

FRAUD DETECTION THEORY DOC

MICHAEL MCLEOD

1. THE PROBLEM

We want to be able to estimate the probability of a given transaction being a fraud. Each transaction is tested against a set of rules, any number of which it may trigger. The goal is to estimate the fraud probability given the rules which have been triggered, and some aggregated historic data from past transactions. This historic data contains the total number of transactions, the total number of fraudulent transactions, and for each rule in the rule set the number of times the rule was triggered by fraudulent transactions and the number of times the rule was triggered by genuine transactions.

2. A SIMPLE BAYESIAN APPROACH

Let f be the proposition that the current transaction is a fraud,

$$(1) \quad f = \begin{cases} 1, & \text{transaction is a fraud} \\ 0, & \text{transaction is genuine.} \end{cases}$$

We want to calculate

$$(2) \quad P(f|S, H),$$

where S is the set of rules triggered and H is the historic data. Let

$$(3) \quad R = \{r_0, \dots, r_n\}$$

be the set of rules available in the system; then

$$(4) \quad S \in 2^R.$$

Using Bayes' theorem we then want to calculate the posterior probability of a fraud

$$(5) \quad P(f|S, H) = \frac{P(S|f, H)P(f|H)}{P(S|H)}.$$

We shall tackle each term in turn.

2.1. The Prior.

The prior term is

$$(6) \quad P(f|H),$$

the probability of the transaction being a fraud given the historic data, but in lieu of rule information (or, rather, integrated over it). Here we will make an assumption that there exists some overall fraud rate F with which frauds occur, such as $F \in [0, 1]$ represents the fraction of fraudulent transactions. Then

$$(7) \quad P(f|F) = F.$$

$$(8) \quad P(f|H) = \int_0^1 P(f|F)P(F|H) dF = \int_0^1 F P(F|H) dF$$

Since F is unknown to us we shall assign it a flat prior, $P(F) = 1$, and therefore

$$(9) \quad P(F|H) = \frac{P(H|F)P(F)}{P(H)} = \frac{P(H|F)}{\int_0^1 P(H|F)dF}$$

Assuming that transactions are all independent then

$$(10) \quad P(H|F) = \binom{n_f + n_g}{n_f} F^{n_f} (1 - F)^{n_g} dF,$$

where n_f is the total number of fraudulent transactions, and n_g is the total number of genuine transactions. The total number of transaction overall is therefore $n_f + n_g$.

Then we have for the denominator

$$(11) \quad \int_0^1 P(H|F)P(F) = \binom{n_f + n_g}{n_f} \int_0^1 F^{n_f} (1 - F)^{n_g} dF = \binom{n_f + n_g}{n_f} B(n_f + 1, n_g + 1),$$

where $B(x, y)$ is the Beta function. Expanding the Beta function in terms of the Γ -function, and substituting the Γ -function for factorial we have:

$$(12) \quad \int_0^1 P(H|F)P(F) = \frac{(n_f + n_g)!}{n_f! n_g!} \frac{n_f! n_g!}{(n_f + n_g + 1)!} = \frac{1}{n_f + n_g + 1}$$

Substituting equation 12 and equation 9 into equation 8 leaves us with

$$(13) \quad P(f|H) = (n_f + n_g + 1) \binom{n_f + n_g}{n_f} \int_0^1 F^{n_f+1} (1 - F)^{n_g} dF.$$

This is integrated similarly to yield

$$(14) \quad P(f|H) = (n_f + n_g + 1) \frac{(n_f + n_g)!}{n_f! n_g!} \frac{(n_f + 1)! n_g!}{(n_f + n_g + 2)!} = \frac{(n_f + 1)}{(n_f + n_g + 2)}.$$

The probability of a genuine transaction \bar{f} is calculated in the same way to yield

$$(15) \quad P(\bar{f}|H) = \frac{(n_g + 1)}{(n_f + n_g + 2)},$$

which satisfies the equation $P(f|H) + P(\bar{f}|H) = 1$ as necessary.

2.2. The Likelihood.

This term is

$$(16) \quad P(S|f, H),$$

the probability of the given rule set S (and only S) being triggered by the transaction, given that the transaction is fraudulent and the historic data.

Assuming that each rule is triggered independently, the

$$(17) \quad P(S|f, H) = \prod_{r \in S} P(r|f, H) \times \prod_{r \in R-S} (1 - P(r|f, H)),$$

i.e. the product of the probability of being triggered for each triggered rule (in S) and the probability of not being triggered for each rule not triggered (in R but not in S).

Again we do not know the probability of a fraudulent transaction triggering a given rule, so we take the same approach as we did above and assume that such a rate exists, and then marginalise over it.

$$(18) \quad P(r|f, H) = \frac{n_{r,f} + 1}{n_f},$$

where $n_{r,f}$ is the number times rule r was triggered by a fraudulent transaction.

2.3. The Evidence.

This term is

$$(19) \quad P(S|H) = P(S|f, H)P(f|H) + P(S|\bar{f}, H)P(\bar{f}|H)$$

The term $P(S|\bar{f}, H)$ is calculated in the same way as the likelihood, but for the probability of rules being triggered by a genuine transaction,

$$(20) \quad P(S|\bar{f}, H) = \prod_{r \in S} P(r|\bar{f}, H) \times \prod_{r \in R-S} (1 - P(r|\bar{f}, H)),$$

$$(21) \quad P(r|\bar{f}, H) = \frac{n_{g,r} + 1}{n_g + 2}.$$

The posterior may be calculated from a combination of the terms above.