

# Detecção de Spammers na Rede de Origem

Pedro Henrique B. Las-Casas<sup>1</sup>, Humberto T. Marques-Neto<sup>1</sup>

<sup>1</sup> Departamento de Ciência da Computação  
Pontifícia Universidade Católica de Minas Gerais (PUC Minas)  
30.535-901 - Belo Horizonte - MG - Brasil

pedro.casas@sga.pucminas.br, humberto@pucminas.br

**Abstract.** *The volume of unsolicited messages (spam) sent over the Internet represents more than 85% of all e-mails. Even with the evolution of the filtering techniques such as the analysis of the message content and the blocking of IPs, network resources are wasted given that such a filtering is usually performed at the e-mail destination server. This paper proposes a method for detecting spammers in the origin network using a supervised classification technique with metrics which do not require inspection of message contents. Results show that the adopted method is efficient, being able to correctly identify most spammers still in the origin network, saving network resources.*

**Resumo.** *A quantidade de mensagens não-solicitadas (spams) enviadas na Internet representa mais de 85% de todos os e-mails. Mesmo com a evolução de técnicas de filtragem como a análise do conteúdo de mensagens e o bloqueio de IPs, recursos da rede são desperdiçados, uma vez que essa filtragem é realizada normalmente no servidor de destino dos e-mails. Este trabalho propõe um método para detecção de spammers na rede de origem utilizando uma técnica de classificação supervisionada composta por métricas que não requerem a inspeção do conteúdo das mensagens enviadas. Os resultados mostram que o método utilizado é eficaz, sendo capaz de identificar a maioria dos spammers ainda em sua rede de origem, preservando assim, os recursos da rede.*

## 1. Introdução

Relatórios recentes [MessageLabs 2010] indicam que cerca de 85% das mensagens eletrônicas (e-mails) que circulam pela Internet são mensagens indesejadas, muitas vezes caracterizadas como spam. Para reduzir o impacto dessas mensagens indesejadas sobre o serviço de correio eletrônico, provedores utilizam filtros de e-mails que as classificam, descartam ou colocam em quarentena, evitando que elas sobrecarreguem as caixas postais dos destinatários. Entretanto, esses filtros não evitam o desperdício de recursos da rede, pois as mensagens recebidas geram tráfego nos links e consomem CPU para serem encaminhadas e processadas.

Uma possível forma de evitar esse desperdício de recursos causado pelos *spams* seria complementar a filtragem no servidor receptor com o uso de técnicas de filtragem prévia, capazes de evitar o envio do *spam* e, com isso, o desperdício de recursos associado. Tais técnicas podem ser aplicadas, por exemplo, em provedores de acesso à Internet de banda larga que, através da análise do tráfego SMTP (*Simple Mail Transfer Protocol*), poderiam detectar a ação de possíveis *spammers*, bloqueando o tráfego na sua origem.

Este artigo descreve a proposta de um novo método para detecção de spammers na rede de origem chamado SpaDeS - *Spammer detection at the Source* - que é baseado em um algoritmo de classificação supervisionada e explora apenas métricas que não requerem a inspeção do conteúdo das mensagens. Para o desenvolvimento da proposta deste método, utilizou-se um estudo prévio do tráfego da rede de origem de um provedor de Internet de banda larga [Castilho et al. 2010]. Este estudo apresenta a caracterização do tráfego SMTP presente na rede de origem, e define métricas capazes de serem utilizadas para distinguir spammers de usuários legítimos. O aluno de Iniciação Científica (IC) Pedro Las Casas participou ativamente deste trabalho, que foi desenvolvido no contexto do projeto ReBu: Sistemas de Redes Robustos - Modelos e Ferramentas. A partir deste trabalho, o aluno de IC desenvolveu a proposta do método de detecção de spammers na rede de origem. Maiores detalhes do método aqui apresentado se encontram em [Las-Casas et al. 2011], artigo científico publicado no XXIX Simpósio Brasileiro de Redes de Computadores e Sistemas Distribuídos (SBRC 2011), sendo o aluno de IC o seu autor principal.

O método proposto foi aplicado e validado utilizando dois conjuntos de dados reais contendo informações agregadas e anonimizadas de transações SMTP de usuários de um provedor brasileiro de Internet de banda larga residencial coletadas em 2009 e 2010. Os resultados apresentados mostram que a utilização do SpaDeS é capaz de diferenciar *spammers* de usuários legítimos ainda na rede de origem, sem inspecionar o conteúdo de suas mensagens. Com o uso da técnica de classificação supervisionada, validada através da comparação com uma base de dados real de denúncias de *spam* e da inspeção de uma amostra dos usuários classificados, estima-se que cerca de 98% dos usuários legítimos e 94% dos *spammers* foram classificados corretamente. Ou seja, as taxas de falsos positivos e de falsos negativos foram, respectivamente, de 2% e 6%. O estudo mostra também que classes de usuários legítimos representaram cerca de 83% dos usuários e realizaram cerca de 1,6% do total de transações SMTP observadas nos dados de 2010. Enquanto isso, os usuários classificados como *spammers*, (cerca de 17%) originaram mais de 98% de todas as transações SMTP no período observado.

O restante deste artigo está organizado da seguinte forma: a Seção 2 apresenta o método de detecção de *spammers* proposto e a Seção 3 discute os resultados mais relevantes do trabalho. Finalmente, a Seção 4 apresenta conclusões e sugere trabalhos futuros.

## **2. SpaDeS: Detector de Spammers na Rede de Origem**

O método proposto para detecção de *spammers* na rede de origem, denominado SpaDeS (*Spammer Detection at the Source*), tem como principal componente um algoritmo de classificação supervisionada, que “aprende” um modelo de classificação de usuários a partir de um conjunto de exemplos (usuários) previamente classificados (conjunto de treino). O classificador recebe como entrada o número de classes distintas  $C$  e exemplos de usuários de cada uma. Após a fase de aprendizado, o modelo derivado pode então ser aplicado para classificar novos usuários (conjunto de teste) nas classes pré-definidas. A Seção 2.1 apresenta as classes de usuários consideradas assim como o modelo de representação dos mesmos. O algoritmo de classificação utilizado é apresentado na Seção 2.2, enquanto a Seção 2.3 discute como obter o conjunto de treino.

## 2.1. Modelo de Representação de Usuários e suas Classes

Cada usuário é representado por um vetor de  $N$  atributos que conjuntamente descrevem seu comportamento quanto ao uso do protocolo SMTP. Para detectar *spammers* na rede de origem com eficiência, foram utilizados  $N=6$  atributos, que são métricas que não envolvem processamento do corpo da mensagem. As métricas são: número de transações SMTP realizadas, número de remetentes distintos, número de servidores SMTP distintos acessados, tamanho médio das transações SMTP, distância geodésica média entre origem e destino e tempo médio entre transações consecutivas (aqui referenciado como IATs, *inter-arrival times*).

A escolha das métricas foi inspirada nos resultados do trabalho [Castilho et al. 2010], em que o aluno de IC participou implementando e executando algoritmos, como também analisando resultados. Utilizando o algoritmo de agrupamento *X-means*, demonstrou-se que essas características podem ser utilizadas para distinguir 4 classes de comportamento, sendo que duas refletem padrões de usuários legítimos enquanto as outras refletem padrões abusivos, potencialmente de *spammers*. Por exemplo, o número de transações por usuário é útil para distinguir usuários que fazem pouco uso de SMTP daqueles que o utilizam com grande intensidade. Mais ainda, enquanto o uso de poucos servidores de SMTP é o esperado para usuários legítimos, o acesso a um número muito grande pode indicar a operação de *open proxies* ou de *open mail relays* sendo explorados para o envio de spam por usuários maliciosos ou *bots*. O uso da distância geodésica como métrica se baseia na hipótese de que conexões SMTP de *spammers* tendem a ocorrer entre endereços IPs mais distantes que as conexões de usuários legítimos, já que *spammers* tendem a ocultar sua presença usando máquinas em outros países [Guerra et al. 2008].

## 2.2. Algoritmo de Classificação Supervisionada

O algoritmo de classificação utilizado neste trabalho é o *Lazy Associative Classifier* (LAC) [Veloso et al. 2006], que tem ótima escalabilidade, com complexidade de tempo polinomial. Diferentemente de outros classificadores, o LAC fornece uma estimativa da confiança na predição feita em cada caso. Essa confiança, que pode ser interpretada como uma probabilidade de acerto da classificação, será explorada na geração do conjunto de treino (Seção 2.3). O LAC explora o fato de que, frequentemente, há fortes associações entre os valores dos atributos e as classes. Tais associações estão geralmente implícitas no conjunto de treino e, quando descobertas, revelam aspectos que podem ser utilizados para prever as classes dos usuários.

## 2.3. Coleção de Treino

O funcionamento de qualquer método de classificação supervisionada depende primariamente de um conjunto de treino contendo usuários pré-classificados. A obtenção desse conjunto para a classificação de usuários em *spammers* e legítimos é um grande desafio, uma vez que tais dados tipicamente não estão disponíveis publicamente. Um fator complicador é que almeja-se detectar *spammers* ainda na rede de origem. Logo, faz-se necessário um conjunto de treino coletado naquele ponto do sistema. Caso contrário, os padrões levantados poderiam não generalizar para o conjunto de teste, resultando em um desempenho pobre do classificador. São propostas duas estratégias para geração da coleção de treino, uma baseada em informação externa ao algoritmo de aprendizado, enquanto a segunda utiliza a informação de confiança provida pelo LAC.

A primeira estratégia parte das 4 classes de usuários identificadas em [Castilho et al. 2010] e acrescidas de informações sobre usuários apontados como *spammers* por mecanismos de relato de abusos do sistema de correio. Para cada uma das duas classes representando usuários legítimos identificadas, foram selecionados os  $M$  usuários mais próximos do centróide de cada classe, de forma a obter bons representantes de cada uma. As duas classes de usuários abusivos (possíveis *spammers*) identificadas naquele trabalho apresentam uma variabilidade maior de comportamentos. Por esse motivo, optou-se por não usar o mesmo mecanismo de seleção, mas basearmos a escolha em uma informação confiável externa. Para isso, utilizou-se a identificação de usuários cujas máquinas foram apontadas como origem de *spam* por relatos oriundos de outros provedores. Tais relatos, enviados para o endereço *abuse* do provedor que forneceu os dados utilizados nesse trabalho, são gerados por provedores tanto a partir de reclamações de seus usuários (como o recurso “*Report spam*” do Gmail) ou por mecanismos automáticos, como listas de bloqueio. Como será visto, um número razoável de usuários denunciados estavam presentes naquelas duas classes, sendo portanto bons representantes das mesmas.

A segunda estratégia parte do pressuposto que o SpaDeS deve ser aplicado continuamente, em diferentes conjuntos de teste (p.ex: dados referentes a diferentes semanas, meses ou anos). Logo, propôs-se retreinar o LAC a partir do resultado da sua execução anterior, explorando as confianças reportadas naquela classificação. Ou seja, considerando sucessivos conjuntos de teste  $t_1, t_2 \dots t_n$ , seleciona-se como treino do LAC para a classificação do teste  $t_i$ , os usuários do conjunto  $t_{i-1}$  que foram classificados com uma confiança superior a um certo limiar. Para a classificação de  $t_1$ , um conjunto de treino inicial é necessário, podendo ser obtido pela primeira estratégia. O algoritmo 1 apresenta a estratégia utilizada. Ele garante que pelo menos  $\alpha\%$  dos usuários de cada classe sejam selecionados, mantendo uma confiança mínima uniforme entre todas as classes.

---

Algoritmo 1: Dados os usuários classificados pelo LAC na iteração anterior, faça:

1. Ordene os usuários de cada classe em ordem decrescente de confiança;
  2. Selecione  $\alpha\%$  dos usuários de cada classe, ordenados anteriormente;
  3. Seja  $c_i^{min}$  a menor confiança dos usuários selecionados da classe  $i$  ( $i = 1..4$ );
  4. Seja  $c = \min(c_1^{min}, c_2^{min}, c_3^{min}, c_4^{min})$ ;
  5. Selecione para o conjunto de treino todos os usuários que possuem confiança  $\geq c$ , mantendo, para cada um, a classe definida pelo LAC na iteração anterior.
- 

Neste trabalho utilizou-se a primeira estratégia, baseada no algoritmo de agrupamento, apenas na iteração inicial e continuou-se o treinamento a partir do resultado da classificação anterior para as iterações seguintes. Esse enfoque se baseia no fato de que, considerando que um número suficiente de bons exemplos de treinos sejam fornecidos, técnicas de classificação supervisionadas (p.ex., LAC) tendem a ser superiores às técnicas não supervisionadas (p.ex., algoritmos de agrupamento) [Veloso et al. 2006]. Note que a coleção de treino pode ser estendida e/ou refinada para incluir exemplos pré-classificados por outros meios, potencialmente mais confiáveis, se tais exemplos estiverem disponíveis. Por exemplo, assim como feito na iteração inicial, havendo conhecimento sobre usuários locais denunciados como *spammers*, os mesmos poderiam ser incluídos no treino das iterações seguintes. Como ficará claro na Seção 3, optou-se por não incluir tais usuários no conjunto de treino nas iterações seguintes para que os mesmos pudessem ser utilizados

para validação do método proposto.

A estratégia proposta é completamente automatizada e não exige esforço manual de classificação. Note que a natureza iterativa do processo, que utiliza resultados da iteração anterior como treino da próxima iteração, pode afetar a classificação ao longo do tempo. Entretanto, nos experimentos realizados, observou-se que os padrões de cada classe de usuários se mantêm estáveis em duas bases de dados incluindo tráfego em 2009 e 2010. Mais ainda, a classificação dos usuários da base de 2010, seguindo a abordagem iterativa descrita, apresentou excelente efetividade (Seção 3). De qualquer maneira, considera-se que, para melhor refletir os padrões de comportamento dos usuários, que podem evoluir com o tempo, e também para interromper uma possível propagação de erros, seja necessário, periodicamente, a aplicação de um conjunto de treino obtido por métodos externos (como a primeira estratégia proposta), reiniciando um novo ciclo de iterações.

### 3. Avaliação e Resultados

Esta seção apresenta uma avaliação do SpaDeS, descrevendo as bases de dados utilizadas e o procedimento de avaliação, e discutindo os principais resultados obtidos.

#### 3.1. Bases de Dados

Este trabalho utiliza 4 bases de dados diferentes, sendo que duas delas refletem o tráfego SMTP de um provedor de Internet de banda larga e duas contêm listas de usuários daquele provedor que foram denunciados como *spammers* através do endereço *abuse* daquele provedor durante o período considerado.

Cada base de dados de tráfego contém um log de tráfego e um log do serviço DHCP do provedor, ambos cobrindo um mesmo período. Os logs de tráfego são formados por *transações* que representam uma conexão TCP ou um fluxo de dados UDP, contendo informações como endereços IP de origem e de destino, serviço/protocolo utilizado, data/hora inicial, duração e volume de bytes enviados e recebidos. Os logs do serviço DHCP permitem associar transações e usuários através do mapeamento dos endereços físicos de suas máquinas (*MAC addresses*) para os endereços IP fornecidos pelo provedor, com base nas informações de data e hora presentes nos dois logs. Vale ressaltar que os dados dos usuários foram anonimizados, por questões de segurança e privacidade.

As bases de dados de tráfego cobrem os períodos de 01 a 28 de março de 2009 e 12 de junho a 09 de julho de 2010. A base de 2009 contém 40,6 milhões de transações associadas a 44,2 mil usuários. Já a de 2010 contém 45,6 milhões de transações associadas a 48 mil usuários. Foram filtradas as transações SMTP, restando 6,3 milhões de transações SMTP associadas a 5.479 usuários na base de 2009, e 5 milhões de transações SMTP associadas a 5.389 usuários na base de 2010.

As duas outras bases de dados contêm denúncias recebidas pelo endereço *abuse* do provedor durante os períodos das bases de tráfego, identificando certos usuários como *spammers*. Os e-mails de denúncia informam o endereço IP de origem do *spam* e a data/hora do seu recebimento e estão no formato ARF (*Abuse Reporting Format*), utilizado para mensagens desse tipo. Foi desenvolvida uma ferramenta de extração para processar essas mensagens e realizar a junção das mesmas com as transações SMTP, possibilitando a identificação de usuários denunciados. Dessa forma, foram identificados 67 e 93 *spammers* nas bases de 2009 e 2010, respectivamente. Para todos esses usuários, os

endereços IPs e as datas/horas listados nas denúncias coincidiram com dados de transações realizadas, listadas nas bases de tráfego utilizadas.

### 3.2. Procedimento de Avaliação

A avaliação consistiu de dois experimentos de classificação, um com as bases de 2009 e outro com as bases de 2010. Para a classificação da base de 2009, foram selecionados  $M=30$  usuários mais próximos do centróide de cada uma das classes de usuários legítimos (Seção 2.1). Além disso, dos 67 usuários denunciados como *spammers* identificados na base de 2009, 40 são da classe 3 e 27 da classe 4. Assim, essas duas classes são formadas principalmente por *spammers* e são o principal alvo do método de detecção. O conjunto de treino foi então composto por esses 127 usuários, e para teste foram utilizados todos os usuários da base de 2009 *que não estavam no conjunto treino*.

Para o segundo experimento de classificação, realizado sobre a base de 2010, foram utilizados como treino usuários selecionados pelo algoritmo 1 considerando o resultado da classificação da base de 2009 (Seção 2.3). Foi utilizado  $\alpha = 20\%$ , o que resultou em uma confiança mínima para todos os usuários selecionados de 64%. No total, foram selecionados 787, 605, 621 e 52 usuários das classes 1, 2, 3 e 4, respectivamente.

As seções seguintes apresentam os principais resultados dos dois experimentos. Para o primeiro experimento, não foi possível quantificar a efetividade da classificação, uma vez que os dados do *abuse*, que poderiam servir como base de comparação, foram usados no treino. Optou-se nesse caso por analisar os padrões de comportamento dos usuários classificados em cada classe. A principal validação quantitativa é feita sobre os resultados do segundo experimento: além da avaliação dos padrões identificados, utilizou-se a lista de usuários denunciados como *spammers* em 2010 e foi feita uma inspeção manual de uma amostra dos demais usuários para estimar a efetividade da classificação.

### 3.3. Classificação da Base de Dados de 2009

O conjunto de treino é composto por usuários selecionados, que representam as quatro classes identificadas anteriormente [Castilho et al. 2010]. Para cada métrica analisou-se a média e o coeficiente de variação CV (razão entre desvio padrão e a média), computados para os usuários de cada classe. As classes 1 e 2 apresentam comportamentos razoavelmente bem uniformes, principalmente com relação aos números de transações, número de remetentes distintos, número de servidores contactados e IAT das transações. A classe 1 compreende usuários que fazem muito baixo uso do correio eletrônico, com apenas uma transação no período coberto pela base. Já a classe 2 compreende usuários com um nível de atividade um pouco mais alto, realizando tipicamente uma transação SMTP a cada 3 dias (IAT médio igual a 70,43 horas). Para ambas as classes, os números de remetentes identificados (1 e 1,70, em média) e de servidores SMTP distintos acessados (1 e 1,53, em média) são baixos, demonstrando um comportamento esperado de usuários legítimos.

Já as classes 3 e 4, representativas de *spammers* e definidas pelos usuários denunciados, revelam padrões bem distintos. Embora ambas apresentem números de transações, remetentes e servidores SMTP superiores aos das classes 1 e 2, a classe 4 revela um padrão muito mais abusivo, com cada usuário enviando 37.841,48 transações SMTP (uma a cada 36 segundos), utilizando 21.660,22 remetentes distintos e acessando

**Tabela 1. Classificação dos Usuários da Base de 2009.**

	Classe 1	Classe 2	Classe 3	Classe 4
Número de usuários	1.422	2.938	927	65
	<b>Média (CV)</b>	<b>Média (CV)</b>	<b>Média (CV)</b>	<b>Média (CV)</b>
Número de transações SMTP	1 (0)	25,71 (9,61)	2.697,54 (1,53)	39.108,21 (1,05)
Número de remetentes distintos	1 (0)	4,14 (1,93)	1.202,50 (1,39)	16.032,23 (0,98)
Número de servidores SMTP distintos	1 (0)	5,33 (2,32)	1.261,37 (1,17)	11.766,04 (0,60)
Tamanho das transações SMTP (KB)	646,81 (4,34)	611,40 (7,51)	9,64 (12,83)	1,67 (0,64)
Distância geodésica entre os IPs (km)	2.657 (1,32)	3.133 (1,1)	8.081 (0,22)	8.352 (0,16)
IAT das transações SMTP (h)	672 (0)	55,34 (1,02)	0,57 (2,45)	0,01 (0,55)

14.704,89 servidores SMTP distintos, em média. De fato, essas classes revelam dois tipos de *spammers* distintos: um envia o maior número de mensagens possível (classe 4) e o outro (classe 3) envia mensagens utilizando, possivelmente, um controle de fluxo com o objetivo de disfarçar sua presença [John et al. 2009]. É interessante notar também os valores médios de IAT muito baixos e as distâncias geodésicas muito mais altas, para ambas as classes, padrões esperados para *spammers*. Por fim, vale também ressaltar os tamanhos de transações muito menores que os das classes 1 e 2, em consistência com estudos anteriores que demonstraram que *spams* tendem a ser menores que mensagens legítimas [Gomes et al. 2007].

Os resultados da classificação são apresentados na Tabela 1, que sumariza as principais características dos usuários em cada classe. A classe 1 é formada por 1.422 usuários, que representam 26,56% do total de usuários na base de 2009, mas são responsáveis por menos de 1% do total de transações SMTP enviadas. Em consistência com o conjunto de treino, esta classe se mostrou a mais uniforme: todos os usuários fizeram apenas uma transação SMTP, utilizando portanto apenas um remetente e acessando apenas um servidor SMTP no período de 28 dias.

A classe 2, composta por 2.938 usuários (54,9% do total), apresenta uma variação muito grande em seus usuários, principalmente com relação ao número de transações SMTP, que varia de 2 a 11.804. Apesar deste limite superior muito alto, notamos que apenas 2% dos usuários desta classe realizaram mais que 100 transações no período. Além disto, cerca de 90% dos usuários acessaram apenas 10 servidores SMTP distintos (5,33, em média) e utilizaram menos de 10 remetentes distintos (4,14, em média). Todas essas características indicam que grande parte dos usuários que compõem esta classe são legítimos. Alguns poucos usuários com um número de transações, número de remetentes e/ou número de servidores SMTP muito altos podem representar *spammers* erroneamente classificados como legítimos (falsos negativos), ou usuários com redes locais com diversos usuários.

A classe 3, formada por 927 usuários (17,32% do total), também apresenta grande variabilidade entre os usuários. O número de transações SMTP, por exemplo, varia entre 2 e 41.227, embora a média seja bastante alta (2.697,68 transações). De fato, mais de 50% dos usuários realizaram mais de 1.000 transações SMTP, utilizaram mais de 500 remetentes distintos e acessaram mais de 700 servidores SMTP. Além disso, o IAT médio apresentado por 90% dos usuários é menor que 10 minutos, o que indica que a cada 10 minutos o usuário realiza uma transação SMTP. Com base nestes dados, podemos supor que esses usuários estejam infectados por *malwares*, agindo como *bots* que utilizam controle de fluxo [John et al. 2009], ou seja, enviam *spams* com uma frequência relativamente

**Tabela 2. Classificação dos Usuários da Base de 2010.**

	Classe 1	Classe 2	Classe 3	Classe 4
Número de usuários	1.821	2.656	836	76
	<b>Média (CV)</b>	<b>Média (CV)</b>	<b>Média (CV)</b>	<b>Média (CV)</b>
Número de transações SMTP	2,78 (9,82)	27,99 (8,79)	2.892,91 (1,53)	34.018,17 (0,70)
Número de remetentes distintos	1 (0)	3,32 (1,37)	1.341,38 (1,73)	22.064,07 (0,73)
Número de servidores SMTP distintos	1 (0)	3,04 (1,81)	1.192,24 (1,52)	14.234,76 (0,50)
Tamanho das transações SMTP (KB)	891,55 (4,53)	826,45 (4,29)	22,64 (5,49)	2,23 (0,78)
Distância geodésica entre os IPs (km)	3.968 (0,91)	4.136 (0,82)	7.653 (0,29)	8.572 (0,05)
IAT das transações SMTP (h)	535,51 (0,49)	58,04 (0,96)	1,00 (2,05)	0,02 (0,50)

baixa (do ponto de vista de ferramentas automatizadas) para dificultar sua detecção por sistemas *anti-spam*. É importante ressaltar que essa maior variabilidade dos usuários nas classes 2 e 3 era esperada, uma vez que elas representam comportamentos de fronteira, que podem ser difíceis de serem distinguidos pelo classificador. Assim como discutido para a classe 2, conjectura-se que alguns usuários classificados como da classe 3 podem de fato serem falsos positivos.

Assim como a classe 1, a classe 4 apresenta pouca variabilidade, com CVs variando entre 0,15 e 1,05. Consistente com o conjunto de treino, foram identificados 65 novos usuários com um padrão muito abusivo de uso do SMTP, claramente relacionados à atividade de envio de *spam*. Embora representem apenas 1,21% dos usuários, eles são responsáveis por quase 40% de todas as transações SMTP realizadas.

### 3.4. Classificação da Base de Dados de 2010

A Tabela 2 mostra os resultados da classificação da base de dados de 2010, utilizando como conjunto de treino usuários selecionados a partir do resultado da classificação da base de 2009, conforme discutido na Seção 3.2. Note que, em termos gerais, os usuários de cada classe mantêm padrões de comportamento bastante semelhantes aos usuários da mesma classe na base de 2009. Por exemplo, as classes 1 e 2, que contabilizam 83% de todos usuários mas apenas 1,6% das transações SMTP realizadas, revelam padrões bastante consistentes com usuários legítimos: números pequenos de transações, remetentes e de servidores distintos e longos períodos de inatividade.

Já as classes 3 e 4, que representam 17% dos usuários e 98,4% das transações, mais uma vez demonstram padrões muito abusivos. Usuários da classe 4 fazem um uso muito mais intenso de tráfego SMTP que os da classe 3. Ainda assim, a classe 3 possui características muito pouco prováveis para usuários legítimos. Por exemplo, dificilmente um usuário legítimo realizaria 2.800 transações em um período de 28 dias (uma média de 100 transações por dia), utilizando 1.300 remetentes distintos e acessando 1.100 servidores SMTP distintos (pouco mais de 2,5 transações para cada servidor).

No geral, podemos concluir que, consistentemente nas duas bases analisadas (2009 e 2010), as classes de usuários legítimos: (i) realizam poucas transações SMTP; (ii) utilizam poucos remetentes distintos; (iii) acessam poucos servidores distintos; (iv) possuem alto intervalo de inatividade entre as transações; (v) possuem alta variabilidade do tamanho médio das transações SMTP, uma vez que usuários legítimos podem enviar tanto mensagens apenas de texto, quanto mensagens com anexos extensos, contendo vídeos e imagens; e (vi) realizam suas transações principalmente para servidores brasileiros ou servidores localizados nos Estados Unidos.

Em contraste, as classes de *spammers*: (i) realizam um número alto de transações SMTP; (ii) utilizam um número elevado de remetentes distintos; (iii) acessam vários servidores SMTP distintos; (iv) efetuam as transações com um período de inatividade muito baixo, sendo muitas vezes de apenas segundos; (v) possuem tamanho médio das transações SMTP baixos, uma vez que, normalmente, *spams* possuem apenas texto; e (vi) tendem a apresentar distância geodésica média maior que usuários legítimos.

Como discutido na Seção 3.2, além da avaliação dos padrões de comportamento detectados, também foi feita uma validação da classificação da base de 2010, utilizando a lista de usuários denunciados em 2010 (não utilizada como parte do treino) e inspecionando manualmente uma amostra aleatória de 5% dos usuários de cada classe (exceto classe 4). Essa taxa de amostragem garante, com confiança de 90%, um erro inferior a 12% nas estimativas. Como trabalho futuro, pretendemos realizar a inspeção manual com uma fração maior dos usuários.

**Tabela 3. Eficácia da Classificação: Estimativas de Taxa de Acerto, Falsos Positivos e Falsos Negativos (Base de Dados 2010).**

Classe Real	Classe Predita			
	1	2	3	4
1 (legítimo)	100%			
2 (legítimo)	0,7%	95,5%	3,8%	
3 ( <i>spammer</i> )		15,9%	84,1%	
4 ( <i>spammer</i> )				100%

Dos 93 usuários denunciados, 31 se encontram na classe 4 e 62 na classe 3. Ou seja, 40,78% e 7,41% dos usuários classificados nas classes 4 e 3, respectivamente, foram de fato denunciados como *spammers*. Foram inspecionados todos os 45 usuários restantes da classe 4, concluindo que todos apresentam um comportamento consistente com *spammers* bastante abusivos. Logo, todos os usuários classificados na classe 4 foram corretamente identificados como *spammers*. Além disto, verificou-se que 35 de 44 usuários selecionados da classe 3 para inspeção manual foram corretamente classificados, enquanto os 7 restantes apresentavam um comportamento aceitável para usuários legítimos da classe 2, sugerindo assim falsos positivos.

Quanto à classificação de usuários em usuários legítimos, apenas 1 de todos os usuários selecionados da classe 1 foi classificado erroneamente, apresentando um comportamento mais condizente com a classe 2. Note que, ainda assim, esse usuário foi classificado como legítimo. Já para a classe 2, 127 dos 132 usuários selecionados foram corretamente classificados, enquanto 5 tinham um padrão mais próximo de *spammers* da classe 3, sendo assim considerados falsos negativos.

Esses resultados, computados sobre os usuários amostrados das classes 1, 2 e 3 e sobre todos os usuários da classe 4, são sumarizados na Tabela 3. Cada linha representa uma classe atribuída aos usuários por inspeção ou pelo abuse (*classe real*), enquanto as colunas representam as classes assinaladas pelo SpaDeS (classes preditas). Os valores indicam as porcentagens das amostras de uma classe real que foram atribuídas à classe predita indicada. Dessa forma, a diagonal indica a taxa de acerto para cada classe e as demais posições indicam predições incorretas. Considerando a classificação de todos usuários nas super-classes “legítimos” (classes 1 e 2) e “*spammers*” (classes 3 e 4), o método SpaDeS apresentou uma excelente taxa de acerto, identificando corretamente 98%

e 94% dos usuários legítimos e *spammers*, respectivamente, apresentando assim taxas de falsos positivos e de falsos negativos de somente 2% e 6%, respectivamente.

Em suma, o SpaDeS apresentou uma excelente efetividade na detecção de *spammers* na rede de origem. Não se conhece nenhum outro método que se proponha a fazer esta detecção tão próximo à fonte dos *spams*, reduzindo assim o tráfego na rede. Uma comparação do SpaDeS com outros métodos de detecção de *spammers* disponíveis na literatura e que abordam o problema em outros pontos do sistema é bastante difícil considerando as bases de dados disponíveis para validação do método proposto, uma vez que tais métodos utilizam informações diferentes das exploradas pelo SpaDeS.

#### 4. Conclusões

Neste trabalho foi apresentado, aplicado e validado com dados reais um método para identificação e detecção de *spammers* na rede origem. Este método, denominado SpaDeS (*Spammer Detection at the Source*), tem como principal componente um algoritmo de classificação supervisionada – *Lazy Associative Classification* (LAC) – e utiliza métricas que não requerem a inspeção do conteúdo da mensagem enviada, para classificar os usuários como sendo legítimos ou *spammers*. O SpaDeS apresentou uma excelente efetividade, com taxa de acerto de 98% para usuários legítimos e 94% na classificação de *spammers*. Como trabalho futuro, propõe-se o aprimoramento do método e a construção e validação de um sistema que viabilize o uso do SpaDeS em tempo real.

#### Agradecimentos

Esta pesquisa foi financiada pelo Projeto REBU(CTInfo/CNPq 55.0995/2007-2) e pelo PROBIC/PUC-Minas.

#### Referências

- Castilho, L. H. D., Las-Casas, P. H. B., Dutra, M. D., Ricci, S. M. R., Marques-Neto, H. T., Ziviani, A., Almeida, J. M., e Almeida, V. (2010). Caracterização de tráfego SMTP na Rede de Origem. Em *XXVIII Simpósio Brasileiro de Redes de Computadores e Sistemas Distribuídos*, Gramado, Brasil.
- Gomes, L. H., Cazita, C., Almeida, J. M., Almeida, V., e Jr., W. M. (2007). Workload Models of Spam and Legitimate E-mails. *Performance Evaluation*, 64(7-8):690–714.
- Guerra, P. H. C., Pires, D. E. V., Guedes, D., Jr., W. M., Hoepers, C., e Steding-Jessen, K. (2008). A Campaign-based Characterization of Spamming Strategies. Em *Proceedings of the Fifth Conference on Email and Anti-Spam*, pág. 1–10, Mountain View, CA, USA.
- John, J., Moshchuk, A., Gribble, S. D., e Krishnamurthy, A. (2009). Studying Spamming Botnets Using Botlab. Em *6th USENIX Symp. on Networked Systems Design and Implementation*, Boston, EUA.
- Las-Casas, P. H. B., Guedes, D., Marques-Neto, H. T., Ziviani, A., e Almeida, J. M. (2011). Detecção de Spammers na Rede de Origem. Em *XXIX Simpósio Brasileiro de Redes de Computadores e Sistemas Distribuídos*, Campo Grande, Brasil.
- MessageLabs (2010). MessageLabs Intelligence: November 2010. Online.
- Veloso, A., Meira, W., e Zakib, M. J. (2006). Lazy associative classification. Em *Sixth International Conference on Data Mining*, Hong Kong, China.