

Detecção de Anomalias de Tráfego usando Entropia Não-Extensiva

Marcelo Luís Monsores, Artur Ziviani, Paulo Sérgio Silva Rodrigues

¹Laboratório Nacional de Computação Científica (LNCC/MCT)
Av. Getúlio Vargas, 333
25651-075 – Petrópolis, RJ

(monsores, ziviani, pssr)@lncc.br

Abstract. *Traffic anomalies are characterized by significant and unusual changes in the traffic patterns of one or multiple network links. Given that Internet traffic usually presents characteristics of long-range dependence, an approach to evaluate systems with this behavior is the adoption of nonextensive entropy, a generalization of the traditional Shannon entropy. This paper proposes the use of the nonextensive Tsallis entropy to detect traffic anomalies in autonomous systems. The experimental results show the flexibility of our proposed approach, enabled by the possibility of tuning the detection sensibility, and the better performance achieved by our proposal in comparison with previous approaches found in the literature.*

Resumo. *Anomalias de tráfego são caracterizadas por alterações significativas e pouco comuns nos padrões de tráfego em um ou múltiplos enlaces da rede, sejam estas alterações intencionais ou não. Dado que o tráfego na Internet comumente apresenta características de longo alcance, uma abordagem para avaliar sistemas com este comportamento é a chamada entropia não-extensiva, uma generalização da entropia tradicional de Shannon. Este artigo propõe o uso da entropia não-extensiva de Tsallis para a detecção de anomalias de tráfego em um sistema autônomo. Os resultados experimentais demonstram a flexibilidade da abordagem proposta, devido ao ajuste da sensibilidade da detecção, e o melhor desempenho da mesma em comparação com o estado da arte.*

1. Introdução

A área de metrologia de redes engloba um conjunto de ferramentas e métodos para inferir e melhor compreender o comportamento, a dinâmica e as propriedades da Internet atual [Brownlee e Claffy 2004, Anderson *et al.* 2004, Ziviani e Duarte 2005]. Nesse contexto de metrologia na Internet, a caracterização eficiente de padrões globais no tráfego de rede é crucial para se identificar uma utilização anômala da rede em um sistema autônomo.¹ Anomalias de tráfego em redes são definidas como alterações significativas e pouco comuns nos padrões de volume de tráfego em um ou múltiplos enlaces da rede [Barford *et al.* 2002], sejam elas intencionais ou não. As causas dessas anomalias de tráfego incluem, por exemplo, ataques distribuídos de negação de serviço em curso [Marchette 2001] e mudanças no encaminhamento IP devido a enganos na

¹Neste artigo, utilizamos domínio IP e sistema autônomo como sinônimos.

configuração de roteadores, falha de equipamentos ou modificações nas políticas de roteamento BGP [Roughan *et al.* 2004]. O diagnóstico de anomalias de tráfego, no entanto, apresenta grandes desafios, pois é necessário extrair padrões anômalos de grandes volumes de dados e as causas de anomalias podem ser bastante variadas.

O conceito de diagnóstico de anomalias de tráfego envolve a detecção, a identificação e a quantificação desses fenômenos [Lakhina *et al.* 2004]. A detecção consiste em determinar os pontos no tempo nos quais a rede enfrenta uma anomalia. A identificação envolve a classificação da anomalia a partir de um conjunto de anomalias conhecidas. A quantificação mede a importância da anomalia ao estimar o volume de tráfego anômalo de um determinado tipo presente na rede. Independente das anomalias presentes na rede terem sido causadas intencionalmente ou não, a sua análise é importante, pois essas anomalias de tráfego podem degradar significativamente o serviço de rede, o que torna a sua detecção de grande valia do ponto de vista dos operadores. Portanto, a detecção robusta e confiável de tais anomalias é essencial para a identificação rápida da ocorrência e para a tomada de ações que as corrijam, se necessário.

Este artigo enfoca a detecção de anomalias de tráfego, primordial no processo de diagnóstico de anomalias de tráfego em sistemas autônomos. Em um sistema autônomo, a distribuição de probabilidade de ocorrência de tráfego nos seus diversos nós de entrada e saída do sistema — os Pontos de Presença (PoPs) do sistema autônomo — pode ser usada para quantificar tais anomalias através da medição de sua entropia [MacKay 2003]. Por outro lado, a longa distribuição espacial e temporal do tráfego na Internet pode definir dependências de longo alcance [Karagiannis *et al.* 2004] e sistemas com essas características podem ser avaliados com entropia não-extensiva [Tsallis 1988], que é uma generalização da entropia extensiva tradicional de Shannon [Shannon 1948].

Nós propomos portanto neste artigo a adoção da entropia não-extensiva de Tsallis [Tsallis 1988] para a detecção de anomalias de tráfego em um sistema autônomo, em contraste com trabalhos recentes [Lakhina *et al.* 2005] que se baseiam na entropia clássica de Shannon. Ao adotarmos a entropia não-extensiva de Tsallis concebida para sistemas com dependências de longo alcance, obtemos maior flexibilidade quando comparados à abordagem anterior devido à possibilidade de ajuste da sensibilidade do mecanismo de detecção de anomalias, como demonstram nossos resultados experimentais na Seção 4.4. Essa maior flexibilidade permite a um administrador de um sistema autônomo ajustar a sensibilidade da detecção em acordo com as suas necessidades. Os resultados experimentais também mostram que a abordagem proposta melhora o desempenho da detecção de anomalias de tráfego, aumentando o número de anomalias detectadas e, como consequência, diminuindo o número de falsos negativos.

Este artigo está organizado da seguinte forma. A Seção 2 apresenta brevemente trabalhos anteriores correlatos e indica as contribuições de nossa proposta com relação a estes. Na Seção 3, introduzimos a detecção de anomalias através do conceito de entropia, a classificação de um sistema autônomo em diferentes padrões de tráfego através do valor da entropia nesse sistema e também propomos a detecção de anomalias usando a entropia de não-extensiva de Tsallis. A Seção 4 apresenta nossa análise de desempenho, realizada com dados experimentais, que fornece resultados comparativos entre ambos os métodos entrópicos de detecção de anomalia. Finalmente, concluímos na Seção 5.

2. Trabalhos relacionados

A detecção de anomalias de tráfego de rede usando entropia é bastante recente. Em [Lakhina *et al.* 2005], onde se estuda a mineração de anomalias, é demonstrada a detecção de anomalias usando a entropia de Shannon de forma similar ao descrito na Seção 3.2, utilizando dados originados no domínio Abilene da Internet II (EUA). Os dados adotados são relativos a quantidade de fluxos IP, volume dos fluxos em pacotes e bytes, bem como endereços e portas de origem e de destino de fluxos. Isso possibilitou a análise da detecção de anomalias pela entropia de Shannon e também a identificação e classificação de algumas das anomalias detectadas.

A detecção, identificação e classificação da anomalia é feita através do cálculo de entropia para quatro categorias: entropias de portas de origem, de portas de destino, de endereços de origem e de endereços de destino. A partir da correlação entre as quatro categorias, pode-se comparar os valores de entropia nessas categorias para assim classificar a anomalia. Um exemplo disso é um *port scan*, onde há uma correlação entre as entropias de endereço de destino (onde a entropia indica concentração, em que poucas estações recebem pacotes) e de porta no destino (a entropia indica dispersão, em que várias portas recebem os pacotes nas poucas estações de destino). Em [Lakhina *et al.* 2005], a classificação da anomalia é feita através da comparação das entropias em um sistema de coordenadas, cujas dimensões representam as entropias nas categorias investigadas. Os valores de entropia nas diferentes categorias são agrupados entre si, o que facilita a classificação. Dessa forma, busca-se caracterizar uma anomalia de tráfego ao se comparar a indicação da entropia (concentração ou dispersão) para cada categoria.

A entropia de Tsallis considerando sistemas não-extensivos vem sendo adotada com sucesso em diferentes contextos [Rodrigues *et al.* 2005, Tsallis 2006]. No entanto, até onde alcança o nosso conhecimento, a aplicação desse tipo de abordagem à área de detecção de anomalias é inédita, sendo então de interesse observar o seu comportamento nesse novo contexto. Essa é uma das contribuições de nosso trabalho.

3. Detecção de anomalias usando entropia não-extensiva

Nesta seção, descrevemos a proposta de adotar a entropia não-extensiva de Tsallis para a detecção de anomalias de tráfego. Primeiramente, vamos revisar brevemente o conceito da entropia clássica de Shannon. Em seguida, apresentamos a aplicação do conceito de entropia à detecção de anomalias. Por fim, introduzimos a entropia não-extensiva de Tsallis, junto à nova conceituação que a permeia, permitindo a obtenção de uma visão mais completa do sistema com conseqüente melhora no desempenho na detecção entrópica de anomalia.

3.1. O conceito de entropia e a sua relação com o padrão de tráfego

O conceito de entropia, tal como adotado na área de Teoria da Informação [MacKay 2003], foi definido em [Shannon 1948] como uma medida ligada à quantidade de informações e de incerteza em um dado sistema com base na probabilidade de um determinado fenômeno ocorrer. No caso de um sistema de detecção de anomalia de tráfego, a entropia pode ser usada para avaliar o padrão de comportamento do tráfego de dados em um sistema autônomo de acordo com o volume de tráfego observado, onde esse volume pode ser mensurado por exemplo em número

de fluxos IP ou quantidade de bytes transportados. Dessa forma, pode-se caracterizar o comportamento do tráfego de dados em um domínio IP, determinando se o fluxo de dados está concentrado (poucos pontos da rede recebem grande parte do tráfego de dados) ou disperso (o tráfego de dados encontra-se distribuído por vários pontos da rede). Isto significa representar a quantidade de fluxos em cada ponto de um domínio através das probabilidades que esses fluxos têm para passarem por cada um desses pontos. Com essa representação podemos obter a entropia do sistema que pode nos informar sobre o comportamento do tráfego em termos de sua distribuição no domínio.

Formalmente, a entropia de Shannon [Shannon 1948] é definida como

$$H_S = - \sum_{i=1}^N p_i \log_2 p_i, \quad (1)$$

onde N é o número de eventos a serem considerados e, no nosso caso, o número de pontos de ingresso ou egresso de fluxos ao sistema autônomo em análise. O valor resultante do cálculo da entropia varia entre 0 e $\log_2 N$. A entropia mínima $H_S = 0$ indica concentração máxima, ou seja, ao considerarmos a entrada de dados em um sistema autônomo, todo o tráfego de dados ingressa neste sistema através de um único ponto da rede. Por outro lado, para o mesmo caso de entrada de dados, a entropia máxima $H_S^{\max} = \log_2 N$ indica dispersão total dos fluxos, ou seja, o tráfego ingressante no sistema autônomo encontra-se uniformemente distribuído entre os pontos de rede com todas as probabilidades de fluxos iguais a $\frac{1}{N}$, levando a

$$H_S^{\max} = - \sum_{i=1}^N \left(\frac{1}{N} \log_2 \frac{1}{N} \right) = \log_2 N. \quad (2)$$

Isto significa que quanto maior a entropia, mais disperso o tráfego será considerado. Note que a entropia tem relação direta com a variabilidade das probabilidades consideradas. Quanto mais os fluxos que compõem o tráfego possuem probabilidades similares de ocorrerem em determinados pontos de rede, mais disperso será o padrão de tráfego do sistema como um todo. Quanto maior for a discrepância entre os valores de probabilidade associados aos fluxos componentes do tráfego no sistema, mais concentrado o tráfego estará em um ou poucos pontos do sistema.

3.2. Detecção de anomalias de tráfego usando o conceito de entropia

Considere os dados relativos a um conjunto de fluxos qualquer (*e.g.* endereços IP ou volume transportado em bytes). Do ponto de vista do número de fluxos que passam pelos pontos de borda do domínio, esses dados podem ser agrupados como referentes aos fluxos que saem por cada ponto do domínio e aos fluxos que entram por cada ponto do domínio. Em relação a esses dados, calculamos uma distribuição de probabilidades para os pontos de entrada de fluxo do domínio ou sistema autônomo, chamados origens, e outra distribuição para os pontos de saída de fluxo, chamados destinos.

A classificação combinada da origem ou destino do tráfego de dados como concentrada ou dispersa usando o conceito de entropia permite o estabelecimento de quatro categorias importantes para a detecção de anomalias, pois através delas podemos verificar,

por exemplo, se há alteração repentina do padrão de tráfego em um sistema autônomo. O sistema autônomo no que concerne ao padrão de tráfego em seus pontos de origem e de destino é classificado como:

- Origem concentrada e destino concentrado (CC);
- Origem concentrada e destino disperso (CD);
- Origem dispersa e destino concentrado (DC);
- Origem dispersa e destino disperso (DD).

A Figura 1 ilustra os quatro padrões adotados neste artigo para a classificação do tráfego em um sistema autônomo, representado pelos seus pontos de entrada e saída — pontos de presença ou PoPs. As setas indicam a quantidade de tráfego que entra ou sai por cada um dos pontos de entrada e saída do domínio — pontos de presença ou PoPs. Setas maiores representam um maior volume de tráfego em um ou poucos pontos, enquanto setas menores representam pequenos volumes nos demais pontos, o que pode indicar que há concentração de tráfego na entrada ou saída do domínio. Se as setas de entrada ou de saída de fluxo forem todas do mesmo tamanho, todos os pontos recebem e distribuem dados em mesma quantidade, o que indica uma dispersão uniforme de tráfego. Por exemplo, na Figura 1(a), os fluxos entram predominantemente nos pontos de presença

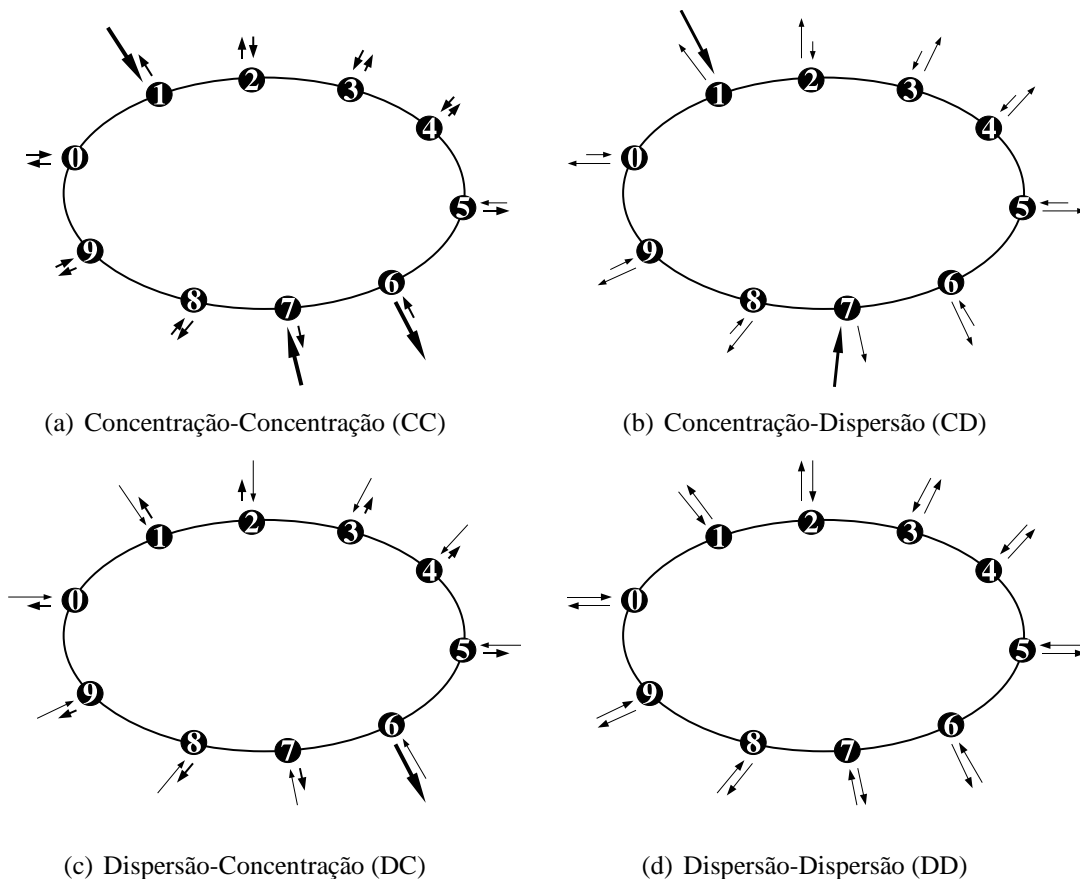


Figura 1. Classificação de um domínio IP quanto à distribuição do tráfego.

1 e 7, ao passo que saem predominantemente pelo ponto 6, havendo uma concentração tanto na entrada quanto na saída de dados no sistema autônomo, caracterizando o padrão de concentração-concentração (CC). De forma similar, um grande volume de tráfego pode chegar predominantemente por poucos PoPs de forma concentrada e, após o roteamento

através do domínio, esse volume pode ser encaminhado a vários outros PoPs para deixar o domínio, caracterizando o padrão Concentração-Dispersão (CD) mostrado na Figura 1(b).

Na Figura 2, ilustramos como os padrões de tráfego podem caracterizar alguns casos particulares que, caso ocorram em maior volume, podem caracterizar uma anomalia passível de detecção. A Figura 2(a) representa o funcionamento de um sistema autônomo durante uma transmissão *multicast* de grande volume, em que um PoP (2) recebe o fluxo *multicast*, que tem sua informação replicada no interior do domínio para seguir seu encaminhamento por vários PoPs de egresso (6, 7, 8, 9). Esse tipo de padrão de tráfego equivale à situação de concentração-dispersão (CD), pois há uma concentração em um ponto na chegada do tráfego, porém a saída do domínio acontece de forma dispersa. Alternativamente, na Figura 2(b), um dado sistema autônomo enfrenta a passagem do tráfego correspondente a um ataque distribuído de negação de serviço (*Distributed Denial of Service* - DDoS) em curso com o objetivo de impedir o acesso à sua vítima. Note que na situação de DDoS em curso, diversos PoPs (1, 2, 3, 4) recebem diversos fluxos destinados a um endereço IP para o qual as tabelas de roteamento indicam convergentemente um único PoP (8) de egresso. Essa situação é classificada como um padrão de tráfego do tipo dispersão-concentração (DC), pois o volume anômalo de tráfego chega ao sistema autônomo por vários PoPs simultaneamente, caracterizando uma dispersão, e é direcionado a um único PoP de destino, indicando uma concentração.

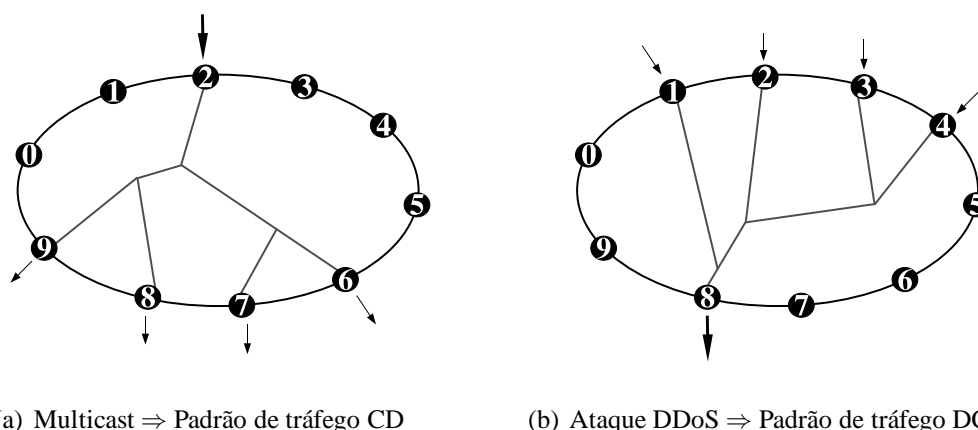


Figura 2. Exemplos ilustrativos de caracterização de anomalias de tráfego.

3.3. Proposta de uso da entropia não-extensiva na detecção de anomalias de tráfego

Nesta subseção, definimos formalmente a entropia não-extensiva de Tsallis. Em seguida, apresentamos como aplicá-la em um mecanismo de detecção de anomalias de tráfego em um sistema autônomo.

3.3.1. Definição da entropia não-extensiva de Tsallis

Em algumas áreas da Física, tais como mecânica e termodinâmica, estudam-se estatisticamente fenômenos microscópicos de um sistema, prevendo-se suas propriedades macroscópicas. Esse estudo pode ser feito pela entropia de Shannon [Shannon 1948], também conhecida nesse contexto como entropia extensiva. Porém, em alguns sistemas de outras áreas da Física, um estudo com método similar gera grandes dificuldades e falhas.

Tais sistemas possuem determinadas características, como dependência de espaço e tempo de longo alcance e comportamento fractal. Para tratar esses casos, Tsallis [Tsallis 1988] propõe o conceito de entropia não-extensiva, generalizando a entropia extensiva convencional de Shannon.

Formalmente, dada uma distribuição de probabilidade $P = \{p_1, p_2, \dots, p_N\}$ com N elementos, onde $0 \leq p_i \leq 1$ e $\sum_i p_i = 1$, a entropia não-extensiva proposta em [Tsallis 1988] é definida por:

$$H_q = \frac{1 - \sum_{i=1}^N p_i^q}{q - 1}. \quad (3)$$

Por outro lado, é interessante avaliar a Eq. (3) quando q tende a 1:

$$\begin{aligned} H_1 = \lim_{q \rightarrow 1} H_q &= \lim_{q \rightarrow 1} \frac{1 - \sum_{i=1}^N p_i^q}{q - 1} \\ &= \lim_{q \rightarrow 1} \frac{\sum_{i=1}^N p_i - \sum_{i=1}^N p_i^q}{q - 1} \\ &= \sum_{i=1}^N p_i \lim_{q \rightarrow 1} \frac{(1 - p_i^{q-1})}{q - 1} \\ &= - \sum_{i=1}^N p_i \ln p_i, \end{aligned} \quad (4)$$

que é equivalente à entropia clássica de Shannon², mostrada em Eq. (1). Esse resultado demonstra que a entropia de Tsallis é uma generalização da entropia de Shannon. Assim, os sistemas entrópicos podem ser caracterizados da seguinte forma:

1. sistemas extensivos ($q = 1$);
2. sistemas subextensivos ($q > 1$);
3. sistemas superextensivos ($q < 1$).

O resultado da entropia de Tsallis pode variar entre 0, que caracteriza a concentração máxima, e H_q^{\max} , indicador de dispersão máxima, onde

$$H_q^{\max} = \frac{1 - N^{1-q}}{q - 1}. \quad (5)$$

Esse novo conceito de entropia se caracteriza pela introdução do parâmetro entrópico q , que está relacionado ao grau de extensividade do sistema e define a escala de

²Na área de física, em particular na mecânica estatística, o conceito de entropia é definido como $S = - \sum_i p_i \ln p_i$, conhecida também como entropia de Boltzman-Gibbs, sendo esta de fato o ponto original da proposta de Tsallis [Tsallis 1988]. Na realidade, a base do logaritmo no cálculo da entropia pode ser definida arbitrariamente, pois a diferença entre as definições é uma constante, ou seja, $S = kH_S$ [Jaynes 1957], onde H_S é a entropia de Shannon. De fato, Shannon [Shannon 1948] escolheu convenientemente a base 2, definindo o conceito de entropia para a teoria de informação, enquanto na física utiliza-se mais comumente o logaritmo natural.

medição da entropia não-extensiva. Para melhor entender o papel do parâmetro q , note que na entropia de Shannon os eventos com probabilidades muito elevada ou muito baixa não possuem grande influência no valor da entropia. Em contraste, na entropia de Tsallis, no caso de $q > 1$, os eventos com maiores probabilidades contribuem mais para o valor da entropia do que eventos de baixa probabilidade. De forma inversa, no caso de $q < 1$, os eventos com menores probabilidades contribuem mais para o valor da entropia do que eventos de alta probabilidade. Portanto, a variação de q modifica a contribuição relativa de um dado evento para a soma total. Isto permite aguçá-la a percepção de um sistema baseado em entropia de Tsallis a eventos de maior ou menor contribuição relativa, o que é decisivo para a detecção de anomalias como mostrado na análise de nossos resultados experimentais na Seção 4.

3.3.2. Aplicação da entropia não-extensiva de Tsallis à detecção de anomalias

Neste artigo, o parâmetro q da entropia de Tsallis desempenha o papel de calibrador do nível de detalhamento (ou sensibilidade) da detecção de um determinado padrão de tráfego: isto é, através dele pode-se detectar mais ou menos anomalias de um certo padrão. Assim, em geral, busca-se achar o valor ótimo de q : ou seja, o $q_{\text{ótimo}}$ que confere ao sistema o maior nível de sensibilidade em relação ao fenômeno estudado — no nosso caso, a detecção de anomalias de tráfego. Como ilustração do efeito da variação do parâmetro q em nosso sistema de detecção de anomalias de tráfego, pode-se fazer uma analogia da variação deste parâmetro com a regulagem de sensibilidade em um sistema de detecção de metais em uma porta de banco: ou seja, variando-se a sensibilidade, regula-se a quantidade de metal que ativará a detecção. De forma similar, ao variarmos o parâmetro q , podemos regular a sensibilidade da detecção de anomalias de tráfego em nossa proposta, como detalhado em nossos resultados experimentais na Seção 4.4.

4. Análise de desempenho

Nesta seção, avaliamos o mecanismo proposto de detecção de anomalias de tráfego usando entropia não-extensiva em um sistema autônomo. Para tanto, são avaliados a capacidade desse mecanismo de identificar padrões de tráfego pré-definidos e o ajuste da sensibilidade do mesmo. Também é realizada uma comparação com propostas anteriores, caracterizando o melhor desempenho de nossa proposta em relação ao estado da arte atual.

4.1. Dados experimentais utilizados

Os dados utilizados em nossa análise de desempenho contêm informações sobre os volumes de fluxos IP e bytes em PoPs do domínio Abilene, pertencente à Internet 2 nos EUA. Vale ressaltar que esses dados nos foram fornecidos pelos autores de [Lakhina *et al.* 2005], pioneiros no uso da entropia clássica de Shannon para detecção de anomalias, o que nos permite uma comparação direta de desempenho de nossa proposta com a literatura.

Os fluxos de dados são analisados através do par origem-destino por onde trafegam para entrar e sair do sistema autônomo em estudo. Nos dados experimentais considerados, são 11 os pontos de origem e destino (*i.e.* os PoPs), totalizando portanto 121

pares origem-destino. Os dados foram coletados entre 7 e 11 de abril de 2003 (segunda-feira a domingo), com coleta dos totais de cada tipo de dado considerado, fluxos IP ou bytes, em intervalos de 5 minutos, totalizando 2016 coletas ao todo. Assim, podemos ter a distribuição dos dados em cada intervalo considerando separadamente a origem ou o destino dos mesmos no sistema autônomo.

4.2. Classificação do tráfego usando entropia

Nesta subseção, investigamos a capacidade de classificação da entropia extensiva de Shannon e da entropia não-extensiva de Tsallis. Ao considerarmos um intervalo de coleta qualquer, as quantidades de fluxo nos pontos de origem são transformadas em probabilidades em relação ao total de tráfego entrante no domínio IP de interesse. Através dessa distribuição, obtemos a entropia de Shannon com respeito aos seus pontos de origem (*i.e.* de entrada) nesse domínio. Esse procedimento é repetido de forma semelhante para a obtenção da distribuição de tráfego com relação aos pontos de destino (saída do domínio). Utilizando essas distribuições de probabilidades, calculamos a entropia correspondente e classificamos o padrão de tráfego em um determinado intervalo de tempo, com relação ao seu par origem-destino no domínio, em Concentração-Concentração (CC), Concentração-Dispersão (CD), Dispersão-Concentração (DC) ou Dispersão-Dispersão (DD) (ver Seção 3.2).

A caracterização do comportamento do sistema autônomo ao longo do tempo é realizada calculando-se as entropias em cada um dos intervalos de tempo. Para a detecção de anomalias de tráfego, normalizamos a entropia observada para o tráfego em relação à entropia máxima, identificando através da entropia normalizada padrões de tráfego disperso e concentrado com relação à origem e ao destino do volume de tráfego nas bordas do sistema autônomo. A partir da entropia normalizada, um limiar é estabelecido para diferenciar entre um padrão concentrado ou disperso. Esse limiar varia de acordo com o tipo de dado e de padrão de concentração ou dispersão que se quer representar. O padrão de tráfego de dados em um domínio pode então ser classificado como concentrado (entropia abaixo do limiar) e disperso (entropia acima do limiar). O valor da entropia normalizada H_{norm} é definido por

$$H_{\text{norm}} = \begin{cases} \frac{H_S}{H_S^{\text{max}}} & \text{se } q = 1, \\ \frac{H_q}{H_q^{\text{max}}} & \text{se } q \neq 1. \end{cases} \quad (6)$$

Todavia, note que o valor máximo da entropia de Tsallis H_q^{max} é dependente do valor de q , como anteriormente apresentado na Eq. (5) e ilustrado na Figura 3 para $N = 11$. De fato, a variação de q , com a conseqüente variação de H_q^{max} , afeta a sensibilidade de detecção de um sistema baseado em entropia não-extensiva e seus efeitos são avaliados na Seção 4.4.

Tanto em volumes de fluxos IP quanto em bytes, em todos os intervalos considerados, os valores de entropia observados se mantiveram acima de $H_{\text{norm}} = 0.86$ tanto na origem quanto no destino, o que pode ser considerada um nível de dispersão elevado tanto na entrada quanto na saída do tráfego. Para valores de entropia abaixo de $H_{\text{norm}} = 0.91$, ainda observamos um número considerável de concentrações de tráfego

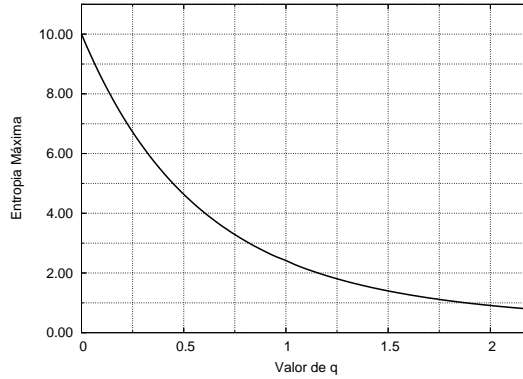


Figura 3. Relação entre o valor de q e a entropia máxima de Tsallis H_q^{\max} ($N = 11$).

em alguns intervalos de tempo, em contraste com a dispersão reinante na maior parte dos 2016 intervalos considerados. Neste trabalho, esse limiar é definido de forma empírica, sendo a otimização desse parâmetro parte de nossos trabalhos futuros.

Com base na observação dos valores relativos de entropia e para verificarmos mais detalhadamente o comportamento em termos de concentração e dispersão de tráfego no sistema autônomo, convencionamos então que o padrão de tráfego é considerado: (i) *concentrado* se $H_{\text{norm}} < 0.91$; ou (ii) *disperso* se $H_{\text{norm}} \geq 0.91$. Assim, há um critério único de classificação, independente do valor (variável) da entropia máxima dos diferentes métodos, que mantém a comparabilidade entre as abordagens. Ao adotarmos essa classificação, o padrão de tráfego considerado normal é o DD, onde os volumes de fluxos IPs ou bytes encontram-se bem distribuídos pelo total de PoPs na origem ou no destino. Dessa forma, as anomalias de tráfego se configuram quando o sistema detecta padrões de tráfego distintos de DD.

Para ilustrar a operação do mecanismo de detecção de anomalias proposto, apresentamos na Figura 4 o comportamento da entropia de Tsallis normalizada ao longo de todo o período de tempo coberto pelos dados experimentais. Nessa ilustração, adotamos arbitrariamente $q = 0.9$, que é o valor mais comumente verificado como $q_{\text{ótimo}}$ nos dados experimentais como detalhado adiante na Seção 4.4. A linha horizontal tracejada indica o limiar de detecção estabelecido de $H_{\text{norm}} = 0.91$ para a classificação entre padrão concentrado (abaixo do limiar) e disperso (acima ou igual ao limiar). Para facilitar o acompanhamento da evolução do padrão do tráfego ao longo da semana de observação, demarcamos por linhas verticais os intervalos correspondentes aos dias da semana e os subdividimos em períodos de 6 horas, resultando em referências para os períodos de madrugada (00:00h–06:00h), manhã (06:00h–12:00h), tarde (12:00h–18:00h) e noite (18:00h–24:00h) de cada dia observado.

A combinação de estado concentrado ou disperso no mesmo intervalo de tempo na origem e no destino caracteriza o padrão de tráfego do domínio. Por exemplo, nas Figuras 4(a) e 4(b), ao meio-dia de sexta-feira, há uma tendência à concentração de fluxos IP na origem e no destino, indicada pela queda da entropia normalizada. Entretanto, essa concentração não é suficientemente elevada para caracterizar uma possível anomalia de tráfego presente. Isso resulta em um padrão DD para fluxos IP neste horário. Por outro lado, ao examinarmos as Figuras 4(c) e 4(d), referentes ao volume de bytes, observamos que em relação aos destinos há uma concentração elevada suficiente para ultrapassar o

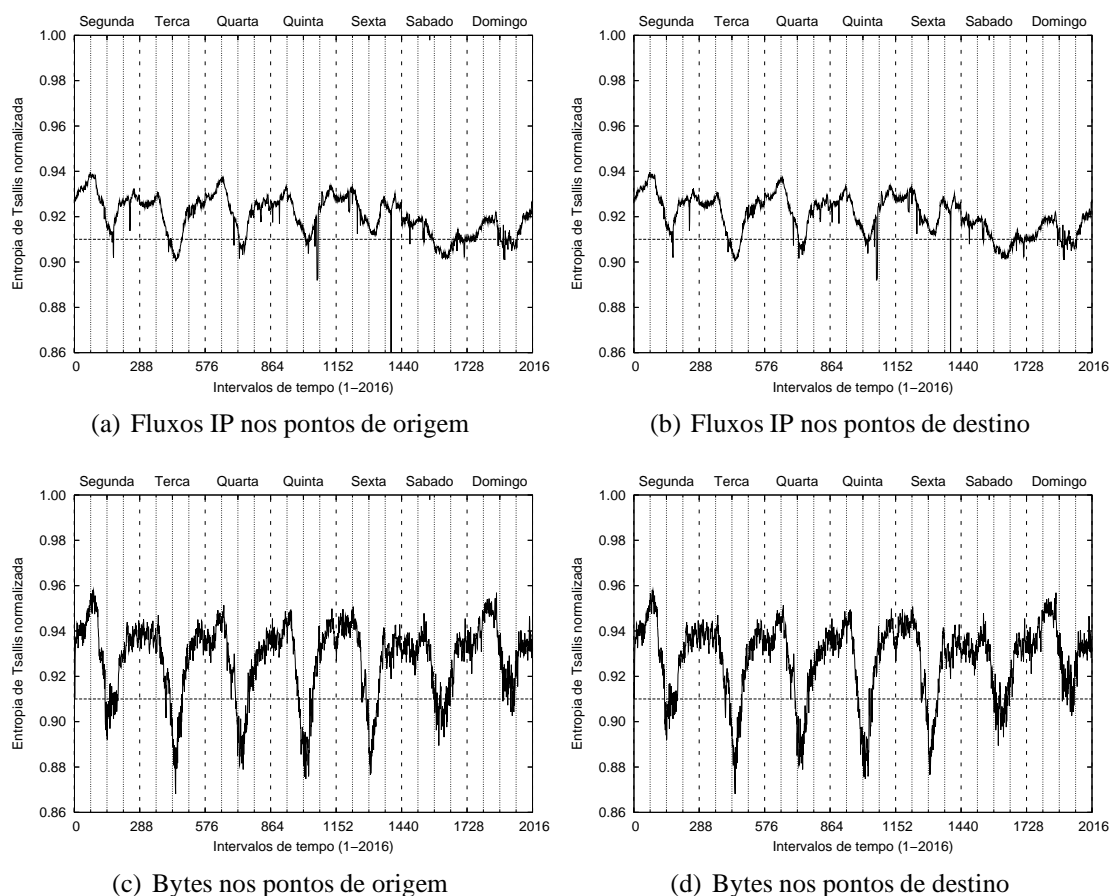


Figura 4. Entropia de Tsallis normalizada ($q = 0.9$) ao longo do tempo.

limiar estabelecido. Como consequência, para bytes, o padrão de tráfego se constitui DC. De forma global, tanto para o volume de fluxos IP quanto para o volume de bytes, o padrão de tráfego se mantém primordialmente em DD, por isso mesmo considerado o padrão esperado de operação. No entanto, é interessante observar a tendência à concentração detectada no volume de tráfego de fluxos IP e de bytes nos períodos diurnos.

4.3. Desempenho em quantidade de anomalias de tráfego detectadas

Nesta subseção, avaliamos o desempenho da proposta de adoção da entropia não-extensiva de Tsallis em um mecanismo de detecção de anomalias. Para tanto, aplicamos ao sistema os nossos dados experimentais, quantificando as ocorrências dos diferentes padrões de tráfego (CC, CD, DC, DD) em número de intervalos de tempo caracterizados com um padrão e, como consequência, a quantidade de anomalias de tráfego detectadas (os padrões que diferem de DD).

A Figura 5 apresenta a quantidade de detecções de cada padrão de tráfego ao longo do período coberto por nossos dados experimentais. Ao realizarmos a experimentação usando $0 \leq q \leq 2$, podemos aferir o valor $q_{\text{ótimo}}$ que proporciona o maior número de detecções de anomalias. Essa análise também permite comparar diretamente o desempenho da entropia não-extensiva de Tsallis ($q \neq 1$) com a de Shannon ($q = 1$), o que equivale a comparar nossa proposta com a adotada em [Lakhina *et al.* 2005] utilizando os mesmos dados experimentais.

A quantidade de detecções de padrões de tráfego CC, CD e DC ao longo dos dados

experimentais é mostrada nas Figuras 5(a), 5(b) e 5(c), respectivamente. Esses padrões são considerados os indicadores da presença de algum tráfego anômalo no domínio IP, visto que representam uma distinção em relação ao padrão de tráfego dominante DD. Para ambas as métricas de volume de tráfego, há consistentemente um valor $q_{\text{ótimo}}$ que proporciona à entropia não-extensiva de Tsallis um melhor desempenho em relação à entropia clássica de Shannon. Em geral, o volume de bytes apresenta um maior número de detecção de estados de concentrações do que os fluxos IPs. A Tabela 1 resume a melhoria de desempenho possibilitada pelo uso da entropia de Tsallis na detecção de um número maior de anomalias, podendo o ganho em relação à entropia de Shannon alcançar 63% nos dados experimentais considerados, como em fluxos IP para um padrão de tráfego CD.

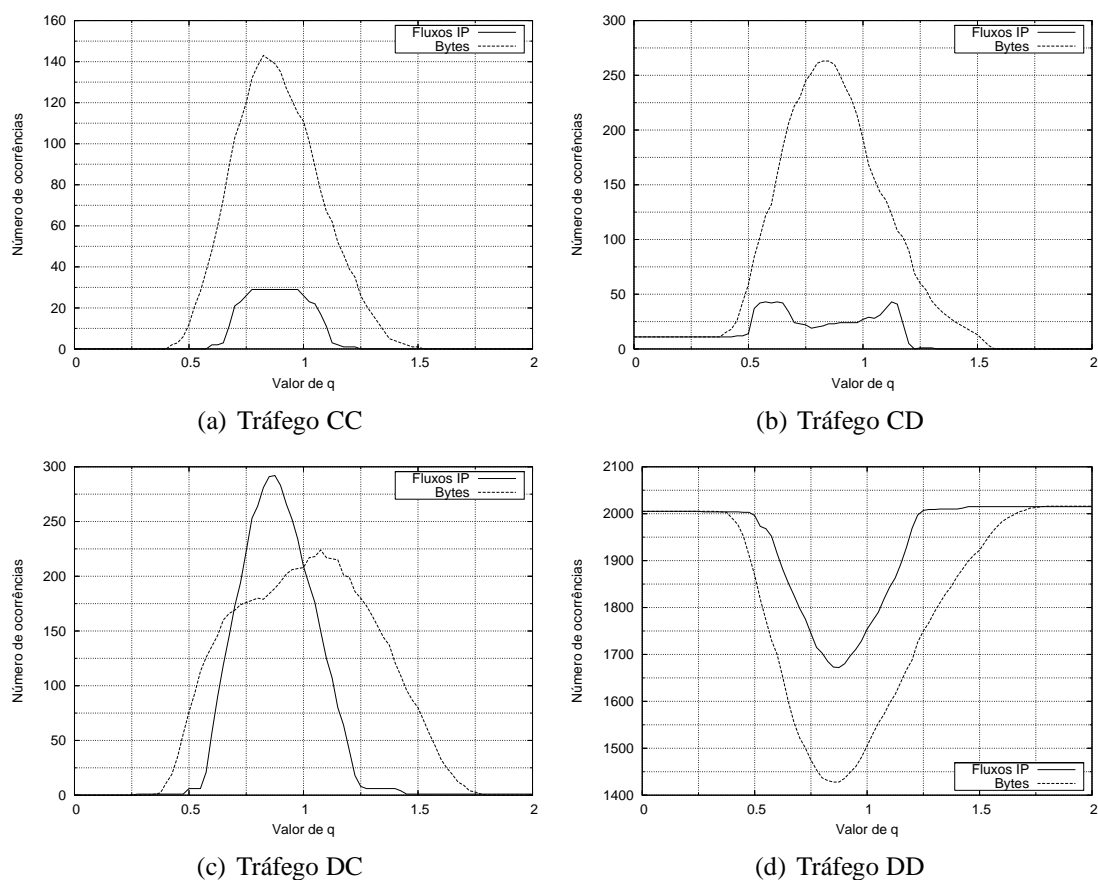


Figura 5. Quantidade de identificação de cada padrão de tráfego.

Tabela 1. Comparação de desempenho no número de anomalias detectadas.

| Métrica | Padrão de tráfego | Número de anomalias detectadas | | Ganho (%) |
|-----------|-------------------|--------------------------------|------------------------|-----------|
| | | H_S | $H_{q_{\text{ótimo}}}$ | |
| Fluxos IP | CC | 26 | 29 | 12% |
| Fluxos IP | CD | 27 | 44 | 63% |
| Fluxos IP | DC | 209 | 294 | 41% |
| Bytes | CC | 111 | 143 | 29% |
| Bytes | CD | 191 | 263 | 38% |
| Bytes | DC | 208 | 224 | 8% |

É interessante, por diferentes aspectos, analisar com atenção a Figura 5(d), que mostra a quantidade de intervalos de tempo que foram identificados como tendo o padrão de tráfego DD. Primeiro, essa figura ilustra a predominância do padrão de tráfego DD ao mostrar que quando não há uma detecção de anomalia, o estado do domínio permanece inalterado em DD. Segundo, ao detectarmos mais anomalias conforme o valor de q tende a $q_{\text{ótimo}}$, o número de intervalos classificados em DD naturalmente diminui, pois intervalos anteriormente percebidos como DD em um certo instante passam a ser identificados como tendo a origem, o destino, ou mesmo ambos com concentração de tráfego, seja para a métrica de fluxos IP, seja pela métrica de bytes. Na realidade, ao compararmos o desempenho da entropia não-extensiva de Tsallis com a entropia de Shannon na Figura 5(d), observamos o número de falsos negativos que são apresentados na entropia de Shannon e evitados ao usarmos a entropia não-extensiva.

4.4. Avaliação da sensibilidade de detecção de anomalias

A sensibilidade da detecção está relacionada ao grau de magnitude da anomalia em relação ao tráfego considerado normal. Através da variação do parâmetro q , a entropia de Tsallis permite aumentar ou diminuir a sensibilidade do mecanismo de detecção, permitindo assim a um operador calibrar o quão sensível o sistema se torna à presença de uma anomalia. Portanto, ao analisarmos a sensibilidade da detecção de anomalias podemos buscar um q ótimo ($q_{\text{ótimo}}$) que confere ao sistema o maior nível de sensibilidade em relação à presença de anomalias de tráfego no sistema autônomo monitorado.

Para analisar de forma controlada a sensibilidade da detecção de anomalias usando entropia não-extensiva, realizamos um experimento, cujo objetivo é investigar o limite de pior caso, buscando-se qual nível de concentração de tráfego faz disparar a detecção no pior caso. Para tanto, uma anomalia artificial (concentração) é inserida progressivamente em um domínio cujo padrão de tráfego esteja em DD com dispersão máxima, o que equivale a ter a entropia H_q^{max} para ($q \neq 1$) e H_S^{max} ($q = 1$). Essa configuração se constitui o pior caso, pois para a detecção da concentração obriga uma maior inserção artificial para que a entropia ultrapasse o limiar de sensibilidade para detecção de uma concentração partindo de uma dispersão absoluta. Dessa forma, podemos simular a condição de concentração exigida para caracterizar um determinado padrão de tráfego considerado anômalo (CC, CD, DC) e verificar a que nível de concentração o padrão artificialmente enxertado pode ser detectado. O volume de dados extra da anomalia artificial é controlado como uma fração do tráfego total no domínio. Por exemplo, para simular uma concentração, o volume extra é acrescido aos dados de um único PoP. Ao variarmos a magnitude da anomalia artificial em relação ao tráfego total do sistema autônomo, verificamos a partir de que magnitude de anomalia o mecanismo de detecção acusa a presença da mesma. Portanto, para cada magnitude de anomalia considerada, é identificado o valor de q que registra a modificação de padrão de tráfego detectado, passando do tipo original para o tipo simulado.

Na Figura 6, mostramos os resultados dessa análise de sensibilidade controlada para o pior caso, o que exigiria mais concentração e assim estabelecemos limites. O ponto mínimo dessa curva representa o valor de q capaz de detectar o menor volume de tráfego artificialmente inserido, caracterizando o limiar de sensibilidade no pior caso. Esse pior caso indica que qualquer anomalia com um volume de pelo menos 40% do tráfego total é detectada variando-se o valor de q .

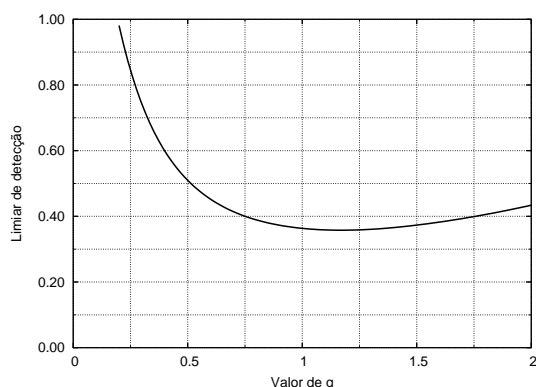


Figura 6. Sensibilidade de detecção de anomalias no pior caso.

Vale ressaltar que o resultado apresentado na Figura 6 refere-se ao pior caso. Observando os dados experimentais podemos constatar limiares de sensibilidade bem menos elevados. Para ilustrar essa constatação, apresentamos na Figura 7 os resultados da análise de sensibilidade para um intervalo de tempo arbitrário de nossos dados experimentais, classificado originalmente como DD usando a entropia de Shannon e Tsallis. O ponto mínimo de cada curva da Figura 7 representa o valor de q capaz de detectar a anomalia que causa o padrão de tráfego simulado (CC, CD, DC) para o menor volume artificialmente inserido, caracterizando o limiar de sensibilidade. Esse valor de q é o chamado $q_{\text{ótimo}}$ para aquele conjunto de dados experimentais naquele intervalo. Note também na Figura 7

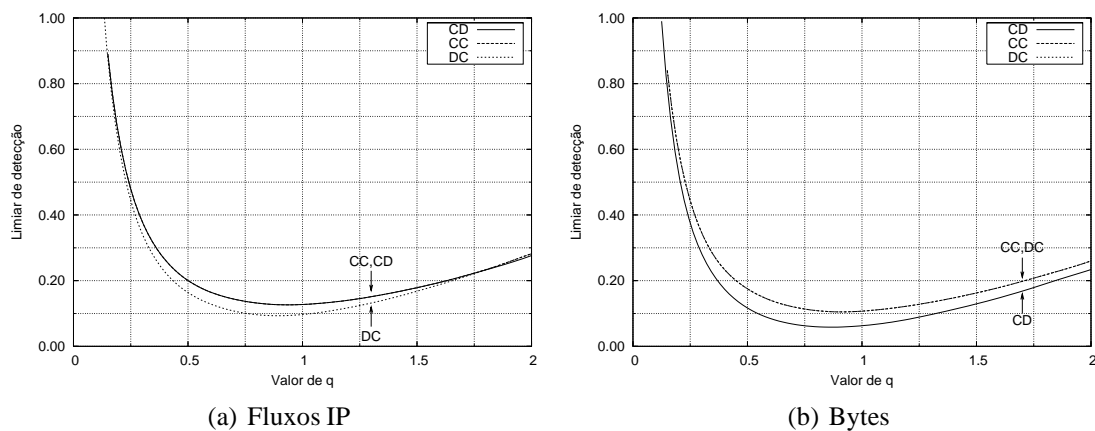


Figura 7. Amostra da sensibilidade de detecção com entropia não-extensiva.

que todas as curvas apresentam uma forma bastante semelhante: decréscimo rápido em exponencial para $q < q_{\text{ótimo}}$ associado a um crescimento mais lento para $q > q_{\text{ótimo}}$, sendo o ponto de inflexão equivalente ao $q_{\text{ótimo}}$. Esse resultado sugere a possibilidade de modelagem dessa relação, o que deixamos para figurar entre nossos trabalhos futuros.

Na Tabela 2, estão listados os valores de $q_{\text{ótimo}}$ identificados para cada métrica de volume considerada e padrão de tráfego associado a um sistema autônomo. Os limiares de sensibilidade dos padrões de tráfego são bastante inferiores ao pior caso, indicando a sensibilidade que pode ser alcançada com a entropia não-extensiva.

Tabela 2. Amostra da sensibilidade alcançada com entropia não-extensiva.

| Métrica | Padrão de tráfego | $q_{\text{ótimo}}$ | Sensibilidade (%) |
|-----------|-------------------|--------------------|-------------------|
| Fluxos IP | CC | 0.90 | 2.12 |
| Fluxos IP | CD | 0.89 | 4.01 |
| Fluxos IP | DC | 0.90 | 2.12 |
| Bytes | CC | 0.90 | 2.54 |
| Bytes | CD | 0.83 | 2.70 |
| Bytes | DC | 0.90 | 2.54 |

5. Conclusão

Este artigo apresentou a detecção de anomalias de tráfego em um domínio IP ou em um sistema autônomo usando a entropia não-extensiva de Tsallis em contraste com a detecção via entropia de Shannon encontrada na literatura. Os resultados mostram que a detecção de anomalias por entropia de Tsallis oferece a vantagem de ser mais flexível que a detecção por entropia de Shannon, levando a um melhor desempenho. Essa flexibilidade permite o detalhamento da detecção através de um ajuste da sensibilidade do sistema com a variação do valor q . Pode-se, por exemplo, avaliar o nível de concentração de tráfego em um domínio IP, em quais instantes de tempo houve maior concentração num domínio, se há alguma mudança brusca de padrão de tráfego ao longo do tempo, entre outros. Como trabalho futuro, visamos a implementação de uma ferramenta que se baseia na entropia não-extensiva de Tsallis para controle e avaliação de um domínio administrativo de rede.

Agradecimentos

Os autores são gratos a A. Lakhina (Boston University, EUA), M. Crovella (Boston University, EUA) e C. Diot (Thomson R&D, França) que gentilmente cederam os dados utilizados em [Lakhina *et al.* 2005] para comparação. Os autores também agradecem a B. Schulze e M. Trindade dos Santos, ambos do LNCC, pelos esclarecimentos sobre a entropia não-extensiva de Tsallis, e a A. T. A. Gomes, também do LNCC, pela revisão criteriosa do texto. M. L. Monsores é graduando do curso de Tecnologia de Informação e Comunicação do Instituto Superior de Tecnologia em Ciência da Computação de Petrópolis (ISTCCP) e bolsista do PIBIC/CNPq no LNCC. P. S. S. Rodrigues é pesquisador bolsista DCR no LNCC com financiamento do CNPq/FAPERJ.

Referências

- Anderson, T., Crovella, M., e Diot, C. (2004). Internet measurements: Past, present, and future. Available at <http://www.ics.uci.edu/~xwy/ics243c/lec-notes/measurements-survey.pdf>.
- Barford, P., Kline, J., Plonka, D., e Ron, A. (2002). A signal analysis of network traffic anomalies. In *Proc. of ACM/SIGCOMM Internet Measurement Workshop – IMW 2002*, Marseille, France.
- Brownlee, N. e Claffy, K. C. (2004). Internet measurement. *IEEE Internet Computing*, 8(5):30–33.
- Jaynes, E. T. (1957). Information theory and statistical mechanics. *The Physical Review*, 106(4):620–630.

- Karagiannis, T., Molle, M., e Faloutsos, M. (2004). Long-range dependence: Ten years of Internet traffic modeling. *IEEE Internet Computing*, 8(5):57–64.
- Lakhina, A., Crovella, M., e Diot, C. (2004). Diagnosing network-wide traffic anomalies. In *Proc. of the ACM SIGCOMM'2004*, Portland, OR, USA.
- Lakhina, A., Crovella, M., e Diot, C. (2005). Mining anomalies using traffic feature distributions. In *Proc. of the ACM SIGCOMM'2005*, Philadelphia, PA, USA.
- MacKay, D. J. C. (2003). *Information Theory, Inference, and Learning Algorithms*. Cambridge University Press, Cambridge, UK.
- Marchette, D. J. (2001). *Computer Intrusion Detection and Network Monitoring*. Springer-Verlag.
- Rodrigues, P. S. S., Giraldi, G. A., e Araújo, A. A. (2005). Using Tsallis entropy into a Bayesian network for CBIR. In *Proc. of the IEEE International Conference on Image Processing – ICIP 2005*, Genova, Italy.
- Roughan, M., Griffin, T., Mao, M., Greenberg, A., e Freeman, B. (2004). IP forwarding anomalies and improving their detection using multiple data sources. In *Proc. of the ACM SIGCOMM'2004 Workshop on Network Troubleshooting*, Portland, OR, USA.
- Shannon, C. E. (1948). A mathematical theory of communication. *Bell System Technical Journal*, 27:379–423 and 623–656.
- Tsallis, C. (1988). Possible generalization of Boltzmann-Gibbs statistics. *Journal of Statistical Physics*, 52(1-2):479–487.
- Tsallis, C. (2006). *Bibliography in Nonextensive Statistical Mechanics and Thermodynamics*. <http://www.cbpf.br/GrupPesq/StatisticalPhys/biblio.htm>.
- Ziviani, A. e Duarte, O. C. M. B. (2005). *Metrologia na Internet*, in Minicursos do XXIII Simpósio Brasileiro de Redes de Computadores - SBRC'2005, pages 285–329. SBC, Fortaleza, CE.