

Consolidated Review of

Understanding the Domain Registration Behavior of Spammers

1. Strengths

The paper systematically studies the domain registration lifecycle of spam domains and attempts to identify distinguishing features that can be used to classify them as malicious at registration time itself. While the results are not conclusive the paper does a great job explaining the methods used and the behaviors observed during the registration process. They categorize domains as brand-new, re-registration, drop catch and retread and explain the lifecycle. They investigate the properties of spam domains along freshness, persistence and black list efficacy. They categorize the DNS servers used for the domains using metrics for toxicity, duplication and association with registrars. They develop a compound Poisson model to allow ease in classification of spam domains during registration time.

Extensive use of statistics from .COM domain (at 5 minute increments over several months) combined with various spam detection techniques.

2. Weaknesses

Although the stated utility is to come up with methods that could potentially detect spam activity during registration process, I am not convinced that the observations in the paper can indeed be integrated in registrants and registrars in a way that will help mitigate these domains.

The insight of attempting to mitigate spam activity by identifying it early in the lifecycle, i.e.; at registration, is a technique already robustly employed by Spamhaus DBL. The paper mentions how Spamhaus proactively black lists domains based on registration information. They authors however do not provide any additional details or compare their methods/results with Spamhaus. The paper does not discuss what are the primary reasons why a small set of registrants are favored by malicious domain registrations. However, I still do believe that the systematic analysis presented in the paper provides some valuable insights to the community.

Lots of sloppy analysis.

The methodology itself is not particularly interesting. But this is a subjective issue and not a reason to reject the paper.

3. Comments

The paper provides new and interesting observations about how spammers register domains.

Overall I like the statistical work but hate its interpretation.

I have the following concerns about the paper:

- ❖ Are the observations really useful? Perhaps one way to check this is to "replay" the registration process, and implement a spam-detection mechanism that is based on your methods. Using this you can study (or approximate) the detection rate, false-positives etc. You provide some insights about this in section 6, but I think this could be extended to encompass all your techniques.

- ❖ Related to the first comment - lets assume that your methods actually work and do get implemented. Spammers can easily overcome them (slow down the registration rate, don't reuse recent URLs, not register in bulks, etc.). Does this problem even have a solution? Perhaps an easier, long-term solution will be to change policies, e.g. require strong identification.
- ❖ Section 5.1, saying actions by a small number of registrars could really help is not satisfying. What could these registrars do differently? Have you read their registration procedures (there are only about a dozen of these folks) and examined how they differ from other registrars? Might you see if, combined with your data, those procedures point to things the registrars should do differently. Perhaps the speed at which domains are used varies by registrar (imagine that 6.1 is right and some registrations are done with fraudulent credit -- that gives a small window for the domains to exist -- are there registrars who put names on-line before credit card fraud could be detected???)
- ❖ The tracking rogue domains by which servers they use is tempting, but how might spammers work around it? Similarly, right now spammers use bulk registration, but given that the majority of domains are used 1 mo + after registration, presumably tracking bulk registrations will simply lead to spammers registering at more regular rates all the time? The paper would be much improved by doing these second level analyses.
- ❖ You study the .com domain. How many spammers actually use .com as opposed to other domains? Eyeballing over the top-50 URIs in uribl.com shows usage of .us, .biz, .info and .net. You should provide the percentage of spam domains that your study covers.
- ❖ Table 3 doesn't contribute more than the text describing it. I'd delete Table 3.
- ❖ Figure 5 shows ABSystems and it has unique distinct characteristic. However it is not discussed in the paper. It would be interesting to understand what causes the observed results
- ❖ 4.1 freshness. As I read the stats, they say that 18%-25% of spam domains are registered in the current month, rising to 32-42% that are 2 months or less old. How does this say "rely heavily on newly registered domains"? My definition of newly registered would be one month -- not 2 or 3 months. And I note the table only runs to 3 months -- at what month do we reach saturation (say 90% of spam domains). Seems to me the stats here suggest that spammers actually camp on domain names for extended periods before using them. The blacklist efficacy in 4.3 supports this notion "domains often remain idle after registration for a number of days...."
- ❖ Note the longevity of domains also blows up the suggestion in 6.1 that bulk registration is done by "fraudulent cards, which will soon be detected." If over 70% of spam domains are > 1 month old, then they were NOT acquired with fraudulent credit. This analysis needs rethinking....

4. Summary from PC Discussion

All the reviewers liked the extent of the dataset and the timeliness of the topic. The discussion mainly focused on the discrepancy in the text versus what is in the plots. We decided to shepherd the paper to fix these concerns.

5. Authors' Response

We regret that the paper as submitted suffered from significant writing issues. These surely complicated the task of the reviewers in assessing the work. One particularly noteworthy point in this regard concerns the lifetime of spammer domains, which was assessed in two different dimensions (time from registration to first use, versus age of domains **given** they are already on blacklists) that the text failed to adequately distinguish, leading to the confusion expressed in the last two bullet points in the review summary. To this end, we believe the criticisms of "sloppy analysis" really in fact instead reflect "sloppy writing" (though the latter makes it difficult to not conclude the former). We stand by the careful nature of our analysis, but clearly the writing of the submission made it difficult to discern. To remedy this, we have undertaken a full rewrite of the text in an attempt to make it much more clear.

Regarding the question of whether the observations are really useful, we agree that this is of prime interest. However, we believe that answering this question requires extensive additional work, and indeed that is where our efforts have focused for several months now.

The issue of the inevitable "arms race" that arises between defenders and adversaries is endemic to the general problem space

of detecting malicious activity conducted by actors motivated to adapt to new defenses. Given that much of spam is profit-driven, our basic goal is to fundamentally increase the cost to attackers for conducting such activities. That some of the features we identify plausibly reflect economies-of-scale employed by attackers (bulk registrations, reliance upon a relatively small number of registrars) gives some promise that countermeasures that drive attackers away from these conveniences will indeed add friction to the spam business enterprise. Even so, we certainly don't argue that such techniques will by themselves defeat spam.

One reviewer states that the analysis boils down to "a technique already robustly employed by Spamhaus DBL" and criticizes the paper on that basis. Here we would call out: just how do we know that the technique is robustly employed (i.e., with strong efficacy)? How do we even know that Spamhaus employs it at all, other than from the evidence presented in our paper? The comment highlights the key difference between a technique that industry apparently uses - but that has not seen impartial assessment - versus an attempt at a scientifically methodical study of such techniques. Perhaps the reviewer knows directly about the Spamhaus DBL's operation and its accuracy, but to our knowledge that is not anything available in the open literature.

Finally, a reviewer asks whether .com has relevance as opposed to other domains. While our primary data source has requested that we not explicitly quantify the prevalence of different TLDs in appearing in spam messages, we can state that .com definitely continues to play a role in this regard.