The essence of the idea is to bias *feature selection* in generating hypotheses during the learning process in favor of those that maximize a cost criterion (for example, the *cost of testing features*, rather than a purely statistical criterion). According to Turney [?] the earliest work here is due to Nunez [?]. Later work by Tan and Schlimmer [?] also incorporates feature costs in the heuristic for searching in a modified decision tree learning algorithm. However, there are costs associated not only with testing features, but also varying costs based upon classifier *misclassification cost* performance. The distinctions are important.

Two alternative features may have the same “test cost” but their predictive outcomes may produce different “misclassification costs.” Hence, we ought to strategically choose “low cost features” that are both cheap to compute and test, and that reduce the misclassification cost of the final model that employs them.

The earliest work on misclassification cost is due to [?, ?]. A more recent paper, [?] by Pazzani et al, discusses several different approaches to minimizing *misclassification cost*. The earliest work on learning strategies that incorporate *both feature costs and misclassification costs* is [?], where empirical evaluations are reported of a proposed hybrid genetic decision tree induction algorithm.

What the cost model for the credit card domain teaches is that there are different costs depending upon the outcome of the predictions of the fraud detectors. This may appear strange but we may want to compute classifiers that are purposely *wrong* in certain cases so that we do not incur their high costs when they predict correctly. Not only are there costs associated with “misclassifications” (False positives/negatives), but also costs are born with *Correct Predictions*, i.e. True Positives also incur costs (*overhead*)! This simple, but perhaps counterintuitive, this observation has not been accounted for in prior work and has been included in our cost models when computing classifiers and evaluating their outcome.

As mentioned, we have performed experiments to generate cost-sensitive classifiers by varying the distribution of training examples according to their costs (*tranamt*). This strategy doesn’t change the underlying algorithm, but rather attempts to bias the outcome of the underlying (statistical-based) algorithm.

This was achieved by two simple methods: replication and biased sampling. In the first case, experiments were performed where training data was “replicated” some number of times based upon the cost of the exemplars. Another strategy

sampled high cost examples and excluded the low cost transactions (those under the *overhead* amount). These “cost-based training distributions” were used in training base models, and meta-classifiers. Unfortunately, the results indicated that the resultant classifiers did not consistently improve their cost performance [?] over varying cost distributions.

Other experiments were performed to directly bias the internal strategy of the learning algorithm. One algorithm we have proposed and studied is a close variant of Singer and Schapire’s [?] AdaBoost algorithm. AdaBoost is an algorithm that starts with a set of “weak hypotheses” of some training set, and iteratively modifies weights associated with these hypotheses based upon the statistical performance of the hypotheses on the *training set*. Elements of the training set are as well weighted, and updated on successive rounds depending upon the statistical performance of the hypotheses over the individual data elements. AdaBoost ultimately, therefore, seeks to generate a classifier with minimum *training error*.

AdaCost [?] is a variant of AdaBoost that modifies its “weight updating rule” by a “cost based factor” (a function of *transamt* and the *overhead*). Here, training elements that are “mis-classified” are re-weighted by a function of the statistical performance of the hypotheses as well as the “cost” of the element. Costlier mis-classifications are “re-weighted” more for training on the next round. All weights are normalized on each round so correct predictions have their weights reduced. However, the new weights of correct predictions are adjusted by the cost model to account for the cost of true positives as well.

It is not possible to change the underlying training distribution according to the credit card cost model because the cost of a transaction is dependent upon the final prediction of the classifier we are attempting to compute, and is not known a priori, i.e., during training. Since the credit card cost model dictates cost even if the classification is correct, adjusting weights of training examples can’t easily reflect that fact. The best we can do here is incorporate the cost for correct predictions on the “current round” during training to produce a different distribution for the “next round” of training.

3.6 AdaCost versus AdaBoost

Experiments here using AdaCost on the credit card data showed consistent improvement in “stopping loss” over what was achieved using the vanilla AdaBoost algorithm. For example, the results plotted in Figures 2 shows the average reduction of 10 months as a percentage cumulative loss (defined as $\frac{\text{cumulative loss}}{\text{maximal loss} - \text{least loss}} * 100\%$) for AdaBoost and AdaCost for all 50 rounds and 4 overheads. We can

Figure 2: Cumulative Loss Ratio of AdaCost and AdaBoost for Chase Credit Card Data Set

Figure 3: Cumulative Loss Ratio and Loss of AdaCost and AdaBoost on Chase Credit Card

clearly see that, except for round 1 with overhead = 90, there is a consistent reduction for all other 398 ($= 50 * 2 * 4 - 2$) runs. The absolute amount of reduction is around 3%. We also observe that the speed of reduction by AdaCost is quicker than that of AdaBoost in all 4 figures. The speed is the quickest in the first few rounds. This means that in practice, we may not need to run AdaCost for many rounds.

Figure 3 plots the ratio of cumulative cost by AdaCost and AdaBoost. The figures are similar to those in Figure ???. We have plotted the results of all 10 pairs of training and test months over all rounds and overheads. Most of the points are below the “Ratio=1” line in the left drawing and above the “y=x” line in the right drawing, both implying that AdaCost has lower cumulative loss in an overwhelming number of cases.

4 Intrusion Detection: Initial results using MADAM

Encouraged by our results in fraud detection¹¹, we shifted our attention to the growing problem of intrusion detection in network based systems. Here the problems are significantly different, although from a certain perspective we seek to perform the same sort of task as in the credit card fraud domain. We seek to build models of “normal” behavior to distinguish between “bad” (intrusive) connections and “good” (normal) connections.

Notice that here again we are considering the specific case of detecting intrusive connections, after intrusion prevention techniques have been applied. Intrusion detection is essentially the last wall of protection from unauthorized users (including perhaps “insiders”).

In the credit card fraud case, the banks over a long period of time, have designed a set of features describing the behavior associated with each credit card

¹¹And under “encouragement” from DARPA

account that seem to work quite well in predicting fraudulent behavior. These features are “easy” to compute and maintain in real time. Fawcett and Provost[?] studied what appears to be the same problem for the cellular phone fraud domain¹², but their work focussed on defining good features to learn detectors.

But it is important to more deeply understand the nature of the learning task for these fraud detection problems. First, the primary “object” whose “behavior” is being modeled and represented by these features is not an individual person *per se*, but rather the “normal credit card account behavior” (in the bank case) and the “normal calling behavior” (in the cell phone case) of a large group of costumers. In the case of Intrusion Detection in network computer systems there are several different categories or types of “attacks” on different “objects” for which one would like to model behavior. For example, in the DARPA/MIT LL datasets described later, there were approximately 48 different known attacks that can be classified into 4 primary attack categories (*Probes*, *Denial of Service*, *User-to-root-access* and *Remote-to-local-access*). Thus, we can consider the learning task to be multi-class, rather than bimodal in the fraud case.

Furthermore, there are distinct “objects” whose normal behavior we wish to model:

1. Services (eg, telnet, http, etc.)
2. Ports
3. IP Addresses
4. System Command (in the case of BSM audit data)

Each of these is a component of a network computer system used in different ways during a *connection* as the *target* of or the *vehicle* of an attack by hacker perpetrators. Generally speaking, hackers are not attempting *per se* to steal money or abscond with free services (like phone calls). However, from the point of view of a site or organization under attack, there are real potential dollar losses due to

¹²Apparently when the US phone companies first fielded their analog based cell phone system without “authentication technology”, a flood of fraudulent calls by various clever means (including “phone cloning”) created heavy losses. The phone companies had to quickly catch up by creating their own fraud detection systems. Provost and Fawcett defined an automated and adaptive learning system to compute and field fraud detectors. However, technology advances in authenticating individual phones essentially rendered these detectors obsolete, or at least less relevant to the new generation of fielded systems.

intrusive connections. Hence, intrusion detection might best be formalized using appropriate cost-based models to detect attacks as in the credit card case.

The general approach here is to model known attacks (called *Misuse Detection*), or model “normal behavior” so that significant deviations from normalcy can warn of attacks (called *anomaly detection*). What constitutes “normal” behavior for each of these potential objects is unclear. Further, determining what “behavior” means, and what features should be computed to represent that behavior is not obvious either.

MADAM (Mining Audit Data for Automated Models for Intrusion Detection) is a set of new data mining algorithms that were developed by our project specifically to process network intrusion and audit data sets. MADAM includes variants of the “association rule” [?, ?] and “frequent episodes” [?, ?] algorithms used to define new *feature sets* that are extracted from labeled *tcpdump* data in order to define training sets for a machine learning algorithm to compute detectors. These features are defined over a set of *connections*. (A connection is the logical entity akin to the “credit card account” or the “cell phone number”.) With respect to *tcpdump* data streams, a connection is a sequence of tcp packets starting and ending at some well defined time, during which data flows to and from a source IP address to a target IP address under some well defined protocol. Connections are labeled as either attack (with the specific attack name or type) or “normal”. We therefore follow the same approach reported by Provost and Fawcett in mining feature sets (what they call *profilers*) from raw “transaction” data. We first determine what patterns of events in the raw stream appear to occur frequently in fraudulent transactions (or attack connections) that do not appear frequently in legitimate transactions (normal connections). These patterns of events define “features” computed for all connections used in training a classifier by some inductive inference or machine learning algorithm. The details of this data mining activity have been extensively reported [?, ?]. (Our previous exploratory work on learning anomalous Unix process execution traces can be found in [?].) Here we report a summary of our results.

4.1 The DARPA/MIT Lincoln Labs ID Evaluation

We participated in the 1998 DARPA Intrusion Detection Evaluation Program, prepared and managed by MIT Lincoln Labs. The objective of this program was to survey and evaluate research in intrusion detection. A standard set of extensively gathered audit data, which includes a wide variety of intrusions simulated in a military network environment, was provided by DARPA. Each participating site

was required to build intrusion detection models or tweak their existing system parameters using the training data, and send the results (i.e., detected intrusions) on the test data back to DARPA for performance evaluation.

We were provided with about 4 gigabytes of compressed raw (binary) *tcpdump* data of 7 weeks of network traffic, which can be processed into about 5 million connection records, each with about 100 bytes. The two weeks of test data have around 2 million connection records. Four main categories of attacks were simulated: DOS, denial-of-service, e.g., syn flood; R2L, unauthorized access from a remote machine, e.g., guessing password; U2R, unauthorized access to local superuser (root) privileges, e.g., various of “buffer overflow” attacks; and PROBING, surveillance and probing, e.g., port-scan.

Using the procedures reported in prior papers [?] we compared the aggregate normal pattern set with the patterns from each dataset that contains an attack type. The following features were constructed according to the intrusion only patterns:

- The “same host” features that examine only the connections in the past 2 seconds that have the same destination host as the current connection, and calculate statistics related to protocol behavior, service, etc.
- The similar “same service” features that examine only the connections in the past 2 seconds that have the same service as the current connection.

We call these the (time-based) “traffic” features of the connection records. There are several “slow” PROBING attacks that scan the hosts (or ports) using a much larger time interval than 2 seconds, for example, one in every minute. As a result, these attacks did not produce intrusion only patterns with a time window of 2 seconds. We sorted these connection records by the destination hosts, and applied the same pattern mining and feature construction process. Rather than using a time window of 2 seconds, we now used a “connection” window of 100 connections, and constructed a mirror set of “host-based traffic” features as the (time-based) “traffic” features.

We discovered that unlike most of the DOS and PROBING attacks, the R2L and U2R attacks don’t have any “intrusion only” *frequent* sequential patterns. This is because the DOS and PROBING attacks involve *many* connections to some host(s) in a very short period of time, the R2L and PROBING attacks are embedded in the data portions of the packets, and normally involves only a *single* connection. Algorithms for mining the unstructured data portions of packets are still under development. Presently, we use domain knowledge to add features

feature	description	value type
duration	length (number of seconds) of the connection	continuous
protocol_type	type of the protocol, e.g. tcp, udp, etc.	discrete
service	network service on the destination, e.g., http, telnet, etc.	discrete
src_bytes	number of data bytes from source to destination	continuous
dst_bytes	number of data bytes from destination to source	continuous
flag	normal or error status of the connection	discrete
land	1 - connection is from/to the same host/port; 0 - otherwise	discrete
wrong_fragment	number of “wrong” fragments	continuous
urgent	number of urgent packets	continuous

Table 7: Intrinsic Features: can be used for general-purpose traffic analysis.

that look for suspicious behavior in the data portion, e.g., number of failed login attempts. We call these features the “content” features.

A complete listing of the set of features defined for the connection records, i.e., its data schema, appears in Table 7 for the initial intrinsic set, while Table 8 describes the “content-based” features, and Table 9 displays the “traffic features” defined by the data mining process.

We then built three specialized models, using RIPPER [?]. RIPPER is a rule learning program that outputs a model quite similar in style to a Prolog program. Each model produced has a different set of features and detects different categories of intrusions. For example, for the “content” model, each connection record contains the “intrinsic” features and the “content” features, and the resultant RIPPER rules employing these features detects U2R and R2L attacks. A meta-classifier was trained and used to combine the predictions of the three base models when making a final prediction to a connection record. Table 10 summarizes these models. The numbers in bold, for example, **9**, indicate the number of automatically constructed temporal and statistical features being used in the RIPPER rules. We see that for both the “traffic” and host-based “traffic” models, our feature construction process contributes the majority of the features actually used in the rules.

Figure 4: ROC Curves on Detection Rates and False Alarm Rates

feature	description	value type
hot	number of “hot indicators”	continuous
num_failed_logins	number of failed login attempts	continuous
logged_in	1 - successfully logged in; 0 - otherwise	discrete
num_compromised	number of “compromised” conditions	continuous
root_shell	1 - root shell is obtained; 0 - otherwise	discrete
su_attempted	1 - “su root” command attempted; 0 - otherwise	discrete
num_root	number of “root” accesses	continuous
num_file_creations	number file creation operations	continuous
num_shells	number of shell prompts	continuous
num_access_files	number of operations on access control files	continuous
num_outbound_cmds	number of outbound commands in a ftp session	continuous
is_hot_login	1 - the login belongs to the “hot” list; 0 - otherwise	discrete
is_guest_login	1 - the login is a “guest”login; 0 - otherwise	discrete

Table 8: Content Features: expert knowledge on intrusion indicators within a TCP connection.

feature	description	value type
count	number of connections to the same host as the current connection in the past 2 seconds	continuous
	<i>the following features refer to these same-host connections</i>	
error_rate	% of connections that have “SYN” errors	continuous
error_rate	% of connections that have “REJ” errors	continuous
same_srv_rate	% of connections to the same service	continuous
diff_srv_rate	% of connections to different services	continuous
srv_count	number of connections to the same service as the current connection in the past 2 seconds	continuous
	<i>the following features refer to these same-service connections</i>	
srv_error_rate	% of connections that have “SYN” errors	continuous
srv_error_rate	% of connections that have “REJ” errors	continuous
srv_diff_host_rate	% of connections to different hosts	continuous

Table 9: Traffic Features: constructed from the “intrusion-only” frequent sequential patterns mined from the audit records.

Model	Feature set	Intrusion categories	# of features in records	# of rules	# of features used in rules
content	“intrinsic” + “content”	U2R, R2L	22	55	11
traffic	“intrinsic” + “traffic”	DOS, PROBING	20	26	4+9
host traffic	“intrinsic” + “host traffic”	Slow PROBING	14	8	1+5

Table 10: Model Complexities

4.2 Off-line Detection Results

We report here the performance of our detection models as evaluated by MIT Lincoln Labs. We trained our intrusion detection models, i.e., the base models and the meta-level classifier, using the 7 weeks of labeled data, and used them to make predictions on the 2 weeks of unlabeled test data. The test data contains a total of 38 attack types, with 14 types in the test data only (i.e., our models were not trained with instances of these attack types).

Figure 4 shows the ROC curves of the detection models by attack category as well as on all intrusions. In each of these ROC plots, the x-axis is the false alarm rate, calculated as the percentage of normal connections classified as an intrusion; the y-axis is the detection rate, calculated as the percentage of intrusions detected (since the models produced binary outputs, the ROC curves are not continuous). We compare here our models with other participants (denoted as Group 1 to 3) in the DARPA evaluation program (these plots are duplicated from the presentation slides of a report given by Lincoln Labs in a DARPA PI meeting. The slides can be viewed on line via <http://www.cs.columbia.edu/~sal/JAM/PROJECT/MIT/mit-index.html>). These groups primarily used knowledge engineering approaches to build their intrusion detection systems. We can see from the figure that our detection models have the best overall performance, and in all but one attack category, our model is one of the best two.

Table 11: Cost types in credit card fraud and network intrusion

Cost Type	Credit Card Fraud	Network Intrusion
Damage	$tranamt(t)$	$DCost(service, attack)$
Challenge	<i>overhead</i>	<i>overhead</i>
Operational	subsumed in <i>overhead</i>	<i>OpCost</i>

5 Formalizing Cost-based Models for Intrusion Detection

In the credit card fraud domain, the notion of costs is inextricably intertwined with the learning task. We seek to learn models of fraudulent transactions that minimizes the overall loss. We believe an analogous cost optimization problem can and should be defined for the IDS domain.

In the arena of IDS, there are at least three types of costs involved (that are derivative of the credit card fraud case):

1. “Damage” cost: the amount of damage caused by an attack if intrusion detection is not available or an IDS fails to detect an attack;
2. “Challenge” cost: the cost to act upon a potential intrusion when it is detected; and
3. “Operational” cost: the resources needed to run the IDS.

(A caveat, we have used dollars (\$) as a measure when discussing the credit card domain, but these costs can be abstracted to some other meaningful unit of measure of utility that may be more appropriate for the IDS case.)

Table 11 illustrates our perspective on the three types of cost in credit card fraud and intrusion detection. In the credit card case, “damage” is the amount of a fraudulent transaction that the bank losses, $tranamt(t)$. In the IDS case, damage can be characterized as a function that depends on the type of service and attack on that service, $DCost(service, attack)$. The challenge cost for both cases is term as *overhead*, which is the cost of acting on an alarm. We did not consider operational cost in the credit card case because we did not have the opportunity to study this aspect of the problem. The banks have existing fielded systems whose total aggregated operational costs have already been considered and are folded

into their *overhead* costs (here called the challenge cost). We shall take a limited view of this by considering the costs of alternative models based upon the “feature costs” used by these models employed in an IDS and we denote this operational cost as *OpCost*. We next elaborate on each of these sources of cost.

5.1 Damage costs

The damage cost characterizes the amount of damage inflicted by an attack when intrusion detection is unavailable (the case for most systems). This is important and very difficult to define since it is likely a function of the particulars of the site that seeks to protect itself. The defined cost function per attack or attack type should be used here to measure the cost of damage. This means, that rather than simply measuring *FN* as a rate of missed intrusions, rather we should measure total loss based upon $DCost(s, a)$, which varies with the service (*s*) and the specific type of attack (*a*). These costs are used throughout our discussion.

For instance, an undetected *Denial Of Service (DOS) attack* may have very high cost for a particular service, while a *Probe attack* may have only minor cost if not detected. That is, $DCost(http, DOS) > DCost(http, Probe)$. This begs a very tough question: how does an organization set its costs for each service, and each type of attack against that service? How does it decide the value of its services that may be attacked and lost, or measure the damage inflicted by a particular attack? Indeed, some attacks may destroy or deny access to all services.

Some may argue, however, that any type of attack on any service of a system has an equally high cost since the cost of repair or recovery is extremely high in any situation. We might, therefore, define a constant cost for all attacks and thus *TP/FN* and *FP/TN* rates tell the whole story about the performance of an IDS. (Recall, this is not the case for financial credit card transactions since transactions have varying dollar amounts.)

5.2 Challenge costs

The challenge cost is the cost to act upon an alarm that indicates a potential intrusion. For IDS, one might consider dropping or suspending a suspicious connection and attempting to check, by analyzing the service request, if any system data have been compromised, or system resources have been abused or blocked from other legitimate users. (Other personnel time costs can be folded in including gathering evidence for prosecution purposes if the intruder can be traced.) These costs can be estimated, as a first cut, by the amount of CPU and disk resources

needed to challenge a suspicious connection. For simplicity, instead of estimating the challenge cost for each intrusive connection, we can “average” (or amortize over a large volume of connections during some standard “business cycle”) the challenge costs to a single (but not static) *challenge cost per potential intrusive connection*, i.e., *overhead*.

5.3 Operational costs

The cost of fielding a detection system is interesting to consider in some detail. In the work on fraud detection in financial systems, we learned that there are a myriad of “business costs” involved in design, engineering, fielding and use (challenge) of detection systems. Each contributes to an overall aggregated cost of detecting fraud. The main issue in operational costs for IDS is the amount of resources to extract and test features from raw traffic data. Some features are costlier than others to gather, and at times, costlier features are more informative for detecting intrusions.

Real-time constraints in IDS. Even if one designs a good detection system that includes a set of good features that well distinguish among different attack types, these features may be infeasible to compute and maintain in real time. In the credit card case, transactions have a 5 second response constraint (a desired average waiting time). That’s a lot of time to look up, update and compute and test features, per transaction. (Throughputs are easily managed by disjoint partitioning of account records among distributed transaction processing sites.) In the IDS case, the desired average response rate should be measured in terms of average connection times, or even by TCP packet rates, a much smaller time frame, so connections can be dropped as quickly as possible before they do damage.

Furthermore, in the credit card case, a bank can (somewhat) easily “drop” a transaction, but still recover simply by reissuing the transaction at the point of service. Dropping connections in a network system might be dangerous. In fact, it might be the intent of a “denial of service” attack that we wish to detect and prevent.

In the case of IDS it is not obvious when an intrusion can be detected, and when an alarm should be issued. Ideally, we would like to detect and generate an alarm during an on-going attack connection in order to disable it, rather than after the fact when damage has already been done. However, certain models of intrusive connections may require information only known at the conclusion of a

connection! Thus, properly designing an intrusion detection system requires that considerable thought be given to the time at which a detection can and should take place.

Therefore, although the problems appear quite similar, the constraints are really much different between the two task domains. The problem seems to be much harder in the IDS case since we have to accommodate in our cost models the *response rate* of the system. It seems evident that a slower IDS should be penalized with a higher cost. (In the credit card case we simply ignored this cost.) This important source of cost however is a major topic of research for IDS, i.e. the *computational costs* for rapid detection.

Choosing effective, low cost features. Here is part of the story on computational costs. One needs to run the IDS on some system. That system running the IDS has a cost relative to the system(s) it is “protecting” on which the services are being provided.

Let’s say we reduce this to two different detection models that inspect an audit stream from the system providing services. Both IDS systems sample from the same audit stream, but then require different computations and storage to decide whether ongoing connections are attacks and whether to raise an alarm.

Let’s say the first IDS uses a model, M_1 , that inspects the audit stream and extracts a set of features, F_1 . Each feature value is extracted from the audit stream and subjected to a series of tests to generate a final outcome prediction, and an alarm decision is made. Taking say one day as the standard period of time by which we measure all IDS’s, we would attribute the cost $OpCost_1$ of M_1 for one day based upon the cost of auditing, extracting features, and evaluating the model. We do the same with M_2 , to evaluate $OpCost_2$.

Here is a specific example in our work. One set of features F_1 is extracted from the audit stream that provides information about an individual connection. F_1 provides information intrinsic to that connection only, eg., source and target IP addresses, the services accessed, port numbers, number of bytes transferred, etc. Model M_1 can detect attacks based only on this information, which we use under test to evaluate and measure its TP costs and FN costs, testing against some standard test set. This produces a final cost savings measurement for M_1 .

But now, M_2 is more sophisticated and may have the ability to better detect different attacks than those detected by M_1 . For example, some of the features we defined in our work involves not only information about a single connection, but also some statistical summaries about connections that occurred immediately

preceding the current connection by some number of seconds, i.e. M_2 is time dependent (what we call the “traffic features”)¹³ But, in order to compute F_2 , we required more storage, and more time to test the model, contributing increased costs to M_2 , and hence its total daily cost $OpCost_2$ may be higher than $OpCost_1$, i.e. the challenge costs vary with each detector. These additional costs may indeed be rational to pay if under the distribution of attacks, and considering TP and FN costs, it maximizes our cost savings.

As a simple first cut at defining the alternative computational costs we may measure storage resources for each model, and the time required to test a model, rather than the size of the model or some other syntactic measure¹⁴.

Here is an example: we may find that an IDS employs a model M that exhibits the best overall detection of DOS attacks, or the highest rated service with maximal TP cost attributed to it. But, M does its job well by first extracting a statistical temporal feature from an archive database of all telnet connections that occurred at the protected site within the last week. This information, coupled with additional information about the current connection being inspected, provides a highly accurate predictor of DOS attacks. Not bad we may say. However, on average, a connection may last, say 5 seconds. After packet filtering and assembly of connection records for all active connections, M has to do its work of accessing the archive data base, computing a statistical feature (say by sequential counting), moving bytes, and then performing a series of tests in order to generate a prediction about the current connection within 1 second, if not faster. And it must do this task for each active connection in our heavily loaded network! The cost of this computation may be very large in that we require a “supercomputer” to guar-

¹³A connection is established, and a login attempt fails, the connection is broken. In this case, M_1 may note the failed login, and assert a possible break-in. But the connection is now gone so nothing is done, no alarm is generated. M_2 however, keeps information about this broken connection, and now observes another connection also with a failed login. But now, the feature set F_2 includes a count of the number of prior connections with failed login attempts. This new simple temporal feature may be sufficient to sound the alarm to an administrator who now challenges the detected break-in connection. The attacker is caught in the act (and we saved some cost as a TP), or instead it turns out it was a legitimate user making a mistake and we lost the FP cost, *overhead*.

¹⁴It is interesting to point out that Stephanie Forrest’s work [?] is in this vein. She is attempting to generate very simple models that are fast to test and cheap in storage, with good detection capability. But we don’t know, accurately, how well her approach is doing with respect to accuracy, and certainly not the savings attributed to detected events. It may turn out that very cheap detectors based upon simple models do very well, or alternatively, a more expensive and complicated model may ultimately do better with respect to overall cost savings (not just “lightweight” costs).

antee the desired response rate for all connections in the network. The cost of the supercomputer naturally contributes to the cost of our model M , and hence the final assigned overhead cost, *overhead*. By the way, if the daily revenue generated by the protected services is sufficiently large, the cost of the supercomputer may indeed be warranted!

More realistically, we want alternative IDS's with alternative models that consume minimum resources but that still generate sufficiently good detection behavior under our cost models to make them worthwhile, i.e. maximize cost savings as best as possible.

Our work in this area is new and ongoing. Details of our initial thoughts here can be found in [?].

5.4 Cost Model for IDS

We just described the three different types of cost in IDS: damage cost, challenge cost, and operational cost. Our cost model for IDS considers these three types of cost. Similar to the credit card case, the IDS cost model depends on the outcomes of the IDS' predictions: false negative (FN), false positive (FP), true positive (TP), and true negative (TN). We now examine the cost associated with each of these outcomes.

FN cost, or the cost of NOT detecting an attack, is the most dangerous case (and is incurred by most systems today that do not field IDS's). Here, the IDS "Falsely" decides that a connection is not an attack and there is no challenge against the attack. This means the attack will succeed and do its dirty work and presumably some service will be lost, and the organization losses a service of some value. The FN Cost is, therefore, defined as the damage cost associated with the particular type of service and attack, $DCost(s, a)$.

TP Cost is the cost of detecting an attack and doing something about it, i.e. challenging it. Here, one hopes to stop an attack from losing the value of the service. There is a cost of challenging the attack, however, that is involved here.

When some event triggers an IDS to correctly predict that a True attack is underway (or has happened), then what shall we do? If the cost to challenge the attack is *overhead*, but the attack affected a service whose value is less than *overhead*, then clearly ignoring these attacks saves cost. Probes are often cited as doing no damage, and so they can be "reasonably" ignored. However, organizations may behave inappropriately and challenge probes at great cost anyway¹⁵.

¹⁵As an example, a famous research institution had at least 2 system administrators spend the

If indeed an organization attributes high costs even to Probes, then challenging probes is warranted. If not, then in evaluating an IDS, it is not wise to give it high scores for detecting all probes, but doing poorly perhaps on other more dangerous attacks. Therefore, for a true positive, if $overhead > DCost(s, a)$, the intrusion is not challenged and the loss is $DCost(s, a)$, but if $overhead < DCost(s, a)$, the intrusion is challenged and the loss is limited to $overhead$.

FP cost. Let's take a quick look at *FP*. When an IDS falsely accuses an event of being an attack, and the attack type is regarded as high cost, a challenge will ensue. We pay the cost of the challenge ($overhead$), but nothing really happened bad except we lost $overhead$ on the challenge. Naturally, when evaluating an IDS we have to concern ourselves with measuring this loss. For this discussion, we define the loss is just $overhead$ for a false positive.

TN cost. Then there is clearly the TN Cost that remains in this discussion. An IDS correctly decides that a connection is normal and Truly not an attack. We therefore bare no cost that is dependent on the outcome of an IDS.

Thus far we have only considered costs that depend on the outcome of an IDS, we now incorporate the operational cost, $OpCost$, that is independent of the IDS' predictive performance. Our notion of $OpCost$ mainly measures the cost of computing values of features in the IDS. We denote $OpCost(c)$ as the operational cost for a connection, c .

We now can describe the cost-model for IDS. When evaluating an IDS over some test set S of labeled connections, $c \in S$, we define the cumulative cost for a detector as follows:

$$CumulativeCost(S) = \sum_{c \in S} Cost(c) + OpCost(c) \quad (6)$$

where $Cost(c)$ is defined (analogous to the credit card case) as:

better part of at least two days tracking down some internal staff member who executed a probe for some reason. Once found, they simply berated the person who did it. How much did their organization spend and lose on their time for them to slap someone's wrist? Naturally, there is some benefit to creating at least an image of protection, but the bottom line here is that the organization lost real money on challenging what ultimately turned out to be a non-issue.

Outcome	$Cost(c)$
Miss (False Negative, FN)	$DCost(s, a)$
False Alarm (False Positive, FP)	$overhead$ if $DCost(s, a) > overhead$ or 0 if $DCost(s, a) \leq overhead$
Hit (True Positive, TP)	$overhead$ if $DCost(s, a) > overhead$ or $DCost(s, a)$ if $DCost(s, a) \leq overhead$
Normal (True Negative, TN)	0

where s is the service requested by connection c and a is the attack type detected by the IDS for the connection.

Note that a higher operational cost, $OpCost(c)$, could be incurred by employing “expensive” features; but this may potentially improve the predictive performance of the IDS and thus lower $Cost(c)$. Hence, in order to minimize $CumulativeCost(S)$, we need to investigate and quantify, in realistic contexts, the trade off between $OpCost(c)$ and $Cost(c)$ in Equation 6. This issue constitutes a major part of our ongoing research in the JAM project.

5.5 Flying Under Mobile Radar: Dynamic Overhead Adjustment

As in the credit card case, we can simplify the IDS cost model by subsuming the operational costs into *overhead* (challenge cost). In this way the cumulative cost of an IDS is highly dependent upon the *overhead*¹⁶ value set at the time models are computed, and certainly when they are evaluated. It is quite possible, and virtually assured, that under different *overheads*, different models will be computed and different cost performance will be exhibited.

In the context of credit cards, thieves have one important learning task to optimize their performance: they need to learn what the *overhead* may be at the time they are transacting to avoid being detected and possibly denied. Under the realistic cost models described in this paper, any transaction they create whose dollar value is under this *overhead* will be immediately approved, and thus they are immediately rewarded, and will naturally transact again. Their optimization problem is to execute as many transactions less than the *overhead* as possible to steal the maximum amount available¹⁷.

¹⁶and we may regard the overhead as the minimum height that a radar system protecting an area can see. Cruise missile programmers have an exceptionally interesting task of assuring their flying missiles stay below this radar to deliver their ordinance!

¹⁷Naturally, banks will not approve any transaction if the available credit has been used entirely.

In some industries (not only the credit card industry), *overheads* are so firmly fixed that fraud is simply modeled as another “cost of doing business” and is simply tolerated¹⁸. The *overhead* amount is defined by a myriad of business costs, but it need not be static when applied at run-time! Thus, it is quite logical to vary the *overhead* limit when operating an IDS, thus changing the challenge cost producing different behavior and cost performance of the detection system. This simple strategy tends to enter “noise” making it difficult for perpetrators to “optimize” their thefts.

But notice that under a changing *overhead*, either up or down, for which detectors had originally been trained, the outcome of cost savings attributed to the detector might vary widely. This change in *overhead* has another fundamental effect: it changes the environment from which our underlying distribution is drawn. This means, that once we lower the *overhead*, thieves might learn to lower their appetite for stealing not to get caught. Concurrently, raising the *overhead* afterwards might then generate large cost savings, until the thieves have learned to return to their former ways of being greedy.

An interesting question, therefore, is whether there is an optimal strategy of dynamically varying the *overhead* in order to maximize savings over a longer period of time. Varying the *overhead* implies that we must concern ourselves with potentially “non-linear effects” in cost savings. A slight reduction may indeed catch more fraud, but may result in far heavier losses due to the real costs of challenging a new found wealth of “cheap fraud”!

This begs further questions and deeper study to determine alternative strategies. Perhaps classifiers ought to be entirely retrained, or meta-classifiers might re-weight their constituent base classifiers under a new changing fraud and cost distribution, and when should we do this? Or, simply measuring daily cost savings performance, or the rate of change thereof, might provide interesting clues to an optimal daily setting? The rate at which we change our *overhead* setting, and/or our models to avoid widely varying oscillations in performance of overall cost savings is not obvious.

It is interesting to note here that one of the design goals of JAM is to provide a scalable, efficient and hence *adaptable* distributed learning system that provides

¹⁸For example, in the auto insurance industry, broken windshields are regarded as an immediately approved expense. Fraud perpetrators will submit insurance charges for bogus repairs of windshields and be assured of payment, simply because the cost of investigation is prohibitively expensive. Here thieves have a different problem. They need to learn the rate at which they submit bogus claims not to draw obvious attention to themselves from human claims processing personnel, the low bandwidth, final detectors.

the means of rapidly learning new classifiers, and distributing (via agent architectures) new detectors to accommodate changing conditions of the environment in which it operates. Another avenue for exploration in JAM is therefore to perhaps dynamically re-weight ensembles of classifiers, our meta-classifiers, to adjust to new *overhead* limits.

5.6 Summary

In our work on intrusion detection, the data mining activity was focussed on uncovering likely features to extract from the streaming TCP packets preprocessed into connection records that are used in preparing training data, computing models and testing those models.

However, much of the traditional research in modeling only considers statistical accuracy, or TP/FP rates of models when comparing approaches. We should now understand that accuracy is not the whole picture. In different real world contexts, “cost” can take on different meanings, and the target application might necessarily be defined as a cost optimization problem.

In the context of IDS, real time performance is crucial. Here cost measures involve throughput and memory resources. It is of no value if one has an IDS that consumes so much resource that services can no longer be delivered on time, or the cost of fielding the IDS is so high that it becomes uneconomical to do so.

The set of features proposed for building models of intrusions should be followed by feature pruning, meaning, some features may turn out to be irrelevant and can be comfortably dropped or they may be too costly to bother with. Dropping features means we do not need to incur the cost of computing them in real time during real time detection, thus reducing the overhead cost of detection. In certain algorithms, computing an underlying probability distribution over all proposed features may produce a better model (better with respect to accuracy), but the resultant model may be too expensive to compute in real time.

We used the RIPPER algorithm that includes a pruning phase when producing its final output model (that essentially throws away features by not including them in the output rules). It was a simple and effective expediency. Bayesian classifiers naturally can produce models that indicate most discriminating features from the max likelihood perspective. But, for IDS (and other related domains, perhaps all) maximum likelihood computed over weighted distributions of all available features needs to be traded-off with cost as well. This is a wide open issue for IDS.

Considering these real-time constraints, Moises Goldszmidt proposes the following interesting approach¹⁹.

In the considerations for real time computation and performance one must consider the time and effort required in computing the features.

Here is an approach taken directly from decision theoretic considerations, namely using value of information. In essence we want to compute whether it is worth to extract the next feature; namely, whether having information from another feature will improve the chances in detecting an intruder.

Lets assume that we have a Bayesian classifier that computes $P(I|F)$ and $P(-I|F)$ and then decides to raise an alarm based on the ratio (namely is it 2 to 1 or more, etc.) $\frac{P(I|F)}{P(-I|F)}$ (1) where I stands for intrusion and F is a vector of features for which we have information (namely we extracted their value).

The exact formula used, whether it is the ratio or the difference, is not important. This kind of formulation is also useful if one would like to take into account the various TP- and FP-costs previously described.

Now F is a subset of the actual features we can compute. We want to assess whether it is worth it to compute or extract feature F^* . For that we compute the expectation of (1), weighted by $P(F^*)$ (the probability that F^* gets a specific value); thus,

$$\sum_{F^*} \frac{P(I|F, F^*)}{P(-I|F, F^*)} \times P(F^*) \quad (2)$$

The value of information (VOI) is (2) - (1) (note that VOI is always > 0).

There are various issues to consider, however. We, of course need P . Now we can only get an estimate of P , and that would require to have a classifier (or density estimator) using all the features to continuously update P . For that we will have an offline estimator running in the background. There may be synchronization issues here however.

The computation of (1) and (2) (depending on the classifier) may need to marginalize over the features not seen which, in turn requires summation over a possibly exponential number of terms. There are ways around this including heuristics and caching.

In summary, the approach will consist of a model that is slowly but surely being updated to provide an approximation of P' and another one that works on line relying on the current estimation provided by P' to perform real-time value of information computation and providing Intrusion Detection services.

¹⁹Extracted from an email exchange nearly verbatim

6 Concluding Remarks

By way of summary, we enumerate a number of open research issues that need to be addressed in the general setting of distributed data mining, but also specific to the important task of detecting intrusions:

1. How does an organization or domain rationally set the costs of its various services and systems it wishes to protect with an IDS, thus defining $Cost(s, a)$ for all services and all attack types? And how do we rationally determine an overhead challenge cost, *overhead* especially under tough real-time constraints?
2. What “cost sensitive” data mining and machine learning algorithms are needed to generate “low cost” models; i.e. models that are cheap to evaluate and operate under (variable) “real-time” constraints, and that also maximize cost savings or minimize loss? AdaCost is one such proposed method of a “cost-sensitive” learning algorithm. Presently, AdaCost does not incorporate “feature costs” into its weight updating formulae.
3. Specifically for network-based intrusion detection, what is the optimal set of features to best model a “good detector” for different environments and platforms? In general, how does one rationally define the costs of models and associated “features” to be extracted from audit trails and used in those predictive models, in order to maximize an objective (cost saving) function for the organization?
4. Can conventional learning algorithms, that use statistical accuracy, information gain, or entropy, as their primary internal measure in strategically or heuristically choosing hypotheses, be easily “biased” by cost measures to generate effective cost saving models?
5. Can “cost-sensitive” meta-learning (in the sense of Chan&Stolfo) provide a means of combining “conventional models” to maximize cost savings or minimize loss?
6. Can strategically biasing training samples (according to different costs of attacks, or costs of services) produce effective models without necessarily altering or inventing new learning algorithms?

7. What are the (computational and predictive) performance tradeoffs between the approaches of computing “cost sensitive models”, or cost-based mixtures of conventional statistical models’?
8. The distribution of attacks, and the various costs associated with services and attacks will naturally change over time. What adaptive strategies might be needed to optimally change models or mixtures of models to improve detection and at what rate of change?
9. Likewise, what strategies may be employed in dynamically adjust overhead challenge costs (*overhead*) to maximize cost savings for a fixed detection system over larger time periods.

In conclusion, a microeconomic theory of intrusion detection poses significant new and interesting problems for data mining.

7 Acknowledgments

Our work has benefited from in-depth discussions with Alexander Tuzhilin of New York University and Foster Provost of Bell Atlantic research, and suggestions from Charles Elkan of UC San Diego.