

The Base-Rate Fallacy and the Difficulty of Intrusion Detection*

Stefan Axelsson
Department of Computer Engineering
Chalmers University of Technology
Göteborg, Sweden
Email: *sax@ce.chalmers.se*

Abstract

Many different demands can be made of intrusion detection systems. An important requirement is that it be *effective* i.e. that it should detect a substantial percentage of intrusions into the supervised system, while still keeping the *false alarm* rate at an acceptable level.

This paper aims to demonstrate that, for a reasonable set of assumptions, the false alarm rate is the limiting factor for the performance of an intrusion detection system. This is due to the base-rate fallacy phenomenon, that in order to achieve substantial values of the Bayesian detection rate, $P(\text{Intrusion}|\text{Alarm})$, we have to achieve a—perhaps in some cases unattainably—low false alarm rate.

A selection of reports of intrusion detection performance are reviewed, and the conclusion is reached that there are indications that at least some types of intrusion detection have far to go before they can attain such low false alarm rates.

1 Introduction

Many demands can be made of an intrusion detection system (IDS for short) such as *effectiveness*, *efficiency*, *ease of use*, *security*, *inter-operability*, *transparency* etc. Although much research has been done in the field in the past ten years, the theoretical limits of many of these parameters have not been studied to any significant degree. The aim of this paper is to discuss one serious problem with regard to the *effectiveness* parameter, especially how the base-rate fallacy may affect the operational effectiveness of an intrusion detection system.

*An earlier version of this paper was presented as [Axe99a]

2 Intrusion detection

The field of automated computer intrusion detection—intrusion detection for short—is currently some nineteen years old [And80], with interest gathering pace in the past ten years.

Intrusion detection systems are intended to help detect a number of important types of computer security violations, such as:

- Attackers using prepacked “exploit scripts.” Primarily outsiders.
- Attackers operating under the identity of a legitimate user, for example by having stolen that user’s authentication information (password). Outsiders and insiders.
- Insiders abusing legitimate privileges, etc.

Early work (see [And80, DN85, Den87, SSHW88]) identified two major types of intrusion detection strategies.

Anomaly detection The strategy of declaring everything that is unusual for the subject (computer, user, etc.) suspect, and worthy of further investigation. The early anomaly detection systems were all self-learning, i.e they automatically formed an opinion of what the subjects normal behaviour was.

Anomaly detection promises to detect abuses of legitimate privileges that cannot easily be codified into security policy, and to detect attacks that are “novel” to the intrusion detection system. Problems include a tendency to take up data processing resources, and the possibility of an attacker teaching the system that his illegitimate activities are nothing out of the ordinary.

Signature detection The detection strategy of deciding in advance what type of behaviour is undesirable, and through the use of predetermined signature of such behaviour, detecting intrusions.

Signature based detection systems promise to detect known attacks and violations easily codified into security policies in a timely and efficient manner. Problems include a difficulty in detecting previously unknown intrusions. If a database containing intrusion signatures is employed it must be updated frequently.

Early in the research it was suggested in [HK88, Lun88] that the two main methods ought to be combined to provide a complete intrusion detection system capable of detecting a wide array of different computer security violations, including the ones listed above.

We have previously written a survey and taxonomy of intrusion detection systems and principles [Axe99b], drawing from that we identify (at least) two other interesting categories in addition to the two above:

Specification based detection These can be described as a signature based system that operates under a *default deny* security policy, instead of the *default*

permit policy of *signature detection* above [KRL97, GWTB96]. That is, the user of the system codifies what the security benign behaviour of the supervised system (or subsystem) is, and the intrusion detection system flags all behaviour that deviates from this set norm as intrusive.

Even though the previous paragraph would indicate that specification based intrusion detection is a sub class of signature based detection, a more fruitful classification would put it with *anomaly detection* systems, since they by their very nature operate with an assumption of what is normal for the system, and flags deviations from this behaviour as intrusive. Thus we can form two major subgroups of the anomaly detection category, namely *self learning* versus *programmed* anomaly detection, specification based detection being an example of the latter.

It is clear that specification based detection shares some of the fundamental properties of anomaly detection, the ability to detect intrusions that are novel to it, for example. The area has not seen much research.

Classical detector For want of a better word. Classical detection and estimation theory teaches us that our detector should operate with knowledge of the characteristics of both the normal and intrusive process [Tre68]. All the detectors above operate with a (more or less clear) picture of either normal behaviour (in the case of anomaly detection) or intrusive behaviour (in the case of signature detection), but not both.

It is natural to assume that a detector that can decide whether some observed action falls in either one or the other class (or more correctly to which degree it falls in either, both, or none of these classes), having a model of both, could exhibit better detection and false alarm behaviour, and perhaps more importantly give a better estimate of its accuracy of detection. To date only one such intrusion detection system, a self learning system that learns by example (being fed flagged intrusions in background data), has seen the light of day [Lee99], and research in this area is still very immature.

We wish to make the division between different principles of detection above, since it is easy to conjecture that these fundamentally different modes of detection will exhibit different characteristics with regard to detection and false alarm rates. They probably also show different performance with regard to other parameters as well, such as runtime efficiency, but a discussion of these parameters fall outside the scope of this paper.

3 Problems in Intrusion Detection

At present, the many fundamental questions regarding intrusion detection remain largely unanswered. They include, but are by no means limited to:

Effectiveness How effective is the intrusion detection? To what degree does it detect intrusions into the target system, and how good is it at rejecting false positives, so called false alarms?

Efficiency What is the run time efficiency of the intrusion detection system, how many computing resources and how much storage does it consume, can it make its detections in real time, etc?

Ease of use How easy is it to field and operate for a user who is not a security expert, and can such a user add new intrusion scenarios to the system? An important issue in *ease of use* is the question of what demands can be made of the person responding to the intrusion alarm. How high a false alarm rate can he realistically be expected to cope with, and under what circumstances is he likely to ignore an alarm? (It has long been known in security circles that ordinary electronic alarm systems should be circumvented during normal operation of the facility, when supervisory staff are more likely to be lax because they are accustomed to false alarms [Pie48]).

Security When ever more intrusion detection systems are fielded, one would expect ever more attacks directed at the intrusion detection system itself, to circumvent it or otherwise render the detection ineffective. What is the nature of these attacks, and how resilient is the intrusion detection system to them?

Inter-Operability As the number of different intrusion detection systems increase, to what degree can they inter-operate and how do we ensure this?

Transparency How intrusive is the fielding of the intrusion detection system to the organisation employing it? How many resources will it consume in terms of manpower, etc?

While interest is being shown in some of these issues, with a few notable exceptions—mainly [HL93]—they remain largely unaddressed by the research community. This is perhaps not surprising, since many of these questions are difficult to formulate and answer.

This paper is concerned with one aspect of one of the questions above, that of *effectiveness*. More specifically it addresses the way in which the base-rate fallacy affects the required performance of the intrusion detection system with regard to false alarm rejection.

In what follows: section 4 gives a description of the base-rate fallacy, section 5 continues with an application of the base-rate fallacy to the intrusion detection problem, given a set of reasonable assumptions, section 6 describes the impact the previous results would have on intrusion detection systems, section 7 considers future work, with section 8 concluding the paper. Appendix A reproduces a base-rate fallacy example in diagram form.

4 The Base-Rate Fallacy

The base-rate fallacy¹ is one of the cornerstones of Bayesian statistics, stemming as it does directly from Bayes' famous theorem that states the relationship between a conditional probability and its opposite, i.e. with the condition transposed:

¹The idea behind this approach stems from [Mat96, Mat97].

$$P(A|B) = \frac{P(A) \cdot P(B|A)}{P(B)} \quad (1)$$

Expanding the probability $P(B)$ for the set of all n possible, mutually exclusive outcomes A we arrive at equation (2):

$$P(B) = \sum_{i=1}^n P(A_i) \cdot P(B|A_i) \quad (2)$$

Combining equations (1) and (2) we arrive at a generally more useful statement of Bayes' theorem:

$$P(A|B) = \frac{P(A) \cdot P(B|A)}{\sum_{i=1}^n P(A_i) \cdot P(B|A_i)} \quad (3)$$

The base-rate fallacy is best described through example.² Suppose that your doctor performs a test that is 99% accurate, i.e. when the test was administered to a test population all of whom had the disease, 99% of the tests indicated disease, and likewise, when the test population was known to be 100% free of the disease, 99% of the test results were negative. Upon visiting your doctor to learn the results he tells you he has good news and bad news. The bad news is that indeed you tested positive for the disease. The good news however, is that out of the entire population the rate of incidence is only 1/10000, i.e. only 1 in 10000 people have this ailment. What, given this information, is the probability of you having the disease? The reader is encouraged to make a quick “guesstimate” of the answer at this point.

Let us start by naming the different outcomes. Let S denote sick, and $\neg S$, i.e. *not* S , denote healthy. Likewise, let P denote a positive test result and $\neg P$ denote a negative test result. Restating the information above; given: $P(P|S) = 0.99$, $P(\neg P|\neg S) = 0.99$, and $P(S) = 1/10000$, what is the probability $P(S|P)$?

A direct application of equation (3) above gives:

$$P(S|P) = \frac{P(S) \cdot P(P|S)}{P(S) \cdot P(P|S) + P(\neg S) \cdot P(P|\neg S)} \quad (4)$$

The only probability above which we do not immediately know is $P(P|\neg S)$. This is easily found though, since it is merely $1 - P(\neg P|\neg S) = 1\%$ (likewise, $P(\neg S) = 1 - P(S)$). Substituting the stated values for the different quantities in equation (4) gives:

$$P(S|P) = \frac{1/10000 \cdot 0.99}{1/10000 \cdot 0.99 + (1 - 1/10000) \cdot 0.01} = 0.00980 \dots \approx 1\% \quad (5)$$

That is, that even though the test is 99% certain, your chance of actually having the disease is only 1/100, because the population of healthy people is much larger than the population with the disease. (For a graphical representation, in the form of a Venn diagram, depicting the different outcomes, turn to Appendix A). This result often surprises people, ourselves included, and it is this phenomenon—that humans in general do not take the basic rate of incidence, the base-rate, into account when

²This example hinted at in [RN95].

intuitively solving such problems of probability—that is aptly named “the base-rate fallacy.”

5 The Base-Rate Fallacy in Intrusion Detection

In order to apply this reasoning in computer intrusion detection we must first find the different probabilities, or if such probabilities cannot be found, make a set of reasonable assumptions regarding them.

5.1 Basic Frequency Assumptions

Let us for the sake of further argument hypothesize a figurative computer installation with a few tens of workstations, a few servers—all running UNIX—and a couple of dozen users. Such an installation could produce in the order of 1,000,000 audit records per day with some form of “C2” compliant logging in effect, in itself a testimony to the need for automated intrusion detection.

Suppose further that in such a small installation we would not experience more than a few, say one or two, actual attempted intrusions per day. Even though it is difficult to get any figures for real incidences of attempted computer security intrusions, this does not seem to be an unreasonable number.

Furthermore, assume that at this installation we do not have the manpower to have more than one site security officer—SSO for short—who probably has other duties, and that the SSO, being only human, can only react to a relatively low number of alarms, especially if the false alarm rate is high (50% or so), see section 5.2.

Even though an intrusion could possibly affect only one audit record, it is likely on average that it will affect a few more than that. Furthermore, a clustering factor actually makes our estimates more conservative, so it was deemed prudent to include one. Using data from a previous study of the trails that SunOS intrusions leave in the system logs [ALGJ98], we can estimate that ten audit records would be affected in the average intrusion.

5.2 Human Machine Interaction in Intrusion Detection

The previous assumptions above are “technical” in nature, i.e. anyone well versed in the field of computer security can make similar predictions, or adjust the ones above to suit his liking. It is a simple matter to verify or predict similar measures. However, the factor of the performance of the human operator does not lend itself to the same technological estimates. Thus, a crucial question is the question of the capacity of the human operator to correctly respond to the output of the system. Especially his capacity to tolerate false alarms.

Unfortunately there have been no experiments concerning these factors in the setting of computer security intrusion detection. There is, however, some research in the context of process automation and plant control, such as would be the case in a (nuclear) power station, paper mill, steel mill, large ship etc.

On the over all level research has shown that a human operator (decision maker) in such an environment has to [Ras86, p. 5]:

Detect the need for intervention and to

observe important data in order to have direction for subsequent activities. He then has to analyse the data in order to

identify the present state of affairs and to

evaluate their possible consequences with reference to operational goals and company policies. Then a

target state into which the system should be transferred has to be chosen, and the

task that the decision maker has to perform is selected from a review of the resources available to reach said target state. When the task has been identified the proper

procedure i.e. how to do it, must be planned and executed.

In our case we have chosen to aid the operator with an intrusion detection system. However we quickly notice the absence of any discussion about the rest of the decision making chain—even though the recovery element has seen some general study—when it comes to the research into human interaction with intrusion detection systems. Most authors don’t even discuss the second step in recovery above, namely that of aiding the operator with *observations* about the state of the system. (The normal state of which is most often not known in our case. No-one knows what the traffic on our computer networks typically looks like, hence the reported difficulty of even deciding if something really is amiss [Sto95].)

More specifically, in this particular case, we are interested in the operators ability to act “correctly” in the presence of false alarms. I.e. how many false alarms an operator can tolerate without losing his vigilance.

This is a difficult question to answer in this particular context, not only because there has been no research into the question. A few difficulties are:

- First, the modeling of the human operator handling such a highly complex and cognitive task as the detection and resolution of a computer security incident is difficult in general terms. It is doubtful that we will ever reach a quantitative model of human performance and limitations in this area. We can make several qualitative statements however [Wic92, pp. 258].
- Second, several different factors influence the performance of the operator at different times, such as previous experience, level of training, work load, external and internal stressors, state of vigilance etc.
- Third, the human operator is prone to several different kinds of bias when making a decision of this kind, biases relating to his inability to correctly make statistical estimates, of making correct logical inferences etc. From our perspective the bias of tending to stay with the original hypothesis (that no intrusion has taken place in our case) and not seek disconfirmatory evidence is especially interesting to us [Wic92, pp. 280].

What previous research in other areas seem to tell us specifically about our situation, is that human operators tend to have a very low tolerance for false alarms. During normal operation, humans have a tendency to overtrust the infallibility of the automated equipment. However once the equipment is seen to malfunction (raise false alarms in our case) humans tend to mistrust the equipment to a larger degree than what would be warranted by its actual performance. “Trust once betrayed is hard to recover” [Wic92, p. 537] Perhaps surprisingly, there has been little empirical research in this area [Nyg94, Wic92, p. 537].

What studies have been made [Nyg94, Dea72], seem to indicate that our required level of false alarms, 50%, is a *very* conservative estimate. Most human operators will have completely lost faith in the device at that point, opting to treat every alarm with extreme scepticism, if one would be able to speak of a “treatment” at all, the intrusion detection system would most likely be completely ignored in a “civilian” setting.

5.3 Calculation of Bayesian Detection Rates

Let I and $\neg I$ denote *intrusive*, and *non-intrusive* behaviour respectively, and A and $\neg A$ denote the presence or absence of an intrusion alarm. We start by naming the four possible cases (false and true positives and negatives) that arise by working backwards from the above set of assumptions:

Detection rate Or *true positive* rate. The probability $P(A|I)$, i.e. that quantity that we can obtain when testing our detector against a set of scenarios we know represent intrusive behaviour.

False alarm rate The probability $P(A|\neg I)$, the *false positive* rate, obtained in an analogous manner.

The other two parameters, $P(\neg A|I)$, the *False Negative* rate, and $P(\neg A|\neg I)$, the *True Negative* rate, are easily obtained since they are merely:

$$P(\neg A|I) = 1 - P(A|I); P(\neg A|\neg I) = 1 - P(A|\neg I) \quad (6)$$

Of course, our ultimate interest is that both:

- $P(I|A)$ —that an alarm really indicates an intrusion (henceforth called the *Bayesian detection rate*), and
- $P(\neg I|\neg A)$ —that the absence of an alarm signifies that we have nothing to worry about,

remain as large as possible.

Applying Bayes’ theorem to calculate $P(I|A)$ results in:

$$P(I|A) = \frac{P(I) \cdot P(A|I)}{P(I) \cdot P(A|I) + P(\neg I) \cdot P(A|\neg I)} \quad (7)$$

Likewise for $P(\neg I|\neg A)$:

$$P(\neg I|\neg A) = \frac{P(\neg I) \cdot P(\neg A|\neg I)}{P(\neg I) \cdot P(\neg A|\neg I) + P(I) \cdot P(\neg A|I)} \quad (8)$$

These assumptions give us a value for the rate of incidence of the actual number of intrusions in our system, and its dual (10 audit records per intrusion, 2 intrusions per day, and 1,000,000 audit records per day). Interpreting these as probabilities:

$$\begin{aligned} P(I) &= 1 \bigg/ \frac{1 \cdot 10^6}{2 \cdot 10} = 2 \cdot 10^{-5}; \\ P(\neg I) &= 1 - P(I) = 0.99998 \end{aligned} \quad (9)$$

Inserting equation (9) into equation (7):

$$P(I|A) = \frac{2 \cdot 10^{-5} \cdot P(A|I)}{2 \cdot 10^{-5} \cdot P(A|I) + 0.99998 \cdot P(A|\neg I)} \quad (10)$$

Studying equation (10) we see the base-rate fallacy clearly. By now it should come as no surprise to the reader, since the assumptions made about our system makes it clear that we have an overwhelming number of non-events (benign activity) in our audit trail, and only a few events (intrusions) of any interest. Thus, the factor governing the *detection* rate ($2 \cdot 10^{-5}$) is completely dominated by the factor (0.99998) governing the *false alarm* rate. Furthermore, since $0 \leq P(A|I) \leq 1$, the equation will have its desired maximum for $P(A|I) = 1$ and $P(A|\neg I) = 0$, which results in the most beneficial outcome as far as the *false alarm* rate is concerned. While reaching these values would be an accomplishment indeed, they are hardly attainable in practice. Let us instead plot the value of $P(I|A)$ for a few fixed values of $P(A|I)$ (including the “best” case $P(A|I) = 1$), as a function of $P(A|\neg I)$ (see figure 1 on the following page). It should be noted that both axes are logarithmic.

It becomes clear from studying the plot in figure 1 that even for the unrealistically high *detection* rate 1.0, we have to have a very low *false alarm* rate (on the order of $1 \cdot 10^{-5}$) for the Bayesian detection rate to have a value of 66%, i.e. about two thirds of all alarms will be a true indication of intrusive activity. With a more realistic *detection* rate of, say, 0.7, for the same *false alarm* rate, the value of the Bayesian detection rate is about 58%, nearing fifty-fifty. Even though the number of events (intrusions/alarms) is still low, it is our belief that a low Bayesian detection rate would quickly “teach” the SSO to (un)safely ignore *all* alarms, even though their absolute numbers would theoretically have allowed a complete investigation of all alarms. This becomes especially true as the system grows; a 50% false alarm rate of in total of 100 alarms would clearly not be tolerable. Note that even quite a large difference in the *detection* rate does not substantially alter the Bayesian detection rate, which instead is dominated by the *false alarm* rate. Whether such a low rate of false alarms is at all attainable is discussed in section 6.

It becomes clear that, for example, a requirement of only 100 false alarms per day is met by a large margin with a *false alarm* rate of $1 \cdot 10^{-5}$. With 10^5 “events” per day, we will see only 1 *false alarm* per day, on average. By the time our ceiling of 100 false alarms per day is met, at a rate of $1 \cdot 10^{-3}$ *false alarms*, even in the best

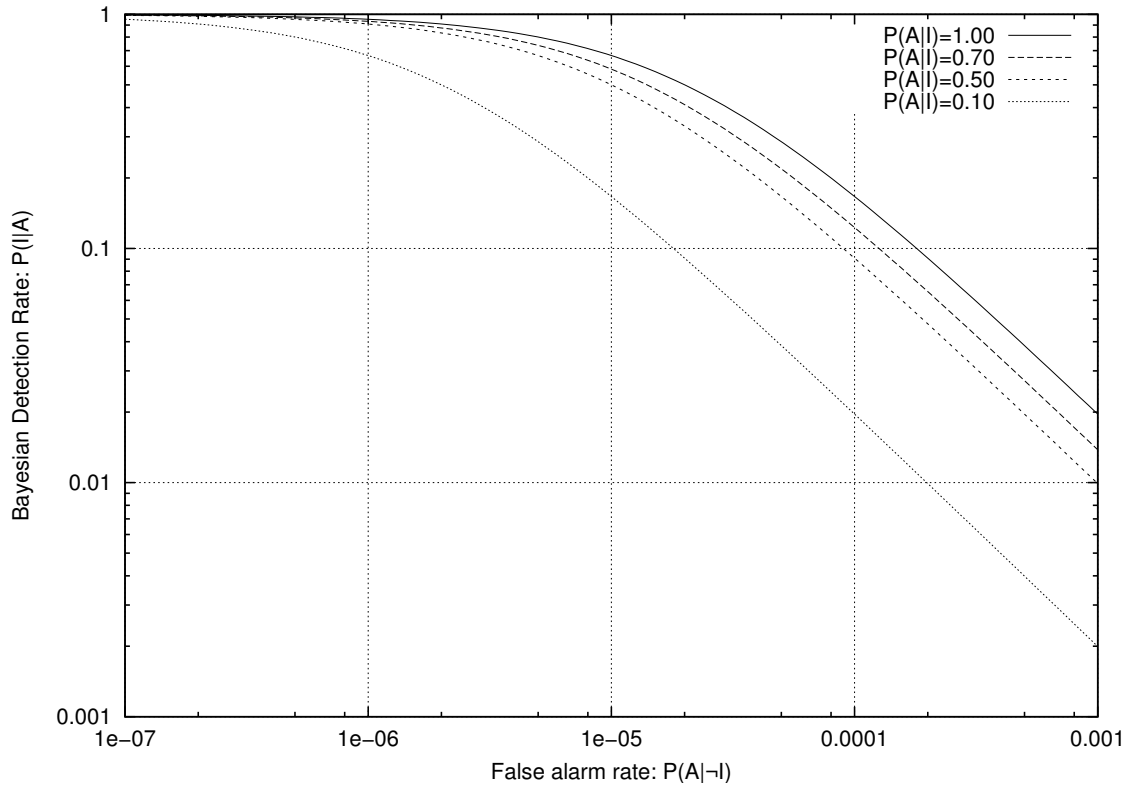


Figure 1: Plot of Bayesian detection rate versus false alarm rate

case scenario, our Bayesian detection rate is down to around 2%,³ by which time no-one will care less when the alarm goes off.

Substituting (6) and (9) in equation (8) gives:

$$P(\neg I|\neg A) = \frac{0.99998 \cdot (1 - P(A|\neg I))}{0.99998 \cdot (1 - P(A|\neg I)) + 2 \cdot 10^{-5} \cdot (1 - P(A|I))} \quad (11)$$

A quick glance at the resulting equation (11) raises no cause for concern. The large $P(\neg I)$ factor (0.99998) will completely dominate the equation, giving it values near 1.0 for the values of $P(A|\neg I)$ under discussion here, regardless of the value of $P(A|I)$.

This is the base-rate fallacy in reverse, if you will, since we have already demonstrated that the problem is that we will set off the alarm too many times in response to non-intrusions, combined with the fact that we do not have many intrusions to begin with. Truly a question of finding a needle in a haystack.

The author does not see how the situation underlying the base-rate fallacy problem will change for the better in years to come. On the contrary, as computers get faster they will produce more audit data, while it is doubtful that intrusive activity will increase at the same rate. In fact, it would have to increase at a substantially higher rate for it to have any effect on the previous calculations, and were it ever to reach levels sufficient to have such an effect—say 30% or more—the installation

³Another way of calculating that differs from equation (10) is of course to realise that 100 false alarms and only a maximum of 2 possible valid alarms gives: $\frac{2}{2+100} \approx 2\%$.

would no doubt have a serious problem on its hands, to say the least!

6 Impact on Intrusion Detection Systems

As stated in the introduction, approaches to intrusion detection can be divided into three major groups, *signature*-based, *anomaly*-based, and *classical detectors*. The previous section developed requirements regarding *false alarm* rates and *detection* rates in intrusion detection systems in order to make them useful in the stated scenario. This section will compare these requirements with reported results on the effectiveness of intrusion detection systems. For a survey of research into the testing of intrusion detection systems, see [AS99].

It can be argued that this reasoning does not apply to anomaly-based intrusion detection. In some cases anomaly-based detection tries not to detect intrusions *per se*, but rather to differentiate between two different subjects, flagging anomalous behaviour in the hopes that it is indicative of a stolen user identity for instance, see for example [LB98], which even though it reports performance figures, is not directly applicable here. However, we think the previous scenario is useful as a description of a wide range of more “immediate,” often network-based, attacks, where we will not have had the opportunity to observe the intruder for an extended period of time “prior” to the attack.

Another paper that discusses the effectiveness of intrusion detection is [Max98]. Unfortunately it is not applicable here.

6.1 ROC Curve Analysis

There are general results in detection and estimation theory that state that the *detection* and *false alarm* rates are linked [Tre68], though the extent to which they are applicable here is still an open question. Obviously, if the *detection* rate is 1, saying that all events are intrusions, we will have a *false alarm* rate of 1 as well, and conversely the same can be said for the case where the rates are 0.⁴ Intuitively, we see that by classifying more and more events as intrusive—in effect relaxing our requirements on what constitutes an intrusion—we will increase our *detection* rate, but also misclassify more of the benign activity, and hence increase our *false alarm* rate.

Plotting the *detection* rate as a function of the *false alarm* rate we end up with what is called a ROC—Receiver Operating Characteristic—curve. (For a general introduction to ROC curves, and detection and estimation theory, see [Tre68].) We have already stated that the points (0; 0) and (1; 1) are members of the ROC curve for any intrusion detector. Furthermore, the curve between these points is convex; were it concave, we would do better to reverse our decision. Nor can it contain any dips, as that would in effect indicate a faulty, non-optimal detector, since a randomised test would then be better. See “Assumed ROC” curve in figures 2 and 3 for the ROC curve that depicts our previous example.

⁴If you call everything with a large red nose a clown, you’ll spot all the clowns, but also Santa’s reindeer, Rudolph, and vice versa.

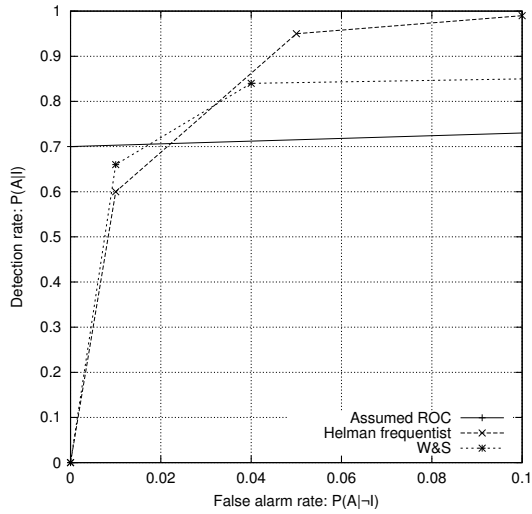


Figure 2: ROC-curves for the “low performers”

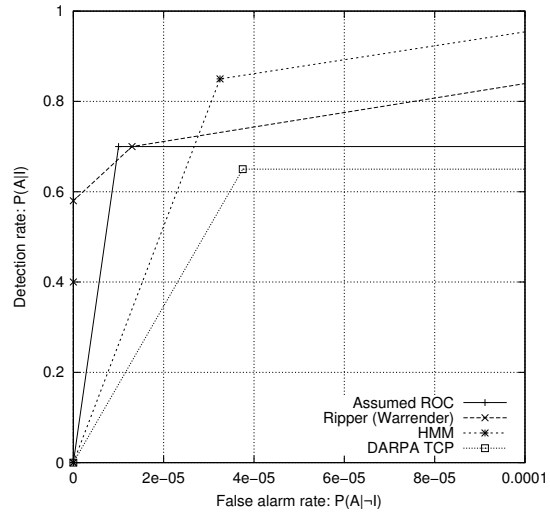


Figure 3: ROC-curve for the “high performers”

We see that the required ROC curve has a very sharp rise from $(0;0)$ since we quickly have to reach acceptable *detection* rate values (0.7) while still keeping the *false alarm* rate under control.

6.2 Previous Experimental Intrusion Detection Evaluations

As previously mentioned, the literature is not overladen with experimental results from tests of intrusion detection systems. Ideally we would like several different results from the four different classes in section 2. Unfortunately there only exists one report of anomaly detection performance in this regard (though with a strong theoretical foundation) [HL93] with no example of specification based intrusion detection, and one independent report of a classical detector [WFP99]. Several signature based detectors have been tested for DARPA by Lincoln labs however [GLC⁺98].

Unfortunately no comprehensive results of the (semi)recent evaluation performed by DARPA by Lincoln Labs at MIT [LGG⁺98, GLC⁺98] have been published, and the data is furthermore unavailable to us for independent evaluation because of U.S. export restrictions.

What has been made known on the web however indicate that the study was conducted using a simulated network of workstations, transmitting simulated traffic. This traffic was generated based on real traffic observed on a large US Air Force base, and a large research institute. This of course lends some credibility to an argument about the generality of the background traffic, but no such argument is made by the authors. Of course, the degree to which the background traffic is representative of the background traffic in the field is a crucial question when it comes to the value of the test as an indicator of false alarm rates during normal usage.

In the test, a number of different attacks were then inserted into the simulated network, including denial of service attacks against the network, and “root” exploits against individual workstations. The experimenters invited several different intru-

sion detectors to participate in the study. These were all signature based detectors operating on either network or host data. Even though there is more going on behind the scenes (the detection rate varies between approximately 20%–90% for the best scoring detector for all attacks) we will limit the presentation the best overall scores for the conglomerate of detectors in the network study, i.e. the detector resulting from combining the four different detectors and choosing the best performer in all instances. Note that this may not be realistic, since it would be difficult to perform this conglomeration in practice, to say the least.

Also not all detectors performed equally well when dealing with all intrusions, and it is a general criticism that in the case of signature based detection, the designer of the signature can easily trade off detection rate for false alarm rate by varying the generality of the signature. The more general, abstract if you will, it is, the more variations of the same intrusive behaviour it will detect, but at the cost of a higher false alarm rate. It is not known to what extent the DARPA evaluation used variations of the attacks presented to the designers of the intrusion detection systems for training purposes, in the final evaluation. This is an important point in that when such systems are commercialised, it will be impossible to keep the detection signatures secret from the would be intruders, and the more savvy among them will of course attempt to vary their techniques to evade the intrusion detection system.

Much more can be said about this evaluation, but we will limit our comments to the above. Of course choosing the best performer makes our comparison more conservative.

The second study [WFP99] lists test results for six different intrusion detection methods that have been applied to traces of system calls made into the operating system kernel by nine different privileged applications in a UNIX environment. Most of these traces were obtained from “live” data sources, i.e. the systems from which they were collected were production systems. The authors’ hypothesis is that short sequences of system calls exhibit patterns that describe normal, benign activity, and that different intrusion detection mechanisms can be trained to detect abnormal patterns, and flag these as intrusive. The researchers thus trained the intrusion detection systems using part of the “normal” traffic, and tested their false alarm rate on the remaining “normal” traffic. They then trained the systems on intrusive scenarios, and inserted such intrusions into normal traffic to ascertain the detection rate. The experimental method is thus close to the one described in sections 4 and 5. This study evaluated as one of the systems the self learning “classical” detector, RIPPER, described by Lee [Lee99].

The third study [HL93] is a treatise on the fundamental limits of the effectiveness of intrusion detection. The authors constructs a model of the intrusive and normal process and investigate the properties of this model from an anomaly intrusion detection perspective under certain assumptions. Their approach differs from ours in that they do not provide any estimates of the parameters in their model, opting instead to explore the limits of effectiveness when such information is unavailable. Of greatest interest here is their conclusion in which the authors plot experimental data for two implementations, one a frequentist detector that—it is claimed—is close to optimal under the given circumstances, and an earlier tool designed by the authors,

Wisdom & Sense [VL89]. Unfortunately, only one type of anomaly detection system, one that operates with descriptive statistics of the behaviour of the subject, is covered. As previously mentioned, specification based intrusion detection is not covered, and furthermore, neither are more “sophisticated” detectors, such as neural network based detectors (such as [DBS92]), that take time series behaviour of the subject into account.

Lack of space precludes a more detailed presentation of these experiments, and the interested reader is referred to the cited papers where available.

The results from the three studies above have been plotted in figures 2 and 3. Where a range of values were given in the original presentation, the best—most “flattering” if you will—value was chosen. Furthermore, since not all the work cited to provided actual numerical data, some points are based on our interpretation of the presented values. In the case of the DARPA study the results were rescaled to conform with our requirements. (The original DARPA test assumes 66,000 events per day instead of our 100,000 events per day.) We feel that these are accurate enough for the purpose of giving the reader an idea of the performance of the systems.

The cited work can be roughly divided into two classes depending on the minimum false alarm rate values that are presented, and hence, for clarity, the presentation has been divided into figures, where the first (figure 2) presents the first class, with larger values for the false alarm rate. These consists solely of the anomaly detection results in this study. In the figure “Helman frequentist,” and “W&S” denote the detection results from [HL93]. It is interesting, especially in the light of the strong claims made by the authors of this evaluation, to note that all of the presented false alarm rates are several orders of magnitude larger than the requirements put forth in section 5.

The second class of detectors, depicted in figure 3, consists of the average results of Ripper [Lee99], a high performance Hidden Markov Model detector (labeled “HMM” in the figure) tested by Warrander et. al. in [WFP99], and the DARPA results. Here the picture is less clear. The authors report false alarm results close to zero for lower detection rates, with one performance point nearly overlapping our required performance point. The HMM detector is also close to what we would require. It is more difficult to generalize these results, since they are based on one method of data selection, and the authors do not make as strong a claim as those made for the previous set of detectors. The DARPA data from [GLC⁺98], show up as “DARPA TCP” in figure 3. They are also in the vicinity of the required performance point, but the question of the generality of the training/test data, and hence the results, remains.

7 Future Work

One sticking point is the basic probabilities that the previous calculations are based on. These probabilities are subjective at present, but future work should include measurement either to attempt to calculate these probabilities from observed frequencies—the *frequentist* approach—or to deduce these probabilities from some model of the intrusive process and the intrusion detection system—the *objectivist* approach. The latter would in turn require real world observation to formulate realistic param-

eters for the models.

Furthermore, this discourse treats the intrusion detection problem as a binary decision problem, i.e. that of deciding whether there has been an “intrusion” or not. The work presented does not differentiate between the different kinds of intrusions that can take place, and nor does it recognise that different types of intrusions are not equally difficult or easy to detect. Thus on a more detailed level, the intrusion detection problem is not a binary but rather an n -valued problem.

Another area that needs attention is that of the SSO’s capabilities. How does the human-computer interaction take place, and precisely which Bayesian detection rates would an SSO tolerate under what circumstances for example?

The other parameters discussed in the introduction (*efficiency*, etc.) also need further attention.

8 Conclusions

This paper aims to demonstrate that intrusion detection in a realistic setting is perhaps harder than previously thought. This is due to the base-rate fallacy problem, because of which the factor limiting the performance of an intrusion detection system is not the ability to identify behaviour correctly as intrusive, but rather *its ability to suppress false alarms*. A very high standard, less than 1/100,000 per “event” given the stated set of circumstances, will have to be reached for the intrusion detection system to live up to these expectations as far as *effectiveness* is concerned.

The cited studies of intrusion detector performance that were plotted and compared indicate that anomaly-based methods may have a long way to go before they can reach these standards, since their false alarm rates are several orders of magnitude larger than what we demand. When we come to the case of signature-based detection methods the picture is less clear. Even though the cited work seems to indicate that current signature intrusion detectors can operate close to the required performance point, how well these results generalise in the field is still an open question. We only have one data point when it comes to the more qualified “classical” detectors, and it seems to perform on par with signature based detectors.

Of course whether some of the more difficult demands, such as the detection of masqueraders or the detection of novel intrusions, can be met without the use of anomaly-based intrusion detection is still an open and interesting question.

Much work still remains before it can be demonstrated that current IDS approaches will be able to live up to real world expectations of effectiveness. However, we would like to stress that, the present results notwithstanding, an equal amount of work remains before it can be proven that they *cannot* live up to such high standards.

9 Acknowledgements

I would like to thank my colleague Ulf Lindqvist and my supervisor Erland Jonsson for valuable insights. I would also like to thank the anonymous reviewers for their comments and suggestions.

This work was funded by the Swedish National Board for Industrial and Technical Development (NUTEK) under project P10435.

References

- [ALGJ98] Stefan Axelsson, Ulf Lindqvist, Ulf Gustafson, and Erland Jonsson. An approach to UNIX security logging. In *Proceedings of the 21st National Information Systems Security Conference*, pages 62–75, Crystal City, Arlington, VA, USA, 5–8 October 1998. NIST, National Institute of Standards and Technology/National Computer Security Center.
- [And80] James P. Anderson. Computer security threat monitoring and surveillance. Technical Report Contract 79F26400, James P. Anderson Co., Box 42, Fort Washington, PA, 19034, USA, 26 February revised 15 April 1980.
- [AS99] Dan Andersson and Heléne Svensson. Testing of intrusion detection systems—A survey. In Louise Yngström and Thomas Svensson, editors, *Proceedings of the fourth Nordic Workshop on Secure IT systems—Encouraging Co-operation*, 99–005, pages 165–179. Department of Computer and Systems Sciences, Stockholm University and Royal Institute of Technology, Sweden, November 1–2 1999. ISSN 1101–8526, ISBN 91–7153–955–7.
- [Axe99a] Stefan Axelsson. The base-rate fallacy and its implications for the difficulty of intrusion detection. In *6th ACM Conference on computer and communications security*, pages 1–7, Kent Ridge Digital Labs, Singapore, 1–4 November 1999.
- [Axe99b] Stefan Axelsson. Intrusion detection systems: A taxonomy and survey. Technical Report 99–15, Department of Computer Engineering, Chalmers University of Technology, SE-412 96 Göteborg, Sweden, 1999. URL: <http://www.ce.chalmers.se/staff/sax>.
- [DBS92] Herve Debar, Monique Becker, and Didier Siboni. A neural network component for an intrusion detection system. In *Proceedings of the 1992 IEEE Computer Society Symposium on Research in Security and Privacy*, pages 240–250, Oakland, CA, USA, May 1992. IEEE, IEEE Computer Society Press, Los Alamitos, CA, USA.
- [Dea72] B. H. Deatherage. Auditory and other sensory forms of information. In HP Van Cott and RG Kinkade, editors, *Human Engineering Guide to Equipment design*. Army, Navy, Air Force, 1972.
- [Den87] Dorothy E. Denning. An intrusion-detection model. *IEEE Transactions on Software Engineering*, Vol. SE-13(No. 2):222–232, February 1987.

- [DN85] Dorothy E. Denning and Peter G. Neumann. Requirements and model for IDES—A real-time intrusion detection system. Technical report, Computer Science Laboratory, SRI International, Menlo Park, CA, USA, 1985.
- [GLC⁺98] Isaac Graf, Richard Lippman, Robert Cunningham, David Fried, Kris Kendall, Seth Webster, and Marc Zissman. Results of DARPA 1998 offline intrusion detection evaluation. <http://www.ll.mit.edu/IST/ideval>, December 15 1998.
- [GWTB96] Ian Goldberg, David Wagner, Randi Thomans, and Eric Brewer. A secure environment for untrusted helper applications (confining the wily hacker). In *Proceedings of the Sixth USENIX UNIX Security Symposium*, San Jose, California, USA, July 1996. USENIX, USENIX Association.
- [HK88] L. Halme and B. Kahn. Building a security monitor with adaptive user work profiles. In *Proceedings of the 11th National Computer Security Conference*, Washington DC, October 1988.
- [HL93] Paul Helman and Gunar Liepins. Statistical foundations of audit trail analysis for the detection of computer misuse. *IEEE Transactions on Software Engineering*, 19(9):886–901, September 1993.
- [KRL97] Calvin Ko, M. Ruschitzka, and K Levitt. Execution monitoring of security-critical programs in distributed systems: A specification-based approach. In *Proceedings of the 1997 IEEE Symposium on Security and Privacy*, volume ix, pages 175–187, Oakland, CA, USA, May 1997. IEEE, IEEE Computer Society Press, Los Alamitos, CA, USA. IEEE Cat. No. 97CB36097.
- [LB98] Terran Lane and Carla E. Brodie. Temporal sequence learning and data reduction for anomaly detection. In *5th ACM Conference on Computer & Communications Security*, pages 150–158, San Francisco, California, USA, 3–5 November 1998.
- [Lee99] Wenke Lee. A data mining framework for building intrusion detection models. In *IEEE Symposium on Security and Privacy*, pages 120–132, Berkeley, California, May 1999.
- [LGG⁺98] Richard P. Lippmann, Isaac Graf, S. L. Garfinkel, A. S. Gorton, K. R. Kendall, D. J. McClung, D. J. Weber, S. E. Webster, D. Wyschogrod, and M. A. Zissman. The 1998 DARPA/AFRL off-line intrusion detection evaluation. The First Workshop on Recent Advances in Intrusion Detection (RAID-98), Lovain-la-Neuve, Belgium, 14–16 September 1998.
- [Lun88] Teresa F Lunt. Automated audit trail analysis and intrusion detection: A survey. In *Proceedings of the 11th National Computer Security Conference*, pages 65–73, Baltimore, Maryland, 17–2 October 1988. NIST.

- [Mat96] Robert Matthews. Base-rate errors and rain forecasts. *Nature*, 382(6594):766, 29 August 1996.
- [Mat97] Robert Matthews. Decision-theoretic limits on earthquake prediction. *Geophys. J. Int.*, 131(3):526–529, December 1997.
- [Max98] Roy A. Maxion. Measuring intrusion-detection systems. Presented to The First Intl. Workshop on Recent Advances in Intrusion Detection (RAID-98), Lovain-la-Neuve, Belgium, *No printed proceedings*, 14–16 September 1998.
- [Nyg94] Else Nygren. Moderna tider: teknikutveckling inom medicinsk service. Technical report, Vårdförbundet SHSTF 42, Stockholm, Sweden, 1994. ISBN91-7043-021-7, ISSN 0349-1757, In Swedish.
- [Pie48] G. McGuire Pierce. Destruction by demolition, incendiaries and sabotage. Field training manual, Fleet Marine Force, US Marine Corps, 1943–1948. Reprinted: Paladin Press, PO 1307, Boulder CO, USA.
- [Ras86] Jens Rasmussen. *Information processing and human-machine interaction, An approach to cognitive engineering*. Elsevier Science Publishing Co., Inc., 52 Vanderbilt Avenue, New York, New York 10017, first edition, 1986.
- [RN95] Stuart J. Russel and Peter Norvig. *Artificial Intelligence—A Modern Approach*, chapter 14, pages 426–435. Prentice Hall Series in Artificial Intelligence. Prentice Hall International, Inc., London, UK, first edition, 1995. Exercise 14.3.
- [SSHW88] Michael M. Sebring, Eric Shellhouse, Mary E. Hanna, and R. Alan Whitehurst. Expert systems in intrusion detection: A case study. In *Proceedings of the 11th National Computer Security Conference*, pages 74–81, Baltimore, Maryland, 17–20 October 1988. NIST.
- [Sto95] Clifford Stoll. *The Cuckoo’s Egg: Tracking a Spy Through the Maze of Computer Espionage*. Pocket Books, July 1995.
- [Tre68] Harry L. Van Trees. *Detection, Estimation, and Modulation Theory, Part I, Detection, Estimation, and Linear Modulation Theory*. John Wiley and Sons, Inc., 1968.
- [VL89] H S Vaccaro and G E Liepins. Detection of anomalous computer session activity. In *Proceedings of the 1989 IEEE Symposium on Security and Privacy*, pages 280–289, Oakland, California, 1–3 May 1989.
- [WFP99] Christina Warrender, Stephanie Forrest, and Barak Perlmutter. Detecting intrusions using system calls: Alternative data models. In *IEEE Symposium on Security and Privacy*, pages 133–145, Berkeley, California, May 1999.

[Wic92] Christopher Wickens. *Engineering psychology and human performance*. HarperCollins Publishers Inc., 10 East 53rd Street, New York, NY 10022, second edition, 1992.

Appendix A Venn Diagram of the Base-Rate Fallacy Example

The Venn diagram in figure 4 depicts the situation in the medical diagnostic example of the base-rate fallacy given earlier.

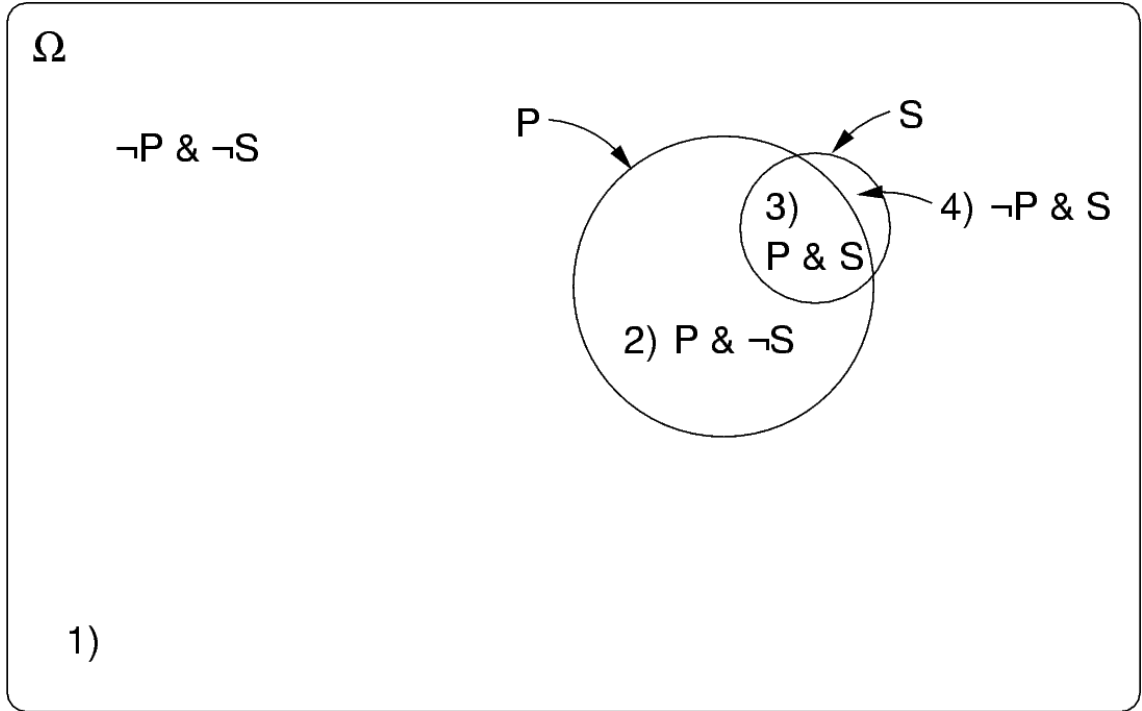


Figure 4: Venn diagram of medical diagnostic example

Although for reasons of clarity the Venn diagram is not to scale it clearly demonstrates the basis of the base-rate fallacy, i.e. that the population in the outcome S is much smaller than that in $\neg S$ and hence, even though $P(P|S) = 99\%$ and $P(\neg P|\neg S) = 99\%$, the relative sizes of the missing 1% in each case—areas 2) and 4) in the diagram—are very different.

Thus when we compare the relative sizes of the four numbered areas in the diagram, and interpret them as probability measures, we can state the desired probability, $P(S|P)$ —i.e. “What is the probability that we are in area 3) given that we are inside the P -area?” It may be seen that, area 3) is small relative to the entire P -area, and hence, the fact that the test is positive does not say much, in absolute terms, about our state of health.

