



UNIVERSIDADE
ESTADUAL DE LONDRINA

MOISÉS FERNANDO LIMA

**DETECÇÃO DE ANOMALIAS UTILIZANDO
ASSINATURA DIGITAL DE SEGMENTO DE REDE**

MOISÉS FERNANDO LIMA

**DETECÇÃO DE ANOMALIAS UTILIZANDO
ASSINATURA DIGITAL DE SEGMENTO DE REDE**

Trabalho de Dissertação de
Mestrado apresentado à Universidade Estadual de
Londrina como parte dos requisitos para
obtenção do título de Mestre em Ciência da
Computação.

Orientador: Prof. Dr. Mario Lemes Proença Jr.

Londrina
2011

**Catálogo elaborado pela Divisão de Processos Técnicos da Biblioteca Central da
Universidade Estadual de Londrina**

Dados Internacionais de Catalogação-na-Publicação (CIP)

L732d Lima, Moisés Fernando.
Detecção de anomalias utilizando assinatura digital de segmento de rede / Moisés Fernando Lima. — Londrina, 2011
78 f.: II

Orientador: Mario Lemes Proença Junior
Dissertação (Mestrado em Ciência da Computação) – Universidade Estadual de Londrina, Centro de Ciências Exatas, Programas de Pós-Graduação em Ciências da Computação, 2011.
Inclui bibliografia

Redes de computação — Anomalias — Teses. 2. Telecomunicações — Tráfego — Teses. 3. Assinaturas digitais — Teses. 4. Redes de computação — Medidas de segurança— Teses. I. Proença Junior, Mário Lemes. II. Universidade Estadual de Londrina. Centro de Ciências Exatas. Programa de Pós-Graduação em Ciência da Computação. III. Título.

CDU 519.68.022

MOISÉS FERNANDO LIMA

**DETECÇÃO DE ANOMALIAS UTILIZANDO ASSINATURA
DIGITAL DE SEGMENTO DE REDE**

Trabalho de Dissertação de Mestrado apresentado à Universidade Estadual de Londrina como parte dos requisitos para obtenção do título de Mestre em Ciência da Computação.

Orientador: Prof. Dr. Mario Lemes Proença Jr.

BANCA EXAMINADORA

Prof. Dr. Mario Lemes Proença Jr.
UEL – Londrina - PR

Prof. Dr. José Valdeni de Lima
UFRGS - Rio Grande do Sul - RS

Prof. Dr. Rodolfo Miranda de Barros
UEL – Londrina - PR

Prof. Dr. Taufik Abrão
UEL – Londrina - PR

Londrina, 8 de fevereiro de 2011.

AGRADECIMENTOS

Primeiramente a Deus pela presença constante em minha vida.

Em especial ao meu orientador Prof. Dr. Mario Lemes Proença Jr. por toda a paciência, apoio e dedicação durante a orientação.

Aos meus pais Valdir e Marina, pelo suporte e confiança demonstrados ao longo de toda minha vida.

A minha noiva Pollyana pelo apoio e por compreender minhas ausências.

Ao Prof. Dr. Taufik Abrão, ao Lucas Dias H. Sampaio e também ao Bruno B. Zarpelão pelo companheirismo e cooperação.

A todos colegas, professores e funcionários do Departamento de Computação que contribuíram para a realização deste trabalho

*Aos meus pais que sempre me deram amor e suporte
A minha noiva pelo apoio e compreensão durante essa jornada*

LIMA, Moisés Fernando. *Detecção de anomalias utilizando assinatura digital de segmento de rede*. 2011. 78 f. Dissertação (Mestrado) - Universidade Estadual de Londrina, Londrina. 2011.

RESUMO

A detecção de anomalias em redes caracteriza-se pela busca de comportamentos incomuns no tráfego, que possam vir a comprometer a segurança, o desempenho e a integridade das informações. Consiste de um problema importante e difícil que tem sido tratado dentro de diversos domínios e áreas de pesquisa, destacando-se principalmente a estatística aprendizagem de máquina e mineração de dados, dentre outras tais como teoria da informação e teoria espectral. Neste trabalho é realizada uma revisão da literatura sobre as técnicas recentemente utilizadas na detecção de anomalias, a fim de criar um embasamento para o desenvolvimento de um Sistema de Detecção de Anomalias (SDA). Deste modo, foi desenvolvido um sistema baseado em uma técnica heurística que analisa dados de tráfego coletados da MIB através do protocolo SNMP. O sistema é baseado na utilização de três modelos, o primeiro é o algoritmo de clusterização *K-means* que consiste de um método de análise de clusters e classificação de dados. O segundo método denominado *Particle Swarm Optimization* (PSO), classificado como um algoritmo evolucionário, consiste de uma ferramenta heurística de otimização numérica altamente eficiente, com baixa complexidade computacional e capacidade de escapar de ótimos locais. O terceiro modelo denominado Assinatura Digital de Segmento de Rede (DSNS), caracteriza-se pela criação de perfis de comportamento normal de tráfego de rede, gerado pela ferramenta de Gerenciamento de *Backbone* Automatizado (GBA), com base em dados históricos da rede. Da combinação das técnicas PSO e *K-means* deu-se origem o algoritmo denominado PSO-Cls, o qual é a base do SDA desenvolvido. A fim de avaliar a qualidade do sistema desenvolvido, foram realizados diferentes experimentos sob a perspectiva de diferentes cenários. Também foram realizados estudos a respeito da complexidade e otimização dos parâmetros do PSO-Cls. Foram utilizados nos experimentos dados reais coletados na rede da Universidade Estadual de Londrina. A fim de comparar os resultados obtidos pelo sistema desenvolvido foram implementados dois SDA's, um baseado em algoritmo determinístico e outro baseado na Análise de Componentes Principais (do inglês *Principal Component Analysis*, PCA). Finalmente, os resultados obtidos através dos experimentos são apresentados e discutidos.

Palavras-Chaves: Detecção de anomalias. DSNS. PSO. K-means. PCA.

LIMA, Moisés Fernando. *Detecção de anomalias utilizando assinatura digital de segmento de rede*. 2011. 78 f. Dissertação (Mestrado) - Universidade Estadual de Londrina, Londrina. 2011.

ABSTRACT

The anomaly detection in networks is characterized by the unusual behaviors presented in the network traffic and may compromise the security, the network performance and the information integrity. It consists of an important problem that has been treated in various domains and research areas, emphasizing mainly statistical, machine learning and data mining, among others, such as information theory and spectral theory. In this work a review of current literature and recent techniques used in anomaly detection in order to create a foundation for developing an anomaly detection system (SDA). Thus, a system was developed based on a heuristic technique that analyzes traffic collected from the MIB using the SNMP protocol. The system is based on the use of three models, the first is the clustering algorithm K-means which consists of a method of cluster analysis and data classification. The second method called Particle Swarm Optimization (PSO), classified as an evolutionary algorithm, consist in a heuristic tool of a highly efficient numerical optimization with low computational complexity and ability to escape local optima. The third model called the Digital Signature of Network Segment (DSNS), characterized by building profiles of normal behavior of network traffic generated by the tool Automatic Backbone Management (GBA), based on historical data from the network. The combination of techniques PSO and K-means, the algorithm called PSO-CIs was constructed, which is the basis of the SDA developed. To assess the quality of the developed system, several experiments were conducted from the perspective of different scenarios. Were also conducted studies regarding the complexity and optimization of the parameters of PSO-CIs. Were used in the experiments real data collected on the network of State University of Londrina. In order to compare the results obtained by the developed system were also implemented an ADS based on a deterministic algorithm, and one based on Principal Component Analysis (PCA). Finally, the results obtained from the experiments are presented and discussed.

Keywords: Anomaly detection. DSNS. PSO. K-means. PCA

LISTA DE FIGURAS

Figura 3.1 Estrutura funcional da ferramenta GBA versão 8	29
Figura 3.2 Tráfego e DSNS gerados pela ferramenta GBA	30
Figura 3.3 Estrutura do banco de dados Amostras	31
Figura 4.1 Passos para a detecção de anomalias utilizando o PSO-CIs	36
Figura 4.2 Evolução da função custo ao longo de 500 iterações, média em 500 realizações. Diferentes combinações de parâmetros φ_1 e φ_2 , $M =$ 20	38
Figura 4.3 Evolução da função custo ao longo de 500 iterações, média em 500 realizações. Diferentes tamanhos de população M	39
Figura 4.4 Média de operações matemáticas para diferentes números de objetos da MIB	41
Figura 4.5 Modelo de referência do sistema de detecção de anomalias	42
Figura 4.6 Alarmes gerados pelo sistema de detecção de anomalias	42
Figura 4.7 Desempenho do sistema de detecção para diferentes intervalos de análise	44
Figura 5.1 Passos para a detecção de anomalias utilizando o PCA	48
Figura 5.2 Exemplo da aplicação do PCA	49
Figura 5.3 Operações matemáticas para diferentes números de objetos da MIB	50
Figura 6.1 a) Conjunto de dados original b) Clusters e centróides iniciais	52
Figura 6.2 a) Iteração final do algoritmo PSO-CIs b) Convergência da função custo (eq.4.1)	52
Figura 6.3 Gráficos com o tráfego e DSNS utilizado nos testes com o objeto <i>ipInReceives</i> no período de 20 – 24/04/2010 no servidor web da UEL	54
Figura 6.4 Alarmes gerados pelo SDA desenvolvido no período de 20 – 26/04/2009, objeto <i>ipInReceives</i> e $\lambda = 1500$	55
Figura 6.5 Taxa de detecção \times taxa de alarmes falsos e $\lambda \times$ taxa de alarmes falsos, gerados pelo SDA desenvolvido, no cenário 2	56
Figura 6.6 Taxa de detecção \times taxa de alarmes falsos e $\delta \times$ taxa de alarmes falsos para SDA baseado no algoritmo determinístico, cenário 2	57
Figura 6.7 Tráfego e DSNS do dia 08/02/2010 do objeto <i>ifInOctets</i> do servidor web da UEL	58

Figura 6.8	Alarmes disparados pelo SDA para o objeto <i>ifInOctets</i> , no dia 08/02/2010	58
Figura 6.9	Resultados obtidos pelo SDA no dia 08/02/2010 no objeto <i>ifInOctets</i> do servidor web da UEL	59
Figura 6.10	Tráfego e DSNS dos dias 05–09/04/2010, do objeto <i>ipInReceives</i> do servidor web da UEL	60
Figura 6.11	Desempenho do SDA utilizando-se a classe 1, para definição de anomalias, cenário 3b	61
Figura 6.12	Desempenho do SDA utilizando-se a classe 2, para definição de anomalias, cenário 3b	62
Figura 6.13	Alarmes gerados para o período de 05 – 09/04/2010, do objeto <i>ipInReceives</i> no servidor web da UEL, cenário 3b	63
Figura 6.14	Gráficos da semana de tráfego do período de 05/04/2010 a 11/04/2010, objeto <i>ipInReceives</i> , cenário 4	64
Figura 6.15	Taxa de detecção × Taxa de alarmes falsos, para $\gamma = 5$, $\gamma = 15$ e $\gamma = 25$. Objeto <i>ipInReceives</i> . SDA PSO-CIs	65
Figura 6.16	Alarmes disparados pelos ADS baseado no PSO-CIs no período de 05/04/2010 a 11/04/2010, objeto <i>ipInReceives</i>	66
Figura 6.17	Taxa de detecção × Taxa de alarmes falsos, para $\gamma = 5$, $\gamma = 15$ e $\gamma = 25$. Quatro objetos SNMP simultâneos. SDA baseado no algoritmo PSO-CIs	67
Figura 6.18	Variância percentual explicada pelos componentes principais	68
Figura 6.19	Alarmes gerados pelo SDA baseado no algoritmo PSO-CIs, no período de 05/04/2010 a 11/04/2010, para os objetos <i>ipInReceives</i> , <i>ipInDelivers</i> , <i>tcpInSegs</i> e <i>ifInOctets</i>	69
Figura 6.20	Alarmes gerados pelo SDA baseado no PCA	70

LISTA DE TABELAS

Tabela 3	Número de operações por iteração	39
Tabela 4	Cenários de teste.....	51
Tabela 5	Parâmetros utilizados no algoritmo PSO-Cls	53

LISTA DE SIGLAS E ABREVIATURAS

BLGBA	Baseline para gerenciamento de backbone automático
DSNS	Digital Signature of Network Segment
GBA	Gerenciamento de Backbone Automático
MIB	Management Information Base
PCA	Principal Component Analysis
PDU	Packet Data Unit
PSO	Particle Swarm Optimization
PSO-Cls	Particle Swarm Optimization with clustering
SDA	Sistema de detecção de anomalias
SNMP	Simple Management Protocol
TAF	Taxa de alarmes falsos
TD	Taxa de detecção

CONVENÇÕES E LISTA DE SÍMBOLOS

Na notação das equações, as seguintes convenções foram utilizadas:

- letras minúsculas em negrito expressam vetores, exemplo: **x** e **p**;
- letras maiúsculas em negrito expressam matrizes, exemplo: **X**;

Os seguintes símbolos serão utilizados:

Símbolo	Descrição
$\alpha\%$	Fator amplitude percentual de uma anomalia.
γ	Fator tempo, de duração de uma anomalia.
Λ	<i>Threshold</i> para a que uma amostra de tráfego seja classificada como anomalia.
σ	Desvio padrão
$var(\cdot)$ ou σ^2	Variância
μ	Média
ω	Peso da inércia no movimento as partículas
φ_1	Coefficiente de aceleração local
φ_2	Coefficiente de aceleração global
$v_i[n]$	Velocidade da partícula i na iteração n
$p_i[n]$	Posição da partícula i na iteração n
$U_i[n]$	Matriz diagonal com elementos randômicos e distribuição uniforme
I	Número máximo de iterações
V_{max}	Constante positiva, com o valor máximo permitido para $v_i[n]$
M	Tamanho da população de partículas
K	Número de clusters
S	Número de amostras de um conjunto
$(\cdot)^T$	Operador de transposição
a	número de componentes principais no subespaço
P	Matriz de projeção no subespaço dos componentes principais
E	Matriz residual
Λ	λ Autovalores

SUMÁRIO

1 INTRODUÇÃO	13
2 TRABALHOS RELACIONADOS	17
3 CARACTERIZAÇÃO DE TRÁFEGO: BLGBA E DSNS	28
3.1 GERENCIAMENTO DE BACKBONE AUTOMÁTICO (GBA)	29
4 SISTEMA DE DETECÇÃO DE ANOMALIAS UTILIZANDO A ASSINATURA DIGITAL DE SEGMENTO DE REDE	33
4.1 <i>K-MEANS</i> CLUSTERING E <i>PARTICLE SWARM OPTIMIZATION</i>	33
4.2 OTIMIZAÇÃO DOS PARÂMETROS DE ENTRADA	36
4.3 ANÁLISE DE COMPLEXIDADE	39
4.4 DESCRIÇÃO DO SISTEMA DE DETECÇÃO DE ANOMALIAS (SDA	41
4.5 AVALIAÇÃO DO TAMANHO DO INTERVALO DE ANÁLISE.	43
5 ANÁLISE DE COMPONENTES PRINCIPAIS (PCA)	45
5.1 ANÁLISE DE COMPLEXIDADE	50
6 RESULTADOS	51
6.1 CENÁRIO 1	51
6.2 CENÁRIO 2	53
6.3 CENÁRIO 3	57
6.4 CENÁRIO 4	60
6.5 CENÁRIO 5	62
7 CONCLUSÕES	71
7.1 CONTRIBUIÇÕES	72
REFERÊNCIAS	74

1 INTRODUÇÃO

Os avanços contínuos da tecnologia têm impulsionado a expansão de uma grande variedade de serviços sobre as redes de comunicação, como por exemplo a convergência de serviços (tráfego de voz, vídeo e dados), o comércio eletrônico, as redes sociais, correio eletrônico, serviços de busca, digitalização de serviços públicos, dentre muitos outros. Boa parte dos serviços disponibilizados através das redes necessitam de garantias de operação, pois consistem de tarefas críticas que não podem sofrer interrupção ou queda no desempenho, pois isso resultaria em um impacto direto na arrecadação e na qualidade do serviço prestado. Além disso, as anomalias podem ser resultado de ações mal-intencionadas que além de prejudicar o desempenho da rede podem ocasionar a quebra de segurança e integridade das informações. De modo geral, anomalias de rede prejudicam seu bom funcionamento e podem estar relacionadas ou não a algum tipo de ação maliciosa, bem como falhas físicas. Sendo assim, a detecção de anomalias consiste de um processo essencial na tarefa de manutenção do desempenho, garantia de segurança da informação e disponibilidade dos recursos e serviços.

A detecção de anomalias é um problema altamente relevante que tem sido tratado dentro de diversos domínios e áreas de pesquisa, tais como a estatística, aprendizagem de máquina, mineração de dados, teoria da informação, teoria espectral e remete ao problema de encontrar padrões em dados que não correspondam ao comportamento esperado [1][2]. Esses comportamentos fora do padrão são muitas vezes referidos como anomalias, *outliers*, alterações, exceções, anormalidades, surpresas, novidades, peculiaridades ou contaminantes de acordo com o domínio de aplicação.

Em redes de computadores, uma anomalia é classificada como um desviorepentino em relação ao modelo de funcionamento padrão da rede. Eventualmente anomalias de tráfego podem ser resultados de ações mal-intencionadas (por exemplo, ataques de negação de serviço), mas são principalmente originadas de falhas na interação entre os componentes do sistema, especialmente sob grande carga de processamento; fontes comuns de anomalias são: erros de software (por exemplo, estouro de memória, consultas de banco de dados mal formuladas), mau funcionamento de hardware (por exemplo, falhas de disco) ou falhas nos18 subsistemas (por exemplo, link de acesso da rede de comunicação quebrado). Além disso, os operadores da rede

devem ser informados em tempo real quando da ocorrência de um problema, para que ele possa ser analisado e as medidas adequadas possam ser tomadas.

O desempenho de uma solução para a tarefa de detecção de anomalias é avaliada em termos de precisão de detecção considerando alarmes falsos e anomalias detectadas. Segundo Hodge [3], existem três abordagens principais utilizadas na solução do problema de detecção de anomalias:

Técnicas supervisionadas, que modelam os comportamentos normais e anormais, e exigem dados rotulados. Esta é a abordagem tradicional.

Técnicas semi-supervisionadas, que modelam apenas o comportamento normal, e são mais aplicáveis do que a técnica anterior, uma vez que apenas os rótulos para os dados normais são necessários.

Técnicas não supervisionadas, que não exigem nenhum conhecimento prévio dos dados, mas devem assumir que as anomalias são de algum modo distante dos dados normais.

Em algumas pesquisas como as de Patcha [1] e Chandola [2], observa-se que a detecção de anomalias é tradicionalmente tratada com a utilização de dados históricos, e que a maioria das técnicas exigem todo o conjunto de dados de teste, antes de se executar a detecção de anomalias, e menciona poucas técnicas que funcionam em tempo real. De fato, os algoritmos mais atuais assumem que a memória principal possui capacidade de armazenar todo o conjunto de dados [4]. Além disso, grande parte das abordagens tem seu foco direcionado para a detecção de intrusões [5].

A detecção de anomalias em grandes sistemas dinâmicos é uma tarefa desafiadora, uma vez que eventos anômalos podem aparecer raramente, e eventualmente podem não possuir assinaturas fixas. A alta dimensionalidade dos dados de tráfego ocorre devido à complexidade dos sistemas e das redes, em conjunto com a diversidade de fontes e freqüentes mudanças de comportamento normal do sistema, resultante de comportamento do usuário. Isso faz com que a detecção de anomalias se torne uma tarefa ainda mais difícil. Apesar de não ser uma área de estudos recente, o problema de detecção de anomalias ainda encontra-se em aberto no sentido de que as técnicas existentes não abordam todos os aspectos que abrangem a detecção de anomalias.

Neste trabalho foi realizado um estudo sobre as técnicas utilizadas na detecção de anomalias, a fim de criar um embasamento para o desenvolvimento de um Sistema de Detecção de Anomalias (SDA). Devido a grande quantidade de

técnicas existentes¹⁹ e a dificuldade de identificar com rapidez que técnica aplicar nas situações de anomalias em redes com grande volume de tráfego, decidiu-se implementar um SDA baseado em um modelo heurístico, que utiliza tráfego coletado de objetos pertencentes a MIB-II [6] através do protocolo SNMP na rede da Universidade Estadual de Londrina. Esse processo envolve a extração de informação relevante de conjuntos de dados históricos, multidimensionais e ruidosos. Da revisão da literatura realizada, destacaram-se dois modelos heurísticos que apresentaram resultados promissores na detecção de anomalias: o K-means *clustering* e o *Particle Swarm Optimization* (PSO). Deste modo, foi proposto em nosso trabalho a utilização de ambos algoritmos para o desenvolvimento de um sistema de detecção de anomalias (SDA), utilizando também os conceitos criados nos estudos de Proença [7], o qual desenvolveu uma ferramenta para gerar perfis de comportamento normal de segmentos de rede, denominada gerenciamento de backbone automático (GBA). Esses perfis de comportamento são chamados de Assinatura Digital de Segmento de Rede (DSNS) e são parte importante na caracterização do tráfego, como pode ser visto no capítulo 3.

O K-means é um algoritmo para classificação de dados, utilizado em diferentes áreas, principalmente no processamento de imagens. Já demonstrou ser promissor na detecção de anomalias com excelentes resultados [8]. Apesar de ser um algoritmo simples, sofre com a falta de mecanismos para escapar dos ótimos locais, que pode ser superada através da combinação com o PSO. Em contrapartida, o PSO é uma técnica heurística altamente eficiente, com baixa complexidade computacional e capacidade de escapar de ótimos locais. Foi desenvolvido em 1995 por Kennedy e Eberhart [9], cuja idéia surgiu após ter analisado o comportamento das aves e dos cardumes de peixes, constatando que o comportamento de vida em grupo poderia ser explorado como uma ferramenta de busca heurística.

O SDA desenvolvido é baseado na junção dos algoritmos PSO e K-means, denominado PSO-CIs, combinados ao DSNS. Seu objetivo é avaliar a proximidade entre as amostras coletadas dos segmentos monitorados e os centróides dos clusters calculados. A detecção é realizada através da comparação das distâncias entre amostras e centróides, conforme detalhado no capítulo 6.

Essa dissertação é organizada da seguinte forma. O Capítulo 2 apresenta os trabalhos relacionados ao desenvolvimento da nossa pesquisa. No Capítulo 3 é apresentado o modelo de caracterização de tráfego utilizado no SDA desenvolvido, enquanto no Capítulo 4 são apresentados os conceitos utilizados em nosso SDA, bem

como um estudo do algoritmo utilizado, e a descrição do funcionamento do sistema. No capítulo 5, são apresentados os conceitos da técnica PCA, uma abordagem que tem se mostrado promissora para a detecção²⁰ de anomalias. Os resultados obtidos através dos experimentos sob a perspectiva de diferentes cenários, são apresentados no Capítulo 6. Finalmente no Capítulo 7, temos a conclusão do trabalho.

2 TRABALHOS RELACIONADOS

As anomalias de rede são classificadas como eventos inesperados no tráfego, que não coincidem com um comportamento normal previsto. Estes eventos são muitas vezes referidos como anomalias, *outliers*, alterações, exceções, anormalidades, dentre outros nomes, variando de acordo com o domínio da aplicação. Esses eventos podem ser originados de uma variedade de situações, tais como falhas físicas em equipamentos da rede ou ataques, que levam a uma mudança abrupta e imprevisível no fluxo de dados, comprometendo o desempenho da rede. Com base nessa definição, observa-se uma crescente sofisticação das ameaças nas redes de comunicação, que fizeram com que a detecção de anomalias se tornasse um tema de grande importância para a manutenção do desempenho e segurança das grandes redes, e conseqüentemente, foco de muitas pesquisas na área acadêmica [10] [11] dentro de diversos domínios, tais como estatística, aprendizagem de máquina, mineração de dados, teoria da informação e teoria espectral [1][2].

Embora utilizando abordagens distintas, grande parte das pesquisas desenvolvidas na área de detecção de anomalias possui o mesmo objetivo, que é o de otimizar o *trade-off* entre a taxa de detecção e a taxa de falsos positivos, ou seja, otimizando a detecção de anomalias. Outro ponto em comum, é a dificuldade em se estabelecer um modelo normal para o tráfego, que descreva qualquer tipo de comportamento de rede. Isso ocorre devido às características heterogêneas das redes, e também ao grande fluxo de tráfego agregado, que fazem com que o padrão de comportamento da rede, se modifique constantemente. Esses aspectos são responsáveis por manter a área de detecção de anomalias com grande foco para o desenvolvimento de novas pesquisas.

Diversos trabalhos como [12, 13, 14, 8, 15] por exemplo, abordam a clusterização como estratégia para a detecção de intrusão e anomalias. Isso ocorre devido à dificuldade em classificar e extrair conhecimento de grandes volumes de dados, proveniente das grandes redes e fluxos agregados. Deste modo, a utilização de heurísticas vem ganhando espaço nas pesquisas acadêmicas, sendo combinada com outras técnicas com o objetivo de reduzir a complexidade proveniente do grande volume de informações, e otimizar parâmetros que afetam a qualidade da solução, como visto nos trabalhos [8, 15, 16, 17, 18, 19]. Recentemente, diversos trabalhos têm utilizado abordagens clássicas, para a detecção de anomalias. Uma dessas técnicas é a

análise de componentes principais (PCA), utilizada para redução de dimensionalidade em grandes conjuntos de dados, a qual tem se mostrado bastante promissora na detecção de anomalias como visto nos trabalhos [20, 21, 22, 23, 24]. Apesar de boa parte dos trabalhos desenvolvidos na área de detecção de anomalias utilizarem abordagens orientadas a fluxos IP, a utilização do protocolo SNMP ainda é muito atrativa, uma vez que continua sendo o padrão para a gerência de redes, e está presente em quase todos os equipamentos de redes comerciais.

Xiao *et al.* [8], propôs a combinação dos algoritmos K-means e *Particle Swarm Optimization* (PSO) para detecção de anomalias de rede. O K-means consiste de uma ferramenta para clusterização de dados muito útil, e que provou ser muito promissora como técnica para detecção de anomalias. Contudo, assim como o método subida de encosta [8], se os padrões de entrada do algoritmo não forem definidos corretamente, o método pode não convergir ou vir a convergir para um ótimo local. Para suprimir essa fraqueza, o K-means é combinado ao PSO, que consiste de uma ferramenta heurística para otimização numérica baseada no comportamento social em grupo das aves e peixes. É uma ferramenta relativamente simples, de rápida convergência e com boa capacidade de pesquisa global, muito utilizada para solução de problemas com muitas amostras e espaço de busca multidimensional. O modelo proposto é avaliado utilizando o KDD Cup 1999 *dataset* [25], que consiste de um conjunto de dados, com uma grande variedade de anomalias simuladas em um ambiente de rede militar. Cada amostra no conjunto de dados consiste de um registro de recursos extraídos de uma conexão de rede TCP/IP durante a simulação das anomalias. O primeiro passo para a detecção de anomalias é a seleção dos campos dos registros através de um classificador Bayesiano. Após essa etapa, o conjunto de dados de treinamento é normalizado, e aplica-se os algoritmos K-means e PSO para clusterizar os dados e calcular os centróides. A segunda etapa, após o treinamento dos dados, é avaliar o conjunto de dados a ser testado e verificar sua distância em relação ao conjunto de dados de treinamento. Se a distância exceder a um *threshold*, a amostra é classificada como anômala. Os resultados obtidos mostraram-se bastante satisfatórios, alcançando uma taxa de detecção de 86% × 2, 8%. O sistema de detecção de anomalias (SDA) que nós desenvolvemos, utiliza a mesma combinação de algoritmos para classificar os dados de tráfego. Além desses dois algoritmos, utilizamos também um modelo para caracterização de tráfego denominado *Digital Signature of Network Segment* (DSNS), o que poupa o algoritmo de realizar a etapa de treinamento

dos dados. O algoritmo é utilizado para clusterizar não somente o conjunto de dados, mas também o DSNS, de forma²³ que a detecção de anomalias é realizada através da comparação da distância entre as amostras de um cluster e seu respectivo centróide. Além disso, este modelo é capaz de detectar desvios abaixo e acima de um limiar previamente calculado, o DSNS. Os resultados que obtivemos mostraram-se bastante promissores, ficando muito próximos aos obtidos pelo modelo proposto por Xiao, com a vantagem de utilizarmos dados coletados em um ambiente de rede real, além de desenvolver um mecanismo de detecção de anomalias mais elaborado.

Ensaifi *et al.* [15], também propôs a utilização do algoritmo K-means para a detecção de anomalias de rede. A fim de solucionar o problema de convergência para ótimos locais e a elevada taxa de falsos alarmes provenientes do K-means, os autores propuseram a combinação de duas técnicas de *soft computing*: a lógica *Fuzzy* e o *Particle Swarm Optimization* (PSO). Deste modo, o algoritmo K-means é combinado com a teoria *Fuzzy*, a fim de otimizar a clusterização. Os algoritmos de clusterização *Fuzzy*, consideram cada cluster como um conjunto *Fuzzy*, enquanto uma função de medida de pertinência mensura a possibilidade de cada vetor de treinamento pertencer a um cluster. Deste modo, os clusters serão ótimos no sentido de que a variância dentro do cluster será tão pequena quanto possível. Isso significa que todos os objetos possuem atributos quase iguais, o que significa uma alta densidade e uma pequena distância entre eles no espaço de atributos. O PSO por sua vez é utilizado para otimizar o cálculo da função custo, originada da combinação entre K-means e lógica *Fuzzy*. O método proposto consiste de duas fases: (i) a fase de treinamento, na qual a melhor das partículas é encontrada ao longo das gerações, representando os conjuntos de clusters, e (ii) a fase de detecção, que utiliza a distância euclidiana entre os centróides dos clusters e os dados de entrada para verificar se um ponto de dado de tráfego é normal ou anormal. Os experimentos foram realizados utilizando o KDD Cup 1999 *dataset* [25]. Os resultados indicaram uma alta capacidade para a detecção de anomalias no conjunto de testes utilizado, apesar de apresentar uma elevada taxa de falsos alarmes. Apesar dos ótimos resultados, os autores não apresentam nenhum estudo a respeito do desempenho e complexidade do modelo desenvolvido. Em nosso trabalho realizamos um estudo a respeito desses parâmetros e também um comparativo com o modelo PCA.

No trabalho de Liu [16], é proposto a utilização de uma versão modificada do algoritmo *Quantum-behaved Particle Swarm Optimization* (QPSO) para

a detecção de anomalias, o MQPSO. O QPSO, consiste de uma abordagem do PSO, que trata o algoritmo sob o ponto de vista da mecânica quântica, assumindo que as partículas possuem um comportamento quântico. No espaço quântico, a natureza do estado das partículas agregadas é completamente diferente, de modo que as partículas podem procurar a solução em todo o espaço de busca possível, fazendo com que o desempenho do QPSO seja superior ao PSO24 convencional. A versão modificada do QPSO proposta pelo autor, sugere uma alteração nos limites do espaço de busca do QPSO, reduzindo drasticamente a possibilidade das partículas ficarem presas a um ótimo local, caso os limites do espaço de busca consistam em um ótimo local. O MQPSO é utilizado para treinar os parâmetros de uma rede neural *wavelet* (WNN), que é utilizada para a detecção de anomalias. Um vetor multidimensional composto pelos parâmetros da WNN é associado a uma partícula no algoritmo PSO. A combinação dos parâmetros adequados da rede neural, determina a viabilidade do espaço de busca para obter a solução ótima. A fim de validar a abordagem proposta, foram realizados experimentos utilizando o KDD Cup 1999 *dataset* [25]. Os resultados demonstraram que o algoritmo proposto obteve um desempenho superior na fase de treinamento, uma convergência mais rápida, além de melhor capacidade de detecção de ataques de tipos desconhecidos, em relação ao QPSO. Todavia, a abordagem utilizada pelo autor consiste de uma técnica supervisionada e de alto custo computacional, uma vez o tráfego normal e anormal devem ser modelados. No sistema que nós desenvolvemos, é utilizado uma abordagem semi-supervisionada, o que representa uma grande vantagem, uma vez que somente o tráfego normal necessita ser modelado.

Em [17], Ma *et al.* propôs uma função de rede neural com base radial (RBFNN) para detecção de anomalias de rede. O método consiste em inserir uma função de base radial em uma rede neural de duas camadas *feed-forward*. Para a utilização de redes neurais artificiais na detecção de anomalias, é necessário escolher um algoritmo de treinamento robusto, a fim de se encontrar os melhores parâmetros da rede neural. Dessa forma, uma técnica de inteligência coletiva denominada *Quantum-behaved Particle Swarm Optimization* (QPSO), foi utilizada pelos autores para o treinamento da rede neural RBF. O QPSO consiste de uma abordagem diferente do PSO tradicional, onde o movimento das partículas é assumido como sendo quântico. Nas primeiras gerações do QPSO, as partículas no espaço multidimensional, buscam pela solução de forma tão rápida que correm o risco de ficarem estagnadas, ao passo que quanto mais gerações são produzidas, mais lentamente

evolui a precisão das partículas, fazendo com que exista a possibilidade de não alcançarem um ótimo global. A fim de contornar essa situação o QPSO é associado com a técnica denominada Gradiente Descendente (GD). A fim de validar a abordagem proposta, os experimentos foram conduzidos utilizando o KDD Cup 1999 *dataset* [25]. Os resultados demonstraram que o algoritmo híbrido QPSO-GD obteve melhores resultados na detecção sobre o conjunto de dados utilizados, do que o QPSO convencional.

No trabalho desenvolvido por Gao [18], é proposto um método de otimização baseado no algoritmo PSO, para detectores do tipo *Negative Selection Algorithm* (NSA).²⁵ A detecção de anomalias é realizada através do NSA, que é inspirado na ideia do sistema imunológico natural, onde as células de defesa denominadas linfócitos, atacam e destroem as células que não pertencem ao corpo. Para que isso ocorra, é necessário que as células de defesa passem por um procedimento denominado *negative selection*, que impede que as células pertencentes ao corpo sejam atacadas e destruídas. A modelagem para detecção de anomalias, se dá através da definição de um conjunto de dados normais, que correspondem às células do corpo, de modo que os detectores NSA não irão classificar os novos pontos de dados que combinam com esse conjunto, como anômalos. O ideal é que os detectores cubram ao máximo o espaço de dados não pertencente ao 'corpo', ao mesmo tempo em que apresentem o mínimo de sobreposição entre eles mesmos. Deste modo o PSO é utilizado para encontrar o melhor conjunto de detectores, que atendam a essas exigências, de modo que cada partícula é considerada um detector candidato a ser otimizado. As simulações demonstraram resultados promissores, alcançando taxas de detecção superiores à 80%, porém, os alarmes falsos não foram apresentados no trabalho. Também não foram apresentados resultados comparativos com outras técnicas, a fim de avaliar o desempenho do modelo proposto. No trabalho que nós desenvolvemos a precisão da detecção de anomalias é avaliada levando em consideração o *trade-off* entre taxa de detecção e alarmes falsos, juntamente com um estudo a respeito do desempenho computacional.

Jianzhen e Jinrong [19] realizaram um estudo sobre a aplicação do algoritmo *Particle Swarm Optimization* (PSO) para detecção de anomalias. São avaliados três métodos combinados com a utilização do PSO: clusterização, redes neurais e *Support Vector Machines* (SVM). São apresentadas as vantagens e desvantagens de cada abordagem. O algoritmo PSO combinado com o algoritmo de clusterização, embora

melhore a qualidade da população e do desempenho global da pesquisa, apresenta uma detecção insatisfatória quando a diferença entre comportamento normal e anômalo é mínima, como por exemplo, ataques originados de usuários legítimos. O algoritmo PSO combinado com redes neurais possui a capacidade de otimizar os pesos da rede neural, *thresholds* e demais parâmetros, porém é necessário identificar a estrutura da rede com antecedência. Além disso, os dados de treinamento possuem grande impacto sobre a qualidade da rede neural, tornando o algoritmo instável em alguns casos. A otimização dos parâmetros da função cerne do SVM através do PSO, pode melhorar a capacidade de generalização, aumentando conseqüentemente a precisão da classificação. Como o desempenho da classificação do SVM é dependente das características das amostras de treinamento, se o algoritmo PSO for utilizado para otimizar essas características e os parâmetros da função cerne, o algoritmo vai obter um melhor desempenho de detecção. De modo geral, os autores identificaram que as abordagens de detecção de anomalias baseadas no algoritmo PSO, quase sempre estão combinadas com a utilização de outros métodos, sendo necessário realizar uma análise profunda das características do problema de detecção, a fim de verificar sua aplicabilidade com o modelo do PSO. Em nosso trabalho, combinamos o PSO com o algoritmo de clusterização K-means, a fim de realizar a detecção de anomalias, através da classificação dos dados. O objetivo do PSO é otimizar a função custo do K-means, responsável pela classificação dos dados.

Gaddam [12], propõe a utilização do K-means combinado com a técnica *ID3 decision tree learning* para a detecção de anomalias. A detecção é realizada em duas etapas: no primeiro estágio, o K-means é aplicado em um conjunto de dados de treinamento para obter k clusters disjuntos. Cada cluster representa uma região de instâncias similares, em termos de distância euclidiana entre as instâncias e os centróides. Os autores escolheram o K-means devido a dois aspectos principais: 1) é um método orientado a dados, com relativamente poucas hipóteses sobre as distribuições dos dados subjacentes e 2) a estratégia de busca gulosa do K-means, garante ao menos um mínimo local da função custo, acelerando a convergência dos clusters em grandes conjuntos de dados. Na segunda etapa, o método K-means é utilizado em cascata com a técnica *ID3 decision tree learning*, através da construção de uma árvore de decisão ID3, utilizando as instâncias em cada cluster gerado pelo K-means. Deste modo, consegue-se contornar dois problemas inerentes ao K-means: 1) o problema da atribuição forçada, que ocorre quando o K-means é inicializado com um valor de k , que subestima o

número natural de grupos presente nos dados de treinamento, fazendo com que instâncias de diferentes grupos sejam forçadas a pertencer ao mesmo cluster, provocando um aumento na taxa de alarmes falsos, e 2) o problema da dominância de classe, que surge em um cluster quando os dados de treinamento possuem um grande número de instâncias pertencentes a uma classe, e poucas instâncias pertencentes às classes restantes. Isso acarreta em um erro, ao classificar uma anomalia associada com um conjunto dominado por instâncias normais ou vice-versa. O modelo foi testado utilizando três conjuntos de dados do *DARPA intrusion detection data sets* [26]. Todos os resultados demonstram um ganho de desempenho bastante significativo, do modelo K-means-ID3 em relação ao K-means tradicional. Embora o modelo proposto pelos autores resolva os problemas da atribuição forçada e o problema da dominância de classe, a deficiência do modelo em convergir para ótimos locais não foi abordada. É com objetivo de contornar esse problema, que nós combinamos o PSO ao K-means. Os resultados que obtivemos demonstram a eficiência do modelo em encontrar soluções globais.

No trabalho de Zhang [13], é proposto uma nova versão do algoritmo de clusterização K-means, para a detecção de intrusões em redes. Os autores realizaram um estudo das vantagens e desvantagens das abordagens por detecção de anomalias e detecção baseada em assinaturas (*misuse*), e propuseram um sistema misto. Primeiramente os dados são analisados pelo módulo de detecção de mau uso (*misuse detection*), de modo que os dados suspeitos são enviados ao módulo de detecção de anomalias. Esse módulo por sua vez, é constituído pelo algoritmo de clusterização não supervisionado KD, que consiste de uma versão modificada do K-means que introduz um parâmetro para definir o raio dos clusters a serem calculados. Após a clusterização dos dados de treinamento, os clusters são classificados. Como os dados normais constitui a maioria do conjunto de treinamento, o cluster que apresentar um número de registros maior do que um *threshold*, é marcado como normal. Para a detecção de anomalias em um novo registro d , são calculadas as distâncias com cada centróide, de modo que irá pertencer ao cluster com centróide mais próximo. Se esse cluster for normal, d é classificado como normal, caso contrário é classificado como anômalo. Os testes foram realizados utilizando o KDD Cup 1999 *dataset* [25], e os resultados foram bastante satisfatórios. O algoritmo K-means convencional, obteve melhores resultados para intrusões individuais, em contrapartida o algoritmo KD demonstrou um ganho significativo de precisão para intrusões mistas.

Jianliang [14], propõe a utilização do algoritmo K-means convencional para a detecção de anomalias. O K-means consiste de uma técnica de clusterização muito útil, classificada como uma técnica de aprendizado competitivo, a qual provou-se através de vários trabalhos, ser uma técnica promissora para a detecção de anomalias. A técnica consiste basicamente em dividir um conjunto de dados com n amostras, em um conjunto de k clusters, de modo que a detecção de anomalias ocorre através da avaliação da distância entre novas amostras, e os centróides. Os experimentos foram realizados utilizando o KDD Cup 1999 *dataset* [25]. Os resultados dos testes demonstraram que o método é muito eficaz para o particionamento de grandes conjuntos de dados. O trabalho de Jianliang não aborda um aspecto básico que é a tendência do K-means em encontrar soluções locais. No DAS que desenvolvemos, utilizamos a técnica de otimização numérica denominada *Particle Swarm Optimization* (PSO) para suprimir essa deficiência.

No trabalho de Brauckhoff [20], os autores aprofundam no estudo da aplicação da análise de componentes principais (PCA) para a detecção de anomalias, com o objetivo de preencher algumas lacunas presentes nos trabalhos relacionados que mostram que o PCA é bastante sensível a ajustes de calibração. O desempenho do PCA aplicado a detecção de anomalias é avaliado utilizando um conjunto de dados composto por três semanas de tráfego NetFlow, oriundos de um ISP de médio porte, que compreende uma grande variedade de anomalias de tráfego tais como *network scans*, *denial of service attacks*, *alpha flows*. Os autores demonstraram que o problema do PCA é que ele falha em capturar a correlação temporal. Deste modo é proposto uma solução para lidar com esse problema, substituindo o PCA pela transformada de Karhunen-Loeve. Os autores comprovaram que quando considerado a correlação temporal, os resultados apresentaram uma melhora significativa utilizando a abordagem proposta. Em nosso trabalho foi apresentado um sistema de detecção de anomalias que aplica a abordagem convencional do PCA, sobre conjuntos de dados coletados da MIB através do protocolo SNMP. Os resultados obtidos foram bastante satisfatórios considerando o *trade-off* entre taxa de detecção e alarmes falsos, porém inferiores aos obtidos pela abordagem heurística que utilizamos.

No trabalho desenvolvido por Callegari *et. al* [21], é apresentado uma nova técnica que estende o estado da arte em detecção de anomalias baseado na análise de componentes principais (PCA). Através da análise multi-escala, o modelo proposto é capaz de obter grandes melhoras em relação ao desempenho da abordagem

convencional. Além disso também é apresentado um método para identificar os fluxos responsáveis por uma anomalia detectada em nível agregado. Os dados utilizados para avaliar o sistema são do *data set* público do *Abilene Observatory* [27]. A análise do desempenho demonstrou que para uma escolha adequada dos parâmetros de ajuste, o sistema implementado obtém resultados bastante satisfatórios, sendo capaz de detectar todas as anomalias sintéticas e também as anomalias já presentes no conjunto de dados original, além de ser capaz de identificar os fluxos responsáveis pela anomalia.

Casas *et. al.* [22], introduz um algoritmo ótimo para detecção de anomalias de volume, em que o tráfego livre de anomalias é tratado como uma característica indesejada. O algoritmo é ótimo no sentido de que maximiza a probabilidade de detecção, com uma taxa de alarmes falsos delimitada. Para superar os problemas de estabilidade das abordagens anteriores proposta pelo autor, um novo modelo linear não orientado a tráfego é proposto. Este modelo mantém-se estável no tempo e torna o processo de demanda de tráfego observável a partir de medições SNMP. O modelo pode ser usado de duas maneiras, (i) realizar a estimativa dos volumes de fluxos OD livres de anomalias, ou (ii) eliminar o tráfego livre de anomalias das medições SNMP a fim de proporcionar resíduos sensíveis à anomalias. O desempenho do método proposto é comparado ao obtido com a abordagem de análise de componentes principais (PCA), a qual foi escolhida como referência dada a sua relevância na literatura para detecção de anomalias. Utilizando os dados reais de tráfego do *backbone* da rede Abilene Internet2 [27], foi apresentada uma comparação empírica entre o algoritmo de detecção de anomalias proposto e o PCA. Através desta análise os autores verificaram os parâmetros ótimos do algoritmo proposto, sua estabilidade em relação ao modelo de tráfego e também como o método supera a abordagem PCA no conjunto de dados de teste considerado

Lee *et. al.* [23], apresenta um algoritmo robusto para coleta de tráfego através do protocolo SNMP, capaz de coletar os dados de maneira frequente, sem que ocorra degradação de performance. Isso é realizado através da correlação dos objetos da MIB. O algoritmo utiliza os objetos do grupo IP ao invés dos objetos do grupo interface da MIB, ao coletar o tráfego dos equipamentos. Como os dados não são coletados diretamente da interface, não ocorre uma degradação de performance nos equipamentos da rede ou nos servidores de coleta. O tráfego do grupo interface é estimado através dos objetos do grupo IP. Os objetos utilizados para se realizar a estimativa são: *ipInReceives*, *ipOutRequest*, *ipOutNoRoutes*, *ipFragCreates*,

ipForwDatagrams, *ipFragOKs*, *ipOutDiscards* e *ipFragFails*. O algoritmo é dividido em duas etapas. Na primeira etapa, são filtrados os equipamentos impróprios que não poderão ser utilizados na segunda etapa, isso ocorre porque nem todos os equipamentos possuem informações suficientes para estimar o tráfego dos objetos do grupo interface. Na segunda etapa, são coletados os valores dos objetos do grupo IP da MIB, dos equipamentos selecionados na primeira etapa, e também realizadas as estimativas dos valores dos objetos do grupo interface, a fim de gerar alarmes. A detecção de anomalias se dá através da definição de *thresholds* para o grupo de objetos interface, que podem ser calculados de duas formas: 1) *accumulative threshold*, que pode ser calculado automaticamente utilizando as médias dos tráfegos anteriores e por valores de tolerância definidos pelo usuário para cada interface. O *threshold* pode ser calculado multiplicando o valor da média pelo valor da tolerância definido pelo usuário; 2) um *threshold* constante para cada interface, definido pelos operadores de rede. Foram realizados experimentos, que avaliaram o tráfego IP de 74 equipamentos, incluindo roteadores e switches no *backbone* da rede KORNET. Os resultados do experimento demonstram que a precisão da estimativa do algoritmo é boa, e pode ser melhorada se os *thresholds* forem definidos corretamente. A detecção de anomalias através de simples *thresholds* como proposto pelo autor, não é de grande eficiência, uma vez que o tráfego pode apresentar diferentes ciclos de comportamento, que não serão bem representados por um único *threshold*. É por isso que em nosso trabalho, utilizamos o modelo *Digital Signature of Network Segment* (DSNS), proposto por Proença em [28] para caracterizar o tráfego. Esse modelo leva em consideração a natureza não-estacionária do tráfego, e o fato de que os níveis de tráfego são geralmente mais elevados durante certos períodos de tempo, diferenciando também entre dias úteis e finais de semana.

No trabalho desenvolvido por Zarpelão [24], é proposto um sistema de detecção de anomalias baseado em perfis de comportamento normais, utilizando dados de tráfego da MIB, coletados através do protocolo SNMP. O modelo proposto, aplica algoritmos determinísticos parametrizáveis, sobre os conjuntos de dados coletados dos objetos da MIB e 30 seus respectivos perfis de comportamento normal denominados *Digital Signature of Network Segment* (DSNS), gerados pela ferramenta Gerenciamento de Backbone Automático (GBA), proposta por Proença em [7]. Um sistema de alarmes informa sobre desvios detectados através da comparação entre os DSNS's e o movimento real representado pelos objetos SNMP. O sistema de alarmes é baseado em

3 eventos, que são responsáveis por caracterizar desvios de comportamento significativos, e também por evitar a geração de uma grande quantidade de alarmes falsos. Os alarmes gerados pelos objetos monitorados, são enviados para o sistema de correlação que é responsável por reunir todos alarmes e verificar a ocorrência de anomalias. A detecção de anomalias se dá através da verificação da correlação dos alarmes no grafo de dependência entre objetos SNMP, através de um algoritmo de busca em profundidade. Dois vértices são considerados correlacionados quando eles são adjacentes e há alarmes gerados para ambos os objetos no mesmo *frame* de cinco minutos. Experimentos foram realizados utilizando tráfego coletado na rede da Universidade Estadual de Londrina. Os resultados obtidos foram bastante satisfatórios, ficando muito próximos do que foi estabelecido como objetivo. O modelo proposto por Zarpelão serviu como base para a nossa pesquisa, e embora tenha alcançado ótimos resultados, o modelo não é capaz de detectar anomalias que representam desvios do tráfego que ficam abaixo do DSNS. O sistema de detecção de anomalias que nós desenvolvemos, também utiliza tráfego coletado através do protocolo SNMP e o DSNS, sendo capaz de detectar desvios de tráfego abaixo e acima do DSNS e com resultados superiores aos apresentados no trabalho desenvolvido por Zarpelão, em relação ao *trade-off* entre taxa de detecção e taxa de alarmes falsos. O critério de optimalidade adotado neste trabalho é uma taxa de detecção de pelo menos 80% e uma taxa de alarmes falsos de até 10%

3 CARACTERIZAÇÃO DE TRÁFEGO: BLGBA E DSNS

O primeiro passo para se realizar a detecção de anomalias em redes é a criação ou adoção de um modelo que caracterize o tráfego de rede de forma eficiente, o que representa um grande desafio devido à natureza não-estacionária do tráfego. O comportamento do tráfego de grandes redes é composto por ciclos diários, onde os níveis de tráfego são geralmente mais elevados durante o horário comercial, diferenciando-se também entre os dias úteis e finais de semana. Deste modo, um modelo eficiente de caracterização do tráfego deve ser capaz de representar de forma fiel essas características. Deste modo, neste trabalho a ferramenta GBA é utilizada para gerar diferentes perfis de comportamento normal para cada dia da semana, de acordo com essas exigências. Estes perfis de comportamento são chamados de Assinatura Digital de Segmento de rede (*Digital Signature of Network Segment, DSNS*), proposto por Proença em [28]. O DSNS pode ser definido como um conjunto básico de informações que constituem o perfil de tráfego de um segmento ou equipamento de rede. Essas informações incluem dados como tráfego de volume, número de erros, tipos de protocolos e serviços que são transportados ao longo do segmento durante o dia.

O algoritmo BLGBA, responsável por gerar o DSNS, foi desenvolvido com base em uma variação no cálculo da moda estatística. A fim de determinar um valor esperado para um determinado segundo do dia, o modelo analisa os valores para o mesmo segundo nas semanas anteriores. Estes valores são distribuídos em frequências, com base na diferença entre o maior G e o menor elemento S da amostra, utilizando-se 5 classes. Essa diferença, dividida por cinco, constitui a amplitude h entre as classes, $h = G - S/5$. Em seguida, os limites de cada classe L_{Ck} são obtidos. Eles são calculados por $L_{Ck} = S + h \cdot k$, onde C_k representa a classe k ($k = 1 \dots 5$). O valor com maior quantidade de elementos inserido na classe com frequência acumulada igual ou superior a 80% é incluído no DSNS. As amostras para a geração do DSNS são coletadas segundo a segundo, ao longo do dia, pela ferramenta GBA. Dois tipos de DSNS são gerados: o BL-7 composto por um DSNS para cada dia da semana, e o BL-3 que consiste de um DSNS para os dias úteis, um para o sábado e outro para o domingo. O período ideal de dados históricos para a geração do DSNS pode variar de 4 a 12 semanas, de modo que sua geração é realizada segundo a segundo do dia, todos os dias da semana.

A figura 3.1 apresenta gráficos de uma semana de tráfego e DSNS, gerados pela ferramenta GBA. Os dados foram coletados do objeto *ifInOctets* que pertence ao grupo *interface* da MIB-II [6], do servidor Web Universidade Estadual de Londrina (UEL). O movimento é representado em verde e o respectivo DSNS pela linha azul. A linha vermelha representa que o movimento excedeu o DSNS, mas não implica necessariamente que ocorreu uma anomalia naquele instante. Na figura 3.1 é possível observar um grande ajuste do comportamento do tráfego monitorado em relação ao DSNS.

3.1 GERENCIAMENTO DE BACKBONE AUTOMÁTICO (GBA)

O Gerenciamento de Backbone Automático (GBA), consiste de uma ferramenta com o propósito de auxiliar o gerenciamento de rede, através do monitoramento de objetos da MIB coletados pelo protocolo SNMP. A ferramenta GBA é utilizada para gerar diferentes perfis de comportamento normal, para cada dia da semana, dos equipamentos ou segmento da rede. Sua versão inicial foi desenvolvida em 1998 por Proença [29], e atualmente encontra-se na versão 8.0. A figura 3.2, apresenta os módulos presentes na versão atual da ferramenta.

Figura 3.1 Estrutura funcional da ferramenta GBA versão 8

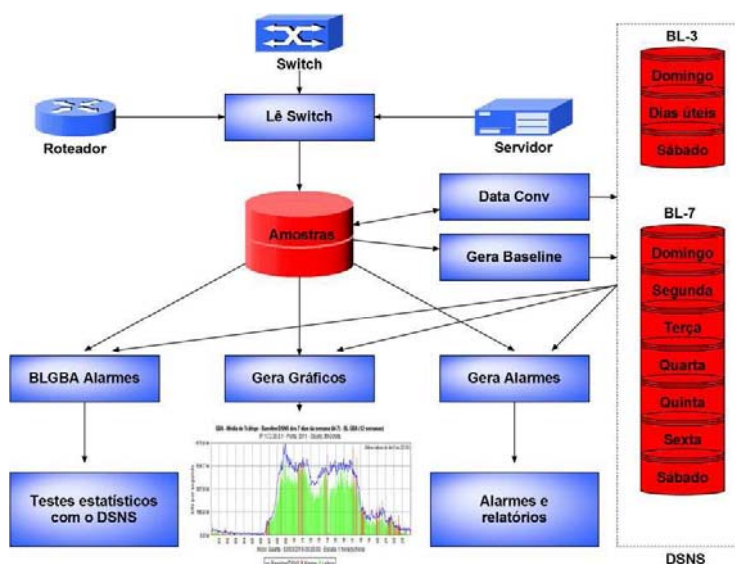
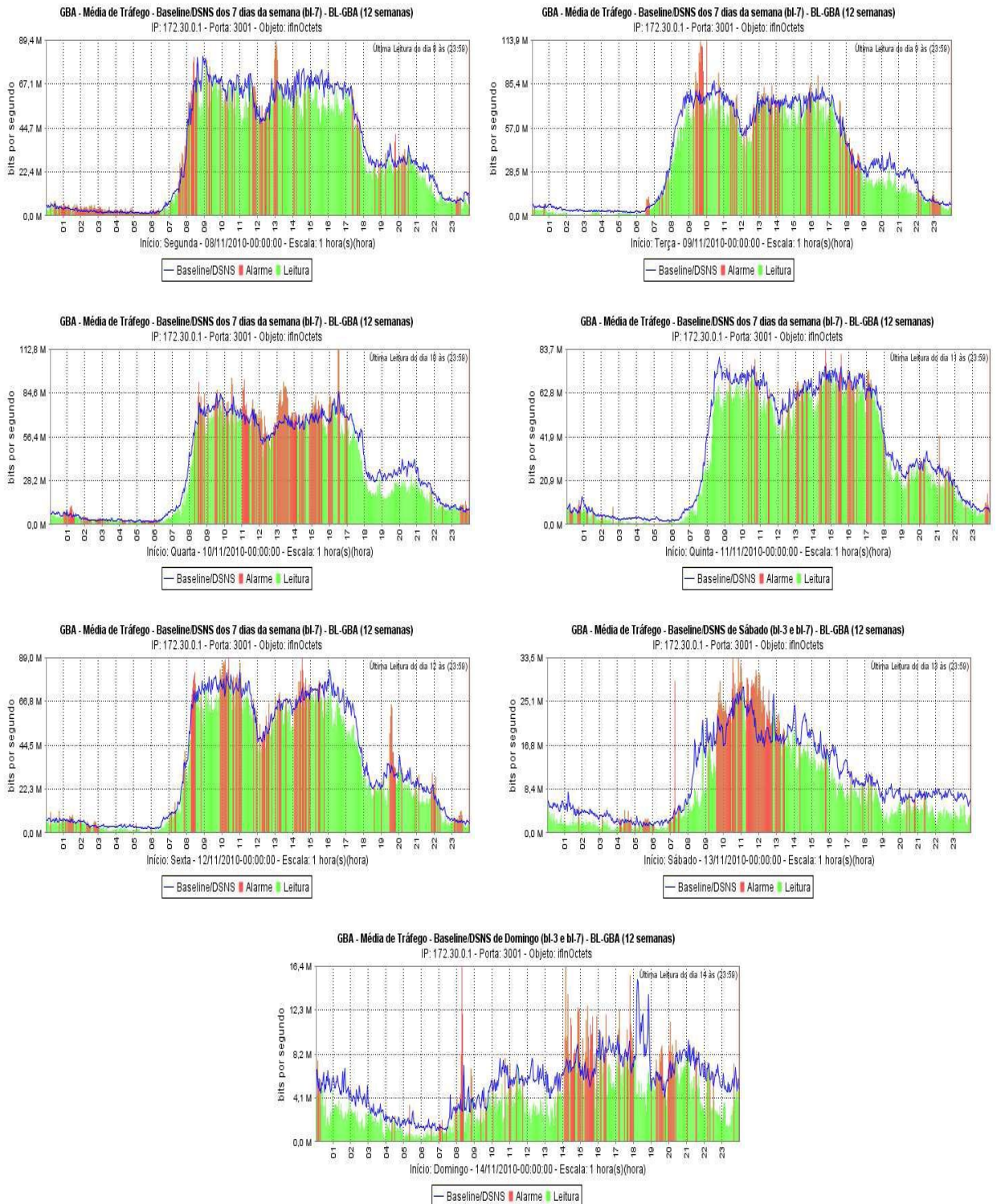


Figura 3.2 Tráfego e DSNS gerados pela ferramenta GBA34



A base de dados de amostras é formada pelas tabelas Equipamento e Objeto como visto na Figura 3.3. A tabela Equipamento é formada pelos atributos IP, comunidade, intervalo e info, onde:

- IP: Consiste do IP do equipamento monitorado;

- comunidade: nome da comunidade. Utilizado para autenticação entre agente e gerenteSNMP, nas versões SNMPv1 e SNMPv2 [30, 31];
- intervalo: consiste no intervalo de coleta ou sondagem (pooling) do equipamento;

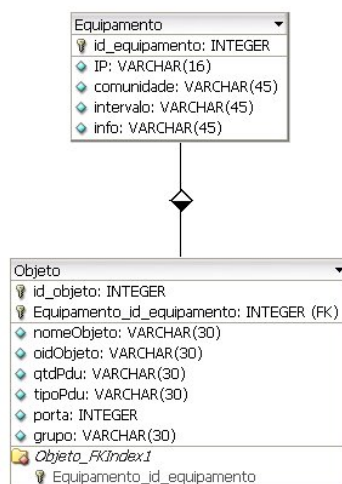
- info: informações adicionais sobre o equipamento monitorado.

A tabela Objeto armazena as informações relacionadas aos objetos SNMP de um determinado equipamento. Os atributos dessa tabela são:

- nomeObjeto: consiste do nome do objeto;
- oidObjeto: indicador numérico do objeto, que informa a posição do objeto dentro da hierarquia de nomes de objetos na MIB;
- qtdPdu: a quantidade de dados (PDU, packet data Unit);
- tipoPdu: o tipo de dado, que pode ser bit, frame, packet e segment;
- porta: a porta que está sendo monitorada;
- grupo: nome do grupo de objetos. Utilizado no modelo de segurança

do SNMPv3

Figura 3.3 Estrutura do banco de dados Amostras



O DSNS é armazenado em forma de arquivos individuais para cada dia da semana, que estão divididos em pastas, de acordo com o grupo que pertencem e objetos SNMP monitorados.

As seguintes funções estão implementadas na versão atual do GBA:

- Módulo Lê Switch: responsável pela leitura e armazenamento dos dados dos segmentos monitorados. Utiliza bibliotecas do protocolo SNMP capazes de acessar os objetos da MIB de switches e roteadores
- Módulo Gera Baseline: responsável pela geração do DSNS para dias úteis, sábados e domingos
- Módulo Gera Alarmes: responsável pela geração de alarmes, baseado na comparação entre os dados coletados e o DSNS
- Módulo Gera Gráficos: responsável por gerar gráficos a partir dos dados coletados, criando uma comunicação visual entre o usuário e a análise realizada pelo sistema, que demonstra o ajuste do tráfego monitorado em relação ao DSNS
- Módulo Data Conv: consiste de um módulo auxiliar desenvolvido para realizar a conversão dos arquivos de dados criados pelo Lê-switch (binários) em arquivos texto, e vice-versa, além de vários outros tipos de conversões nos arquivos de amostras e de DSNS.
- Módulo BLGBA Alarmes: utilizado para o desenvolvimento de estudos e pesquisas, com o objetivo de aperfeiçoar a ferramenta.

4 SISTEMA DE DETECÇÃO DE ANOMALIAS UTILIZANDO A ASSINATURA DIGITAL DE SEGMENTO DE REDE

Neste capítulo são apresentados os conceitos aplicados na construção do algoritmo base do sistema de detecção de anomalias (SDA) desenvolvido. São descritos os módulos e o funcionamento do SDA, que é baseado na combinação dos algoritmos *Particle Swarm Optimization*, K-means clustering e a Assinatura Digital de Segmento de Rede (DSNS) descrito no Capítulo 3. Foram realizados estudos a respeito da complexidade computacional do algoritmo, e também sobre a otimização dos parâmetros de entrada.

4.1 K-MEANS CLUSTERING E PARTICLE SWARM OPTIMIZATION

$$J(\mathbf{p}) = \sum_{k=1}^K \sum_{s=1}^S \|\mathbf{p}^k - \mathbf{c}^k\|^2 \quad (4.1)$$

onde K é o número de clusters, S o número de amostras do conjunto, \mathbf{p}^k o vetor de amostras e \mathbf{c}^k o k -ésimo centróide. O algoritmo do K-means pode ser visto no Pseudocódigo 1.

Nem sempre o algoritmo K-means consegue alcançar um ótimo global ao longo das iterações. Isso ocorre porque o algoritmo utiliza atribuições discretas ao invés de um conjunto de variáveis contínuas, de modo que a solução encontrada pode corresponder a um ótimo local. A fim de superar estas limitações e simultaneamente manter uma baixa complexidade computacional, o K-means pode ser associado ao algoritmo Particle Swarm Optimization (PSO)[8][33][34]. A maior preocupação no problema de detecção de anomalias é a classificação de uma enorme quantidade de dados (multi-objetos) em tempo real. em nosso trabalho a complexidade é um fator preponderante, o que torna a utilização da heurística interessante. Em nosso SDA utilizamos a combinação entre o K-means e a heurística PSO, o qual denominamos algoritmo PSO-CIs.

Pseudocódigo 1: algoritmo K-means clustering

1. Atribuir K pontos dentro do espaço representado pelos objetos que vão ser clusterizados. Esses pontos representam os centróides iniciais.
2. Atribuir cada objeto ao grupo que possui o centróide mais próximo.

3. Quando todos os objetos forem atribuídos, recalculer as posições dos K centróides.

4. Repetir os passos 2 e 3 até que nenhum centróide se mova mais. Isso produz uma separação dos objetos em grupos a partir da métrica escolhida

O PSO é uma técnica de computação evolucionária baseada em inteligência coletiva, criada por Kennedy e Eberhart em 1995, e inspirada no comportamento social dos pássaros e cardumes de peixes. Consiste de uma ferramenta de otimização computacional poderosa, pois é capaz de escapar de ótimos locais enquanto mantém uma estrutura simples e dependente de poucos parâmetros de ajuste. No PSO, as soluções no espaço de busca são denominadas de partículas. Cada partícula possui um valor de aptidão que é medido pela função a ser otimizada e uma velocidade de atualização que guia seu vôo e deslocamento através do espaço de busca. O princípio do PSO é o movimento em grupo das partículas, distribuídas aleatoriamente no espaço de busca, cada uma com sua própria posição e velocidade. A posição de cada partícula é modificada pela aplicação da velocidade, a fim de alcançar um melhor desempenho [34]. A interação entre as partículas é inserida no cálculo da velocidade da partícula. Assim, a cada iteração, a velocidade e a posição de todas as partículas de uma população de tamanho M são atualizadas. Se os melhores valores de posições para soluções locais ou globais forem encontrados, o respectivo melhor vetor candidato é atualizado, onde \mathbf{p}^{BEST} é o melhor valor encontrado até aquela iteração em cada partícula na população de tamanho M , e \mathbf{p}^{BEST} é a melhor posição obtida entre todas as partículas até o momento.

As partículas de ótimo global e local são vetores colunas com dimensão D .

Na estratégia do PSO, cada vetor candidato na iteração n , é definido como $\mathbf{p}_i[n]$ com dimensão $D \times 1$, e é utilizado para o cálculo da velocidade na próxima iteração como:

$$\mathbf{v}_i[n+1] = \omega \cdot \mathbf{v}_i[n] + \varphi_1 \cdot \mathbf{U}_i[n](\mathbf{p}^{\text{BEST}} - \mathbf{p}_i[n]) + \varphi_2 \cdot \mathbf{U}_i[n](\mathbf{p}_i^{\text{BEST}} - \mathbf{p}_i[n]) \quad (4.2)$$

onde ω é o peso da inércia; $\mathbf{U}_i[n]$ e $\mathbf{U}_i[n]$ são matrizes diagonais com dimensão D e elementos randômicos com distribuição uniforme $\square \mathbf{U} \in [0, 1]$, geradas para a i -ésima partícula na iteração $n = 1, 2, \dots, N$; \mathbf{p}^{BEST} e $\mathbf{p}_i^{\text{BEST}}$ são as posições do ótimo local e global encontradas até a n -ésima iteração respectivamente; φ_1 e φ_2 são os

coeficientes de aceleração relacionados à influência das partículas de ótimo local e global na atualização da velocidade respectivamente.

A posição da i -ésima partícula na iteração n é um vetor candidato a clusterização $\mathbf{p}_i[n]$ de tamanho $D \times 1$. A posição de cada partícula é atualizada utilizando o novo vetor de velocidade (4.2) para aquela partícula de acordo com:

$$\mathbf{p}_i[n + 1] = \mathbf{p}_i[n] + \mathbf{v}_i[n + 1], \quad i = 1, \dots, M \quad (4.3)$$

O algoritmo PSO consiste da aplicação repetida das equações de atualização da posição e da velocidade até que um critério de parada seja encontrado. O critério de parada pode ser um número fixo de iterações ou determinado pelo não avanço da solução quando o algoritmo evolui.

A fim de reduzir a probabilidade de que a partícula deixe o espaço de busca, um fator limitante da velocidade máxima V_{max} é adicionado ao modelo do PSO, o qual é responsável por limitar a velocidade na faixa $[\pm V_{max}]$. O ajuste da velocidade permite que a partícula se mova em um subespaço contínuo, mas restrito, sendo definido por

$$v_i[n] = \min \{V_{max}; \max \{-V_{max}; v_i[n]\}\} \quad (4.4)$$

Da equação (4.4) fica claro que se $|v_i[n]|$ exceder o valor de uma constante positiva V_{max} definida pelo usuário, a velocidade da i -ésima partícula será definida como $\text{sign}(v_i[n])V_{max}$, por exemplo, a velocidade das partículas em cada uma das D - dimensões fica restrita a uma magnitude máxima V_{max} . Se pudermos definir o espaço de busca pelos limites $[P_{MIN}; P_{MAX}]$, então o valor de V_{max} será atribuído tipicamente a $V_{max} = \tau(P_{MAX} - P_{MIN})$, onde $0.1 \leq \tau \leq 1.0$. [35].

Em relação ao parâmetro da inércia das partículas, ω , observou-se que um valor relativamente alto contribui para encontrar ótimos globais, minimizando a influência do valor das melhores posições da partícula e do grupo, enquanto um pequeno valor de inércia tende a contribuir para uma melhor convergência. Deste modo, a fim de encontrar um equilíbrio entre as habilidades de busca local e global, foi adotado o seguinte modelo de inércia:

$$\omega[n] = (\omega_{inicial} - \omega_{final}) \cdot \frac{l - m}{l} + \omega_{final} \quad (4.5)$$

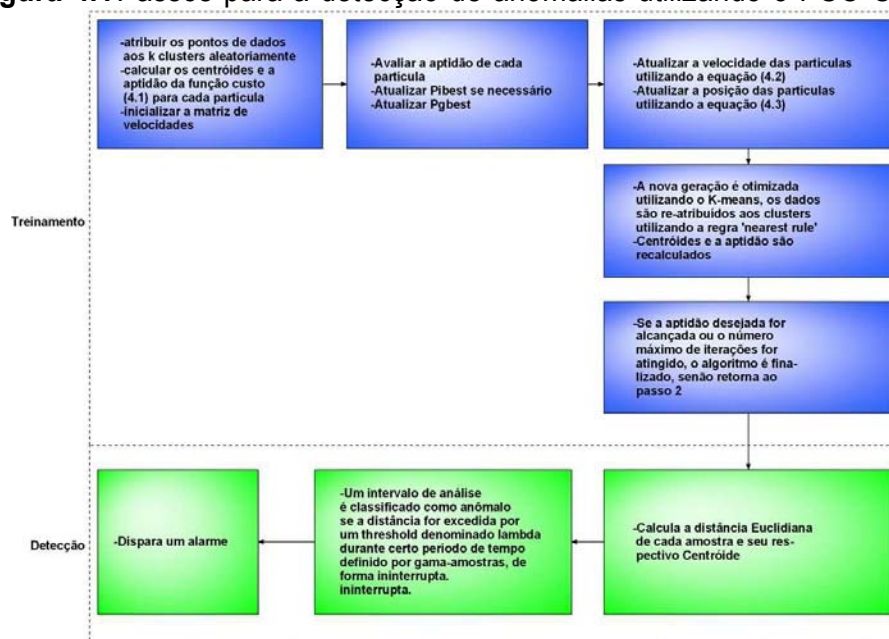
onde $\omega_{inicial}$ e ω_{final} são a inércia inicial e final, respectivamente, $\omega_{inicial} > \omega_{final}$, l é o número máximo de iterações, e $m \in [0, 6; 1, 4]$ o índice da adaptação não linear [36].

O processo de treinamento dos dados e detecção de anomalias é sintetizado na figura 4.1.

4.2 OTIMIZAÇÃO DOS PARÂMETROS DE ENTRADA

Experimentos foram realizados a fim de determinar os melhores valores dos parâmetros de entrada do PSO-CIs para o problema de otimização da equação (4.1). O critério de optimalidade adotado para nosso problema de detecção de anomalias, é aquele que resulta na maior taxa de detecção (TD) e simultaneamente na menor taxa de alarmes falsos (TAF). Foram otimizados os coeficientes de aceleração, φ_1 e φ_2 , a velocidade máxima permitida, V_{max} , a inércia, ω , e o tamanho da população, M . Para problemas de otimização contínuos como o investigado neste trabalho, os resultados numéricos apresentados no capítulo 6 indicam que após um número suficiente de iterações para atingir a convergência, a minimização da função custo foi obtida com valores baixos para os dois coeficientes de aceleração. O fator V_{max} também foi otimizado. A diversidade aumenta com a velocidade da partícula que tende

Figura 4.1 Passos para a detecção de anomalias utilizando o PSO-CIs.



A passar dos limites estabelecidos por (4.4). O valor de V_{max} determina a maior mudança possível na posição de uma partícula a cada iteração. Sem peso inercial, i.e. $\omega = 1$, Eberhart e Shi [37] descobriram que os melhores valores de velocidade máxima permitida estão entre 10 e 20% do intervalo dinâmico de cada dimensão da partícula. A escolha apropriada de V_{max} evita que a partícula tenha sua posição atualizada para um ponto muito distante do espaço de soluções promissoras. Portanto, para o nosso problema de otimização, uma busca não exaustiva indicou que o melhor *trade-off* de desempenho \times complexidade pode ser obtido designando um fator de velocidade máxima de $V_{max} = 0,2 \times (P_{MAX} - P_{MIN})$. Para o valor da inércia, ω , simulações confirmaram que valores altos implicam em uma convergência rápida, em detrimento da diversidade de busca, de tal forma que o algoritmo pode convergir facilmente para um ponto de ótimo local. Por outro lado, um valor pequeno de ω resulta em uma convergência lenta devido à quantidade muito grande de mudanças de posição das partículas em um espaço de busca pequeno. Neste trabalho foi utilizado valores de inércia adaptativos conforme mostra a Equação (4.5), com $m = 2$, e valores de inércia inicial e final de: $\omega_{inicial} = 1$ e $\omega_{final} = 0,01$. Assim, os valores máximos de velocidade ficam limitados pelos valores de inércia inicial e final multiplicados pelo parâmetro V_{max} , resultando em:

$$\omega_{inicial} \times V_{max} = 0,2 \times (P_{MAX} - P_{MIN}), \text{ e } \omega_{final} \times V_{max} = 0,002(P_{MAX} - P_{MIN}) \quad (4.6)$$

Finalmente, o critério de parada do algoritmo pode ser definido como sendo o número máximo de iterações ou um limiar de erro máximo permitido, dado por:

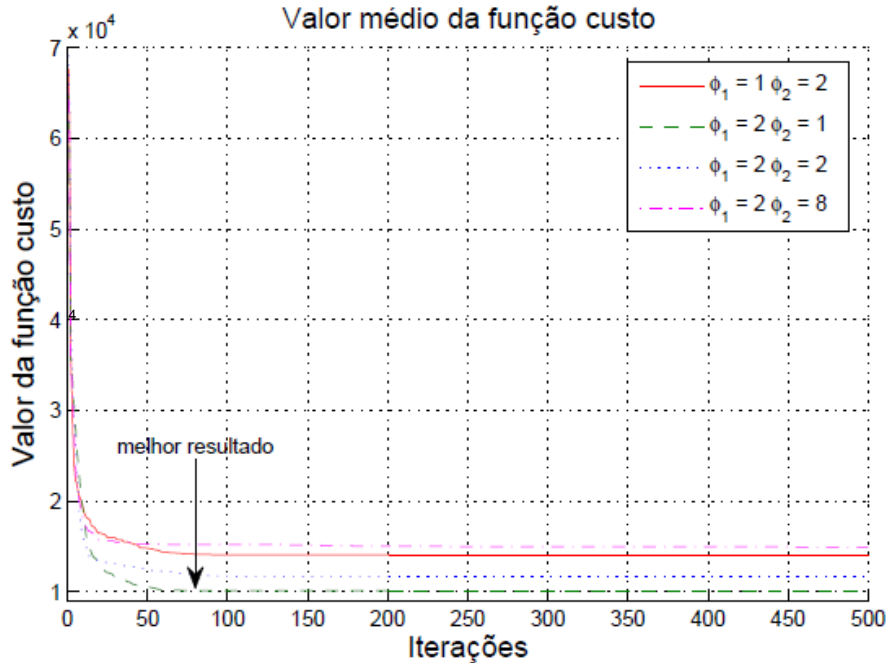
$$|J[I] - J[I - 10]| < \epsilon_{parada} \quad (4.7)$$

onde J é o valor da função custo na iteração I , e $\epsilon_{parada} \in [0,001; 0,01]$. A figura 4.2 apresenta a otimização dos parâmetros φ_1 e φ_2 do algoritmo PSO-CIs, com $M = 20$.

O melhor desempenho obtido foi para $\varphi_1 = 2$ e $\varphi_2 = 1$.

Figura 4.2 Evolução da função custo ao longo de 500 iterações, média em 500 realizações.

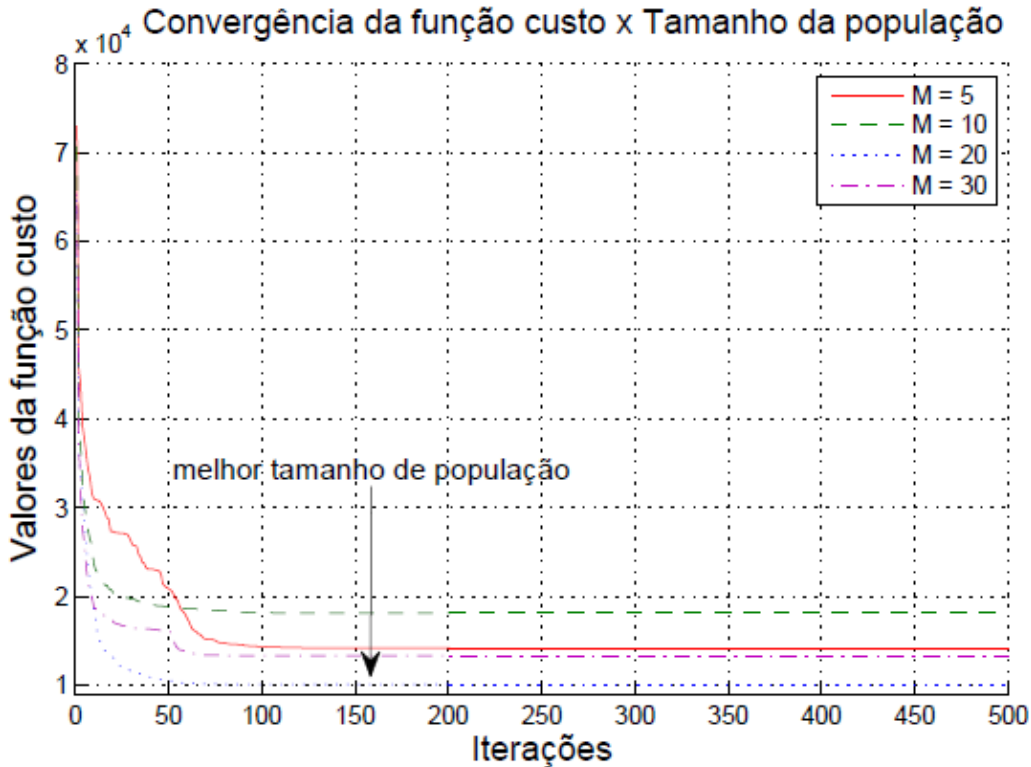
Diferentes combinações de parâmetros φ_1 e φ_2 , $M = 20$.



Na figura 4.3, é apresentado a convergência da função custo, levando em consideração o tamanho de população, fixando $\varphi_1 = 2$ e $\varphi_2 = 1$. O valor que apresentou melhor convergência foi $M = 20$.

O último parâmetro de entrada com influência nos resultados do algoritmo PSO-CIs, é K . Em nosso modelo, diferentemente de outros trabalhos que avaliam a existência de anomalias através do tamanho dos clusters [8], a detecção se dá através da avaliação da distância das amostras de um cluster em relação ao seu centróide. Assim, adotamos um valor fixo de $K = 1$, de modo que através da execução do algoritmo em pequenos intervalos de tempo, como visto na Seção 4.5, o objetivo é calcular apenas o centróide de um conjunto de dados n -dimensional.

Figura 4.3 Evolução da função custo ao longo de 500 iterações, média em 500 realizações. Diferentes tamanhos de população M .



4.3 ANÁLISE DE COMPLEXIDADE

A fim de avaliar a complexidade do algoritmo PSO-CIs, foi construída a Tabela 3 levando em consideração o número de operações matemáticas realizadas em cada equação, a cada iteração, em função do número de objetos da MIB (D) sendo monitorados pelo SDA.

Tabela 3 Número de operações por iteração

PSO-CIs	Eq.	Soma	Multiplicação
Função custo	4.1	$K \times S$	0
Atualização posição	4.3	$D \times M$	0
Atualização velocidade	4.2	$5D \times M$	$(3D + 2) \times M$

A complexidade das equações (4.3) e (4.2) está associada ao tamanho da população M e a quantidade de objetos monitorados D , enquanto a complexidade da função custo (4.1), depende da quantidade de clusters K e do número de amostras. Desta forma, a complexidade assintótica do PSO-CIs é da ordem de $O(DM + KS)$. Conforme descrito na Seção (6.1), o algoritmo PSO-CIs obteve uma média de 62

iterações para convergir totalmente. Sendo assim, a Figura 4.4 apresenta o número médio de operações matemáticas para que o algoritmo PSO-CIs obtenha a convergência, de acordo com diferentes números de objetos.

Pseudocódigo 2 Algoritmo PSO-CIs

Início

Entrada: conjunto de amostras de tráfego, DSNS

Saída: conjunto de amostras e DSNS clusterizados, centróides

1. Dados de entrada são clusterizados aleatoriamente
2. A população é iniciada com uma distribuição uniforme em $U[P_{\text{MIN}}; P_{\text{MAX}}]$

3. **Para** $n = 1$ até N

Para $i = 1$ até M

//calcula da velocidade

$$\mathbf{v}_i[n + 1] = \omega \cdot \mathbf{v}_i[n] + \varphi_1 \cdot \mathbf{U}_i[n](\mathbf{p}^{\text{BEST}} - \mathbf{p}_i[n]) + \varphi_2 \cdot \mathbf{U}_i[n](\mathbf{p}^{\text{BEST}} - \mathbf{p}_i[n])$$

//Limites de velocidade

$$v_i[n] = \min \{V_M; \max \{-V_M; v_i[n]\}\}$$

//Atualização da posição

$$\mathbf{p}_i[n + 1] = \mathbf{p}_i[n] + \mathbf{v}_i[n + 1], i = 1, \dots, M$$

Se $\mathbf{p}_i \in [P_{\text{MIN}}; P_{\text{MAX}}]$

Calcula a aptidão de \mathbf{p}_i e atualiza \mathbf{p}^{BEST}

e \mathbf{p}^{BEST} g

fimSe

fimPara i

//K-means clustering

Os dados de entrada são re-atribuídos aos clusters de acordo com a menor distância euclidiana

Encerra se N for atingido

fimPara n

4. \mathbf{p}^{BEST} determina os centróides

Fim

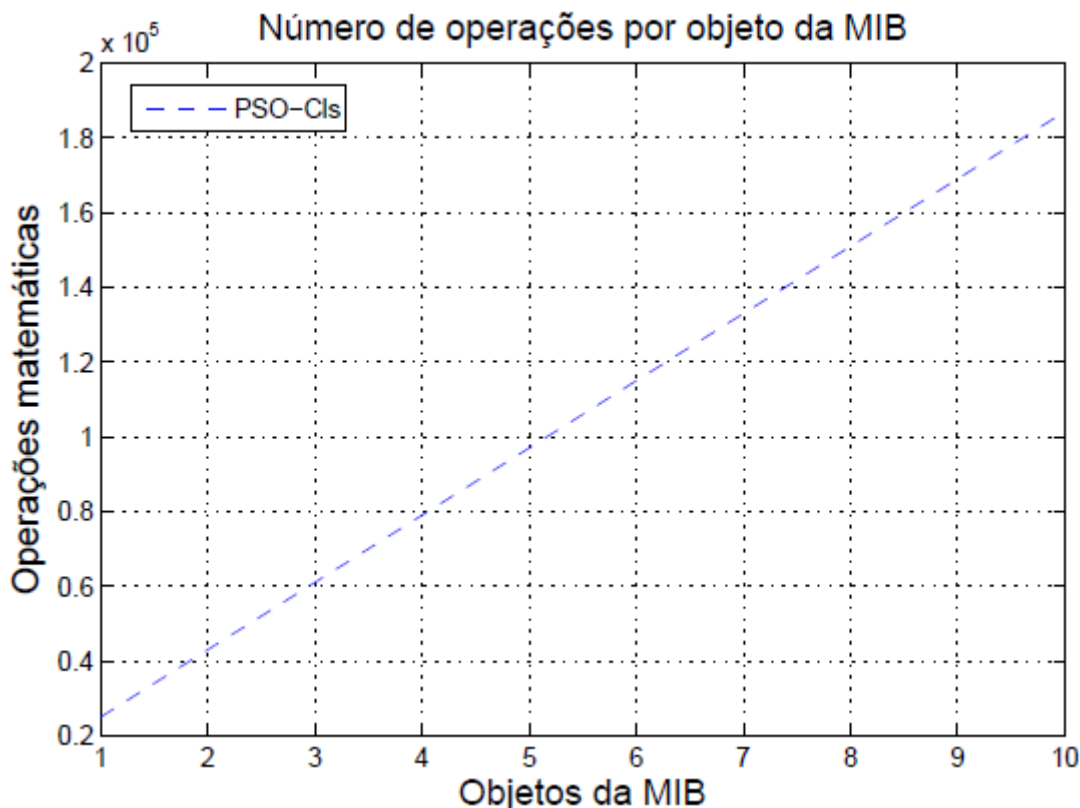
i \mathbf{p}^{BEST} : melhor valor encontrado em cada partícula

g \mathbf{p}^{BEST} : melhor valor encontrado dentre todas as partículas

P_{MIN} : menor valor encontrado nos dados de entrada

P_{MAX} : maior valor encontrado nos dados de entrada

Figura 4.4 Média de operações matemáticas para diferentes números de objetos da MIB.



4.4 DESCRIÇÃO DO SISTEMA DE DETECÇÃO DE ANOMALIAS (SDA)

O SDA desenvolvido é baseado na utilização do DSNS gerado pela ferramenta GBA, e nas informações coletadas dos objetos da MIB através do protocolo SNMP. Esse conjunto de dados, formado pelas amostras de tráfego e respectivos DSNS's, é processado pelo algoritmo PSO-Cls em intervalos de análise de 300 segundos. O tamanho dos dados pode variar de objeto para objeto, de acordo com o intervalo de coleta. De modo geral, os objetos presentes na MIB dos equipamentos da rede da Universidade Estadual de Londrina, possuem um intervalo de coleta de 10 segundos, o que resulta em conjuntos de 8640 amostras por dia. Foi utilizada a ferramenta MATLAB para o desenvolvimento e teste do SDA.

O algoritmo PSO-Cls é a base do módulo PSO-Cls system, responsável por clusterizar os dados e calcular os centróides. Com base nos dados processados pelo PSO-Cls system, o módulo *PSO Alarm System* analisa as amostras de tráfego uma a

uma, em busca de anomalias. A figura 4.5 ilustra a estrutura funcional do sistema de detecção desenvolvido. O processo para detecção de anomalias ocorre em 2 etapas:

1. O módulo *PSO-Cls system* recebe o tráfego coletado do objeto da MIB e seu respectivo DSNS, e processa esses dados em intervalos de análise de 300 segundos. A aplicação do algoritmo PSO-Cls nesse conjunto de dados, resulta nos dados clusterizados e nos seus respectivos centróides, os quais representam o comportamento esperado das amostras pertencentes ao cluster. O resultado dessa etapa é o conjunto de dados divididos em 288 intervalos de análise, e seus respectivos centróides. O pseudocódigo do algoritmo PSO-Cls utilizado para clusterizar os dados e calcular os centróides é apresentado no Pseudocódigo 2.

Figura 4.5 Modelo de referência do sistema de detecção de anomalias

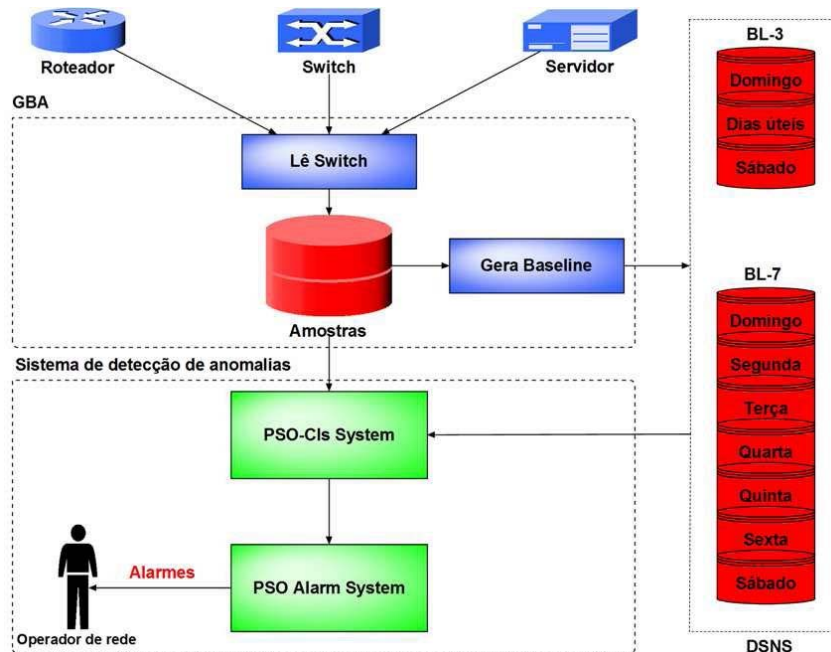
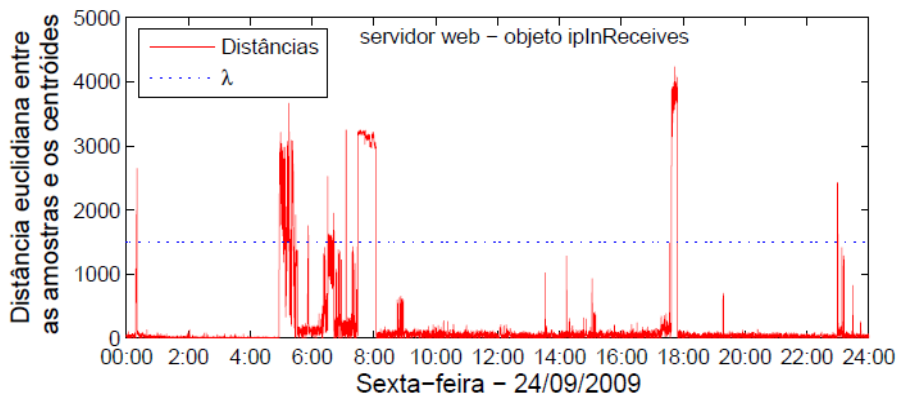


Figura 4.6 Alarmes gerados pelo sistema de detecção de anomalias



2. O módulo *PSO Alarm System* é responsável por analisar os dados processados na etapa anterior, buscando por anomalias nos intervalos analisados. O *PSO Alarm System* verifica a proximidade entre uma amostra de tráfego e seu respectivo centróide do cluster ao qual⁴⁶ pertence. A medida de distância adotada para tal, é a Euclidiana, a qual consiste da distância em linha reta entre dois pontos em um espaço bidimensional. Um intervalo de análise é classificado como anômalo se essa distância for excedida por um *threshold* denominado λ (lambda), durante certo período de tempo, definido por γ -amostras (gama), de forma ininterrupta. Quando isso ocorre, o SDA dispara uma alarme para notificar ao administrador de rede que ocorreu uma anomalia, como é ilustrado na Figura 4.6. No Pseudocódigo 3 podemos observar o algoritmo de detecção utilizado no módulo *PSO Alarm System*.

Pseudocódigo 3 Algoritmo de detecção

Início

Para $s = 1$ até S **faça**

Se $D(\text{amostras}(s), \mathbf{C}) > \lambda$ **então**

cont_anomalias + 1;

Se cont_anomalias > γ **então**

intervalo classificado como anômalo;

dispara um alarme;

cont_anomalias := 0;

fimSe

Senão

cont_anomalias := 0;

fimSe

fimPara

Fim

S : número de amostras, D : distância Euclidiana, \mathbf{C} : centróide

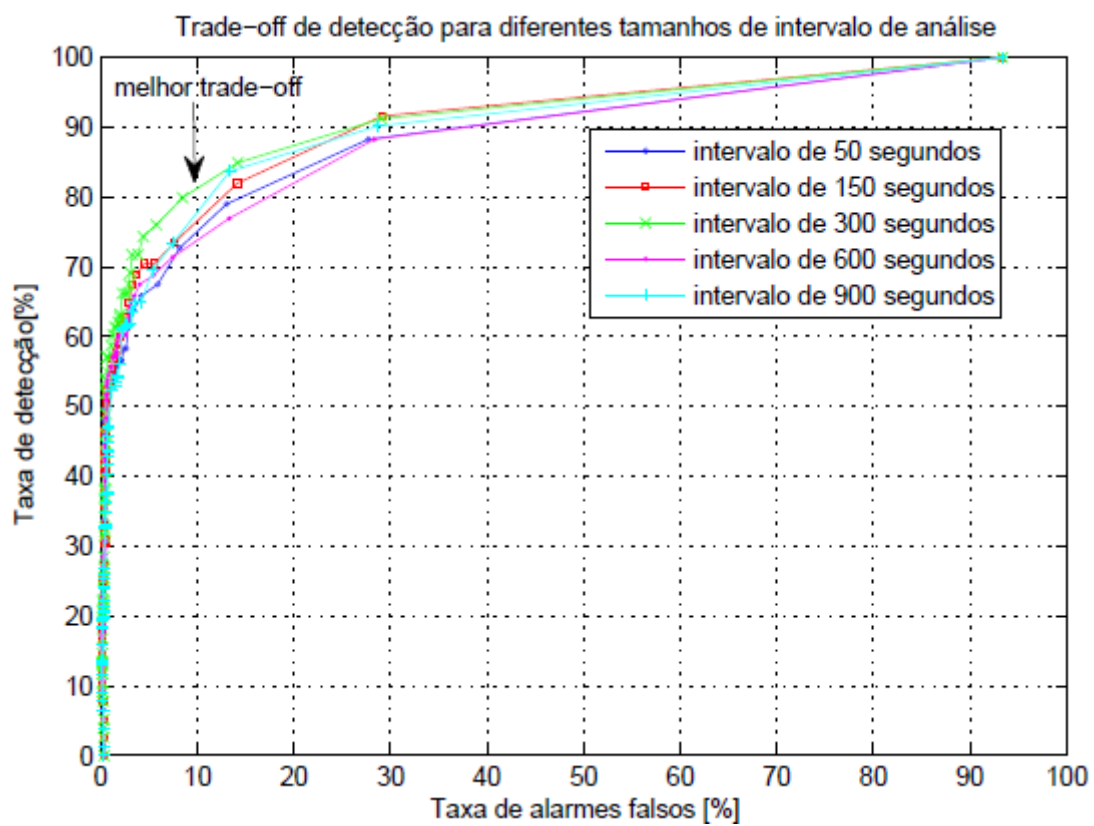
λ : *THRESHOLD* da distância Euclidiana, γ : limite de amostras anômalas

4.5 AVALIAÇÃO DO TAMANHO DO INTERVALO DE ANÁLISE

A fim de avaliar o impacto causado pela escolha do tamanho do intervalo de análise, foram realizados experimentos sob a perspectiva do Cenário 4 descrito na Seção 6.5. A figura 4.7 apresenta a taxa de detecção \times taxa de alarmes falsos, para

diferentes tamanhos de intervalo de análise. Os resultados obtidos através desse experimento indicam que o melhor intervalo de análise é 300 segundos⁴⁷

Figura 4.7 Desempenho do sistema de detecção para diferentes intervalos de análise.



5 ANÁLISE DE COMPONENTES PRINCIPAIS (PCA)

Recentemente, a análise de componentes principais (do inglês *Principal Component Analysis*, PCA), emergiu como uma técnica promissora para a detecção de anomalias. Inventado em 1901 por Karl Pearson, o PCA consiste de uma técnica para redução de dimensionalidade, utilizada para converter um conjunto de observações de variáveis possivelmente correlacionadas, em um conjunto de valores não correlacionados chamados de componentes principais (do inglês *Principal Component*, PC). Atualmente, o modelo é utilizado principalmente como um instrumento de análise exploratória de dados e para a criação de modelos preditivos. O PCA retorna uma representação compacta de um conjunto de dados multidimensional, através da projeção dos dados em um subespaço de menor dimensão, mas mantém a maior parte da variabilidade presente no conjunto de dados original, que é projetado sobre os novos eixos formados pelos componentes principais. Cada PC possui a propriedade de apontar na direção de máxima variância restante nos dados, considerando-se a variância já representada nos componentes anteriores.

A análise de componentes principais aplicada na detecção de anomalias de rede [38] consiste de uma metodologia de dois passos: 1) modelagem livre de anomalias, utilizando a decomposição das medições do tráfego em uma base de componentes principais e 2) a detecção de anomalias nos resíduos do tráfego, isto é, o tráfego não descrito pela decomposição PCA. O modelo consiste de uma transformação linear de coordenadas que mapeia um dado conjunto de dados em um novo sistema de coordenadas, de modo que a maior variação de qualquer projeção situa-se na primeira coordenada PC_1 , a segunda maior variação na segunda coordenada PC_2 , e assim por diante.

Seja $\mathbf{X} \in \mathbb{R}^{n \times m}$ uma matriz de dados contendo n amostras de m objetos da MIB recolhidos em operação normal de rede, ou seja, cada coluna m representa uma série temporal de n amostras para cada objeto SNMP. Esta matriz deve ser normalizada com média zero e variância unitária com o propósito de remover a diferença de escala entre as amostras dos diferentes objetos. Deste modo, o primeiro passo para se calcular o PCA é construir a matriz covariância \mathbf{R} , dada por:

$$\mathbf{R} = \frac{1}{n-1} \mathbf{X}^T \mathbf{X} \quad (5.1)$$

pós isso, deve ser então calculada a decomposição em valores singulares (do inglês Singular Value Decomposition, SVD), de acordo com:

$$\mathbf{R} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T \quad (5.2)$$

onde $\mathbf{\Lambda}$ é uma matriz diagonal que contém em sua diagonal os autovalores de \mathbf{R} ordenados decrescentemente ($\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m \geq 0$). As colunas da matriz \mathbf{V} são os autovetores de \mathbf{R} . A matriz de projeção $\mathbf{P} \in \mathbb{R}^{m \times a}$ é gerada selecionando-se a autovetores ou colunas de \mathbf{V} correspondentes aos a maiores autovalores, que indicam quais são os objetos SNMP com maior variabilidade no conjunto de dados recolhidos em operação normal transforma o espaço original das variáveis em um espaço de dimensão reduzida.

$$\mathbf{T} = \mathbf{X}\mathbf{P} \quad (5.3)$$

As colunas da matriz \mathbf{P} são denominadas loadings e os elementos de \mathbf{T} de scores. Os loadings podem ser entendidos como os pesos para cada variável original no cálculo do componente principal, enquanto os scores são os valores das variáveis originais transformadas no espaço de dimensão reduzida. Operando na equação (5.3), os scores podem ser transformados novamente no espaço original:

$$\hat{\mathbf{X}} = \mathbf{T}\mathbf{P}^T \quad (5.4)$$

de modo que a matriz de resíduos \mathbf{E} pode ser calculada da seguinte forma:

$$\mathbf{E} = \mathbf{X} - \hat{\mathbf{X}} \quad (5.5)$$

O espaço de dados original pode ser reconstruído da seguinte forma:

$$\mathbf{X} = \mathbf{T}\mathbf{P}^T + \mathbf{E} \quad (5.6)$$

A escolha dos a componentes principais é essencial, pois $\hat{\mathbf{X}}$ representa a principal fonte de variabilidade no processo, sendo que \mathbf{E} representa

a variabilidade correspondente ao ruído. Existem diversos métodos utilizados para determinar a quantidade de componentes principais a serem utilizados no modelo PCA, como visto em [39]. Neste trabalho utilizamos a abordagem *Cumulative Percent Variance* (CPV) para determinar a . De acordo com Zumoffen [40] o cálculo do CPV é dado por:

$$CPV(a) = \frac{\sum_{i=1}^a \lambda_i}{\text{TRACE}(\mathbf{R})} 100 \quad (5.7)$$

onde $\text{trace}(\mathbf{R})$ é a soma dos elementos da diagonal da matriz \mathbf{R} .

Uma vez estabelecido o modelo PCA, a detecção de anomalias pode ser reduzida à análise de apenas uma variável denominada Q , ou *square prediction error* (SPE), de modo que a ocorrência de novos eventos pode ser detectada através do cálculo do SPE sobre o resíduo. O escalar Q é uma medida de adequação da amostra ao modelo, e está diretamente associada ao ruído do conjunto:

$$\begin{aligned} Q &= \mathbf{r}^T \mathbf{r}, \text{ onde} \\ \mathbf{r} &= (\mathbf{I} - \mathbf{P}\mathbf{P}^T) \mathbf{y} \end{aligned} \quad (5.8)$$

onde \mathbf{y} corresponde as amostras de todos objetos SNMP no instante n .

O *threshold* de Q é calculado da seguinte forma:

$$Q_\alpha = \theta_1 \left[\frac{h_0 c_\alpha \sqrt{2\theta_2}}{\theta_1} + 1 + \frac{\theta_2 h_0 (h_0 - 1)}{\theta_1^2} \right]^{\frac{1}{h_0}} \quad (5.9)$$

onde

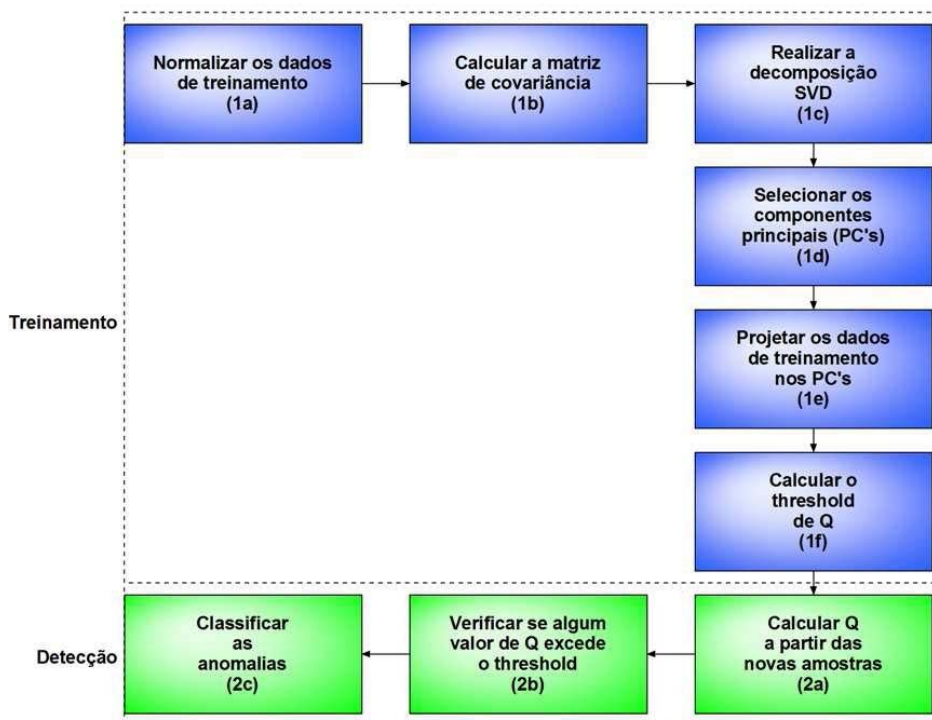
$$\theta_i = \sum_{j=a+1}^m \lambda_j^i \quad e \quad h_0 = 1 - \frac{2\theta_1\theta_3}{3\theta_2^2}$$

onde h_0 , θ_1 , θ_2 e θ_3 são calculados utilizando-se λ^i , que representa os autovalores; c_α é o valor da distribuição normal com α nível de significância. A detecção

de anomalias se dá através de um alto valor na variável Q , que é provocado por uma mudança no modelo de covariância, decorrente de um evento incomum nos dados.

Em síntese, para se realizar a detecção de anomalias utilizando-se o PCA, é necessário seguir os seguintes passos ilustrados na figura 5.1:51

Figura 5.1 Passos para a detecção de anomalias utilizando o PCA



1. Treinamento dos dados

(a) Realizar a normalização da matriz de tráfego \mathbf{X} , de forma que cada coluna possua média zero e variância unitária. Isso garante que as dimensões do PCA capturem uma variância verdadeira, evitando assim, a distorção dos resultados devido a diferenças na ordem de grandeza dos valores entre os objetos. A normalização pode ser realizada através da seguinte equação:

$$\mathbf{X} = (\mathbf{X} - \mathbf{ones}(n, 1) \cdot \mu) / (\mathbf{ones}(n, 1) \cdot \sigma) \quad (5.10)$$

onde $\mathbf{ones}(n, 1)$ é um vetor coluna com n linhas, μ a média e σ o desvio padrão de \mathbf{X} .

(b) Calcular a matriz de covariância de \mathbf{X} , dada pela equação 5.1

(c) Realizar a decomposição em valores singulares SVD, a fim de obter os autovalores de \mathbf{R} , que são os componentes principais.

(d) Selecionar os componentes principais que representam a maior variabilidade no conjunto, através de uma métrica. Em nosso trabalho, utilizamos a técnica *Cumulative Percent Variance* (CPV) [40], que consiste de uma medida de variância percentual ($CPV(a)$, 90%) capturado pelos a primeiros componentes principais.⁵²

(e) Projetar o conjunto de dados de treinamento no novo espaço dos componentes principais, e calcular o resíduo conforme a equação 5.5

(f) Calcular o *threshold* de Q conforme a equação 5.9

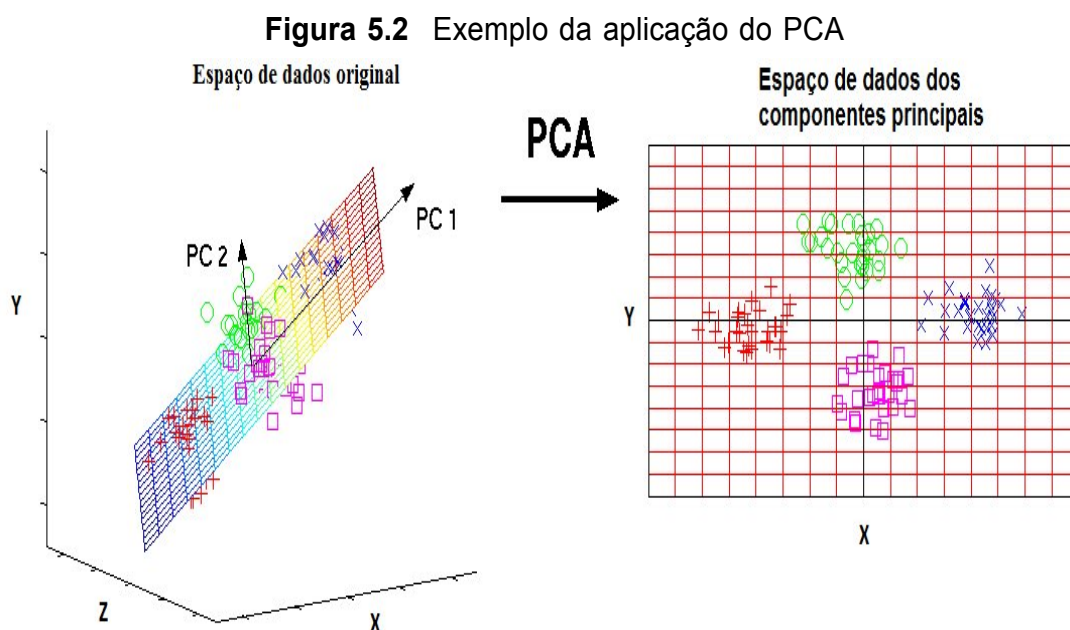
2. Detecção de anomalias

(a) Calcular Q a partir das novas amostras, conforme a equação 5.8

(b) Verificar se algum valor de Q excede o *threshold*

(c) Se uma amostra excede o *threshold* de Q , deve ser classificada como anômala

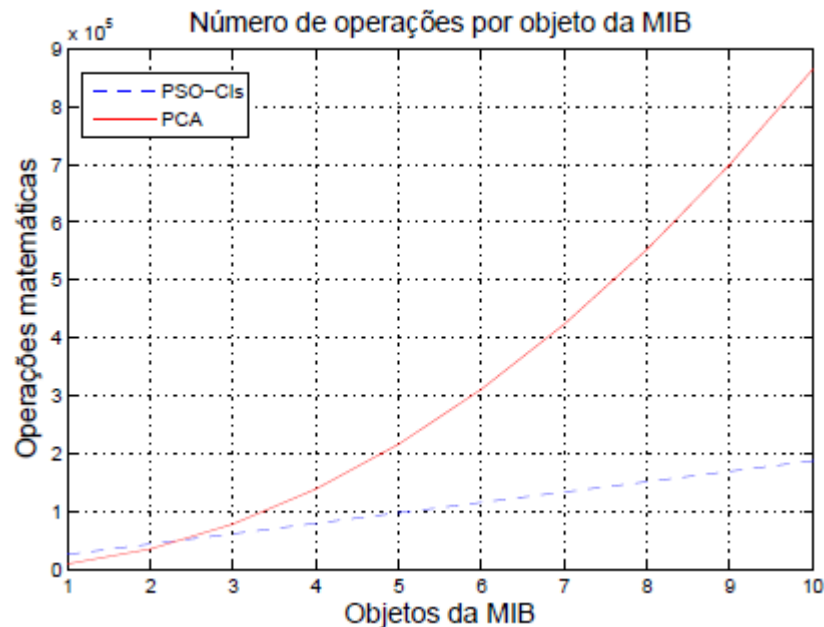
Na figura 5.2, temos um exemplo da aplicação do PCA sobre um conjunto de dados tri-dimensional, onde a maior parte dos dados está localizada dentro de um conjunto bi-dimensional. O PCA é utilizado para representar esses dados através da redução da dimensionalidade: As três variáveis originais são reduzidas a um número menor de duas novas variáveis (componentes principais). No lado esquerdo da figura, utilizando PCA, podemos identificar um plano bi-dimensional que descreve a maior variância dos dados. Este subespaço é o resultado da projeção no conjunto dos componentes principais, e é apresentado como um conjunto bi-dimensional, visto no lado direito da figura.



5.1 ANÁLISE DE COMPLEXIDADE

De acordo com Lakhina [38], computar todos os componentes principais da matriz de tráfego \mathbf{X} , é equivalente a resolver o problema do autovalor simétrico para a matriz de covariância, $\mathbf{X}^T \mathbf{X}$. O processo padrão para resolução desse problema baseia-se em computar a decomposição SVD de \mathbf{X} . A complexidade para se computar a SVD de uma matriz $n \times m$ é de $O(nm^2)$. A figura 5.3 apresenta o número de operações matemáticas necessárias para a execução do PCA, em um conjunto de dados com $n = 8640$ pontos de dados, para diferentes quantidades de objetos SNMP, em comparação com os resultados obtidos pelo PSO-Cls, como visto na Figura 4.4.

Figura 5.3 Operações matemáticas para diferentes números de objetos da MIB.



A reta pontilhada representa a complexidade do algoritmo PSO-Cls, na ordem de $O(DM + KS)$, onde D representa a quantidade de objetos da MIB, M o tamanho da população do PSO-Cls, K o número de clusters e S o número de amostras do conjunto de dados. A curva em vermelho representa a complexidade do PCA que é da ordem de $O(nm^2)$ onde n corresponde ao número de amostras do conjunto e m o número de objetos da MIB.

6 RESULTADOS

Com o objetivo de avaliar o desempenho e precisão dos SDAs desenvolvidos, foram realizados experimentos utilizando-se diferentes conjuntos de dados e objetos SNMP, levando em consideração o impacto da alteração nos parâmetros de ajuste do algoritmo PSO-CIs. Neste capítulo são apresentados e discutidos os resultados obtidos através da aplicação do SDA desenvolvido sobre esses conjuntos de dados, coletados nos equipamentos da rede da Universidade Estadual de Londrina (UEL).

Tabela 4 Cenários de teste

Cenário	Experimento	Objeto	Período
1	Avaliação do algoritmo PSO-CIs	dados simulados	
2	Avaliação do SDA desenvolvido, comparação com os resultados de um modelo determinístico	<i>ipInReceives</i>	20 – 26/04/2010
3	Definição parametrizada de anomalia, define-se uma classe de anomalias	<i>ifInOctets</i>	08/02/2010
4	Definição parametrizada de anomalia, define-se duas classes de anomalias	<i>ipInReceives</i>	05 – 09/04/2010
5	Experimentos com 4 objetos simultâneos, e avaliação do SDA baseado em PCA	<i>ifInOctets</i> <i>ipInReceives</i> <i>ipInDelivers</i> <i>tcpInSegs</i>	05 – 09/04/2010

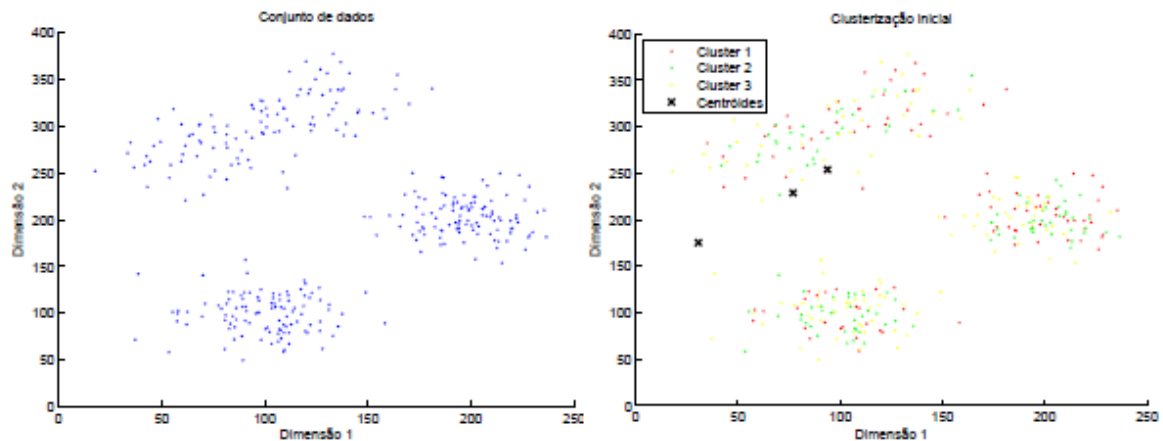
Foram montados cinco cenários, a fim de abordar os aspectos importantes para a validação do SDA. A tabela 4 apresenta uma breve descrição dos experimentos realizados em cada cenário.

6.1 CENÁRIO 1

A fim de avaliar a capacidade do algoritmo PSO-CIs em clusterizar conjuntos de dados e calcular os centróides, foi simulado um conjunto de dados bidimensional dividido em três subconjuntos, gerados através de uma distribuição estatística normal com médias (100; 100), (102, 5; 302, 5) e (200; 200) para os subconjuntos 1, 2 e 3 respectivamente, e 55 desvio padrão $\sigma = 20$. Na tabela 5 são apresentados os parâmetros de entrada do algoritmo PSO-CIs, utilizados na simulação.

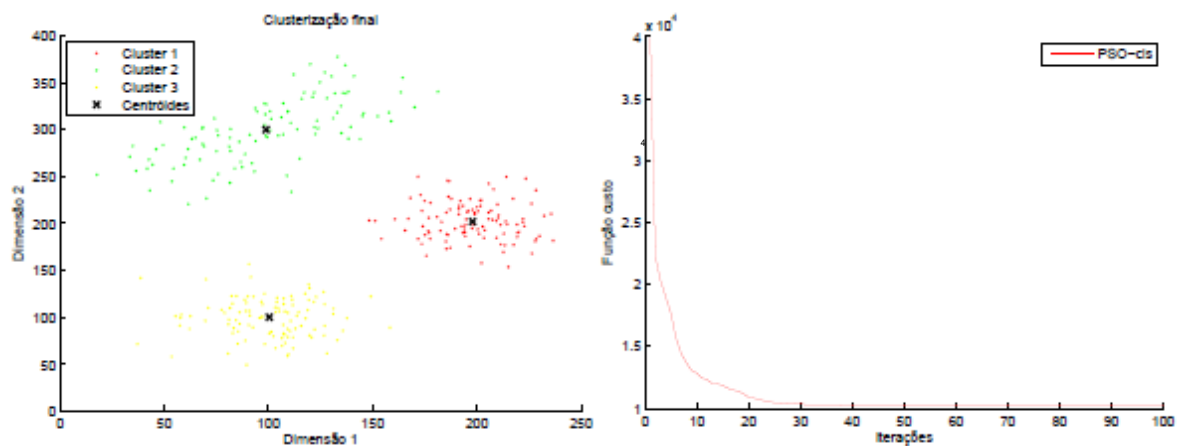
A figura 6.1a apresenta o conjunto de dados original, enquanto a figura 6.1b apresenta a primeira iteração do algoritmo PSO-CIs, denominada de clusterização inicial, onde os dados são aleatoriamente atribuídos aos clusters e os K centróides são atribuídos.

Figura 6.1 a)Conjunto de dados original b)Clusters e centróides iniciais.



A figura 6.2a, apresenta a iteração final do algoritmo, com cada dado do conjunto original atribuído ao cluster com centróide mais próximo. A figura 6.2a apresenta a evolução da função custo ao longo das iterações. Os testes foram realizados 500 vezes, de modo que a média de iterações necessárias para a convergência do algoritmo foi de 62.

Figura 6.2 a)Iteração final do algoritmo PSO-CIs b)Convergência da função custo (eq.4.1).



O algoritmo PSO-CIs mostrou-se eficiente na clusterização e cálculo dos centróides, atingindo uma média de 62 iterações para a convergência

Tabela 5 Parâmetros utilizados no algoritmo PSO-CIs

Parâmetros	Valores adotados
<i>PSO</i>	
Fator peso solução local	$\varphi_1=2$
Fator peso solução global	$\varphi_2=1$
Peso inicial da inércia	$W_{inicial} = 1$
Peso final da inércia	$W_{final} = 0.01$
Tamanho da população	$N = 20$
Número máximo de iterações	$It = 100$
<i>K-means</i>	
Número de clusters	$K = 3$

6.2 CENÁRIO 2

Este cenário utiliza tráfego real da rede da Universidade Estadual de Londrina, coletado durante uma semana no período de 20 – 26/04/2009, do objeto *ipInReceives* da MIB, no servidor web. O objeto *ipInReceives* pertencente ao grupo IP e determina a taxa de datagramas recebidos pelo elemento de rede monitorado. A figura 6.3 apresenta os gráficos do tráfego e dos respectivos DSNS's utilizados nos testes, enquanto a figura 6.4 apresenta os alarmes disparados pelo sistema de detecção de anomalias para esse período.

A avaliação do desempenho do SDA desenvolvido é baseada em duas métricas de desempenho: a taxa de detecção, que consiste na probabilidade de detecção dada por (6.1), e a taxa de alarmes falsos, que representa a probabilidade dos alarmes não apresentarem uma variação significativa entre o tráfego e o DSNS, dada por (6.2). As variáveis utilizadas para calcular a taxa de detecção e a taxa de alarmes falsos são:

- *correctly_detected*: número de anomalias corretamente detectadas.
- *occ_anomalies*: número de anomalias existentes no tráfego.
- *false_positives*: número de alarmes que não corresponde a uma situação anômala.
- *total_observations*: número total de observações

$$detection_rate = correctly_detected/occ_anomalies \quad (6.1)$$

$$false_alarm_rate = false_positives/total_observations \quad (6.2)$$

Figura 6.3 Gráficos com o tráfego e DSNS utilizado nos testes com o objeto *ipInReceives* no período de 20 – 24/04/2010 no servidor web da UEL.

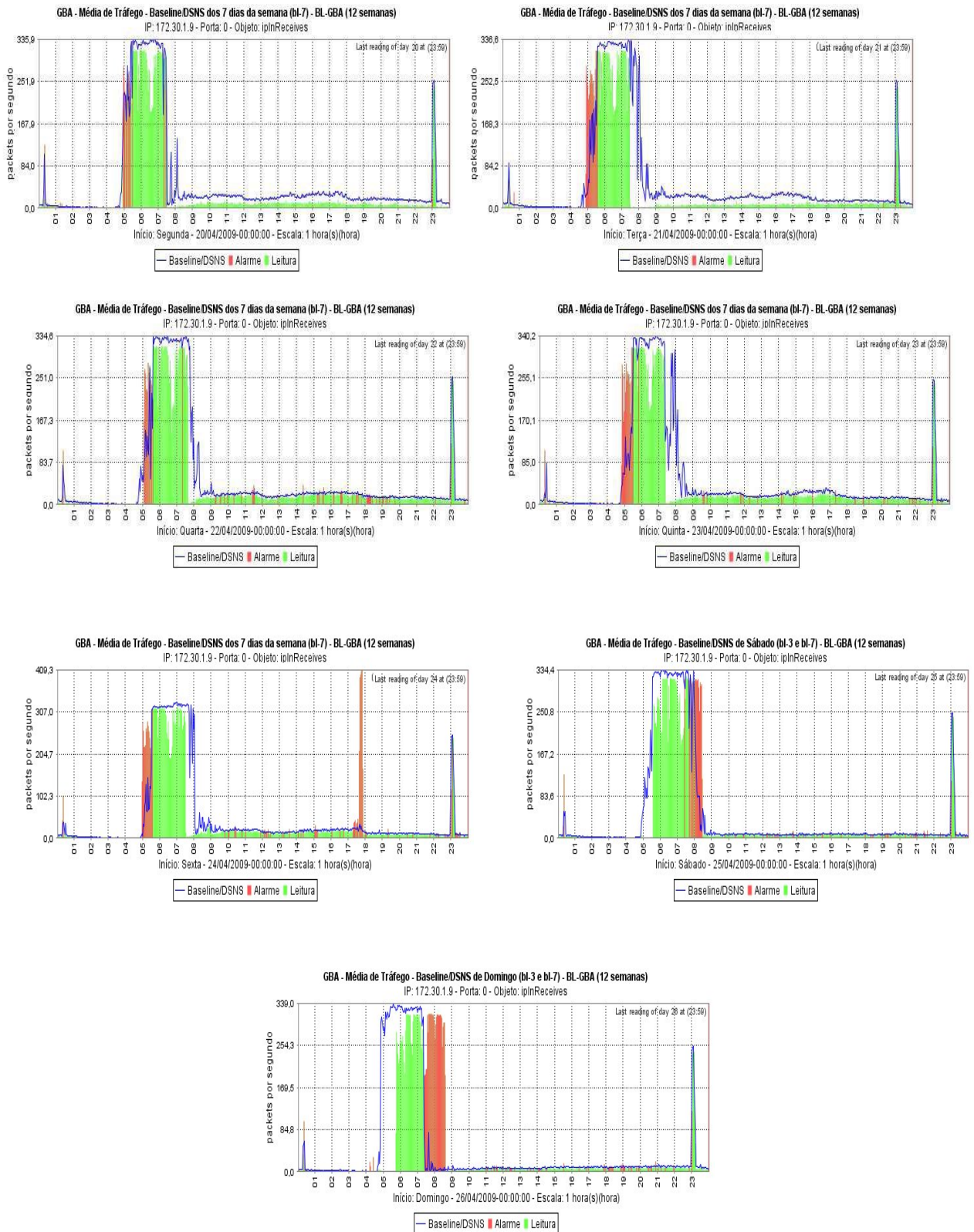
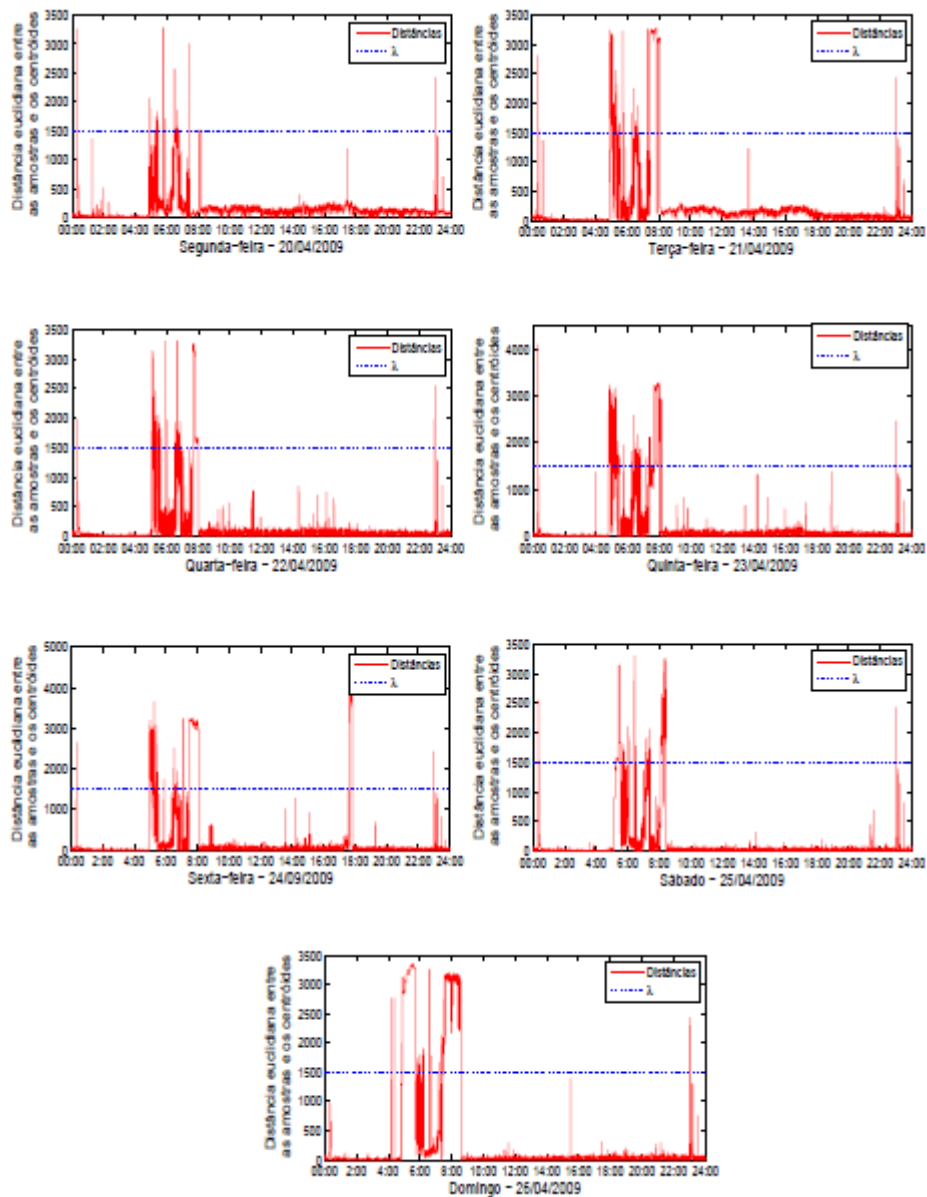
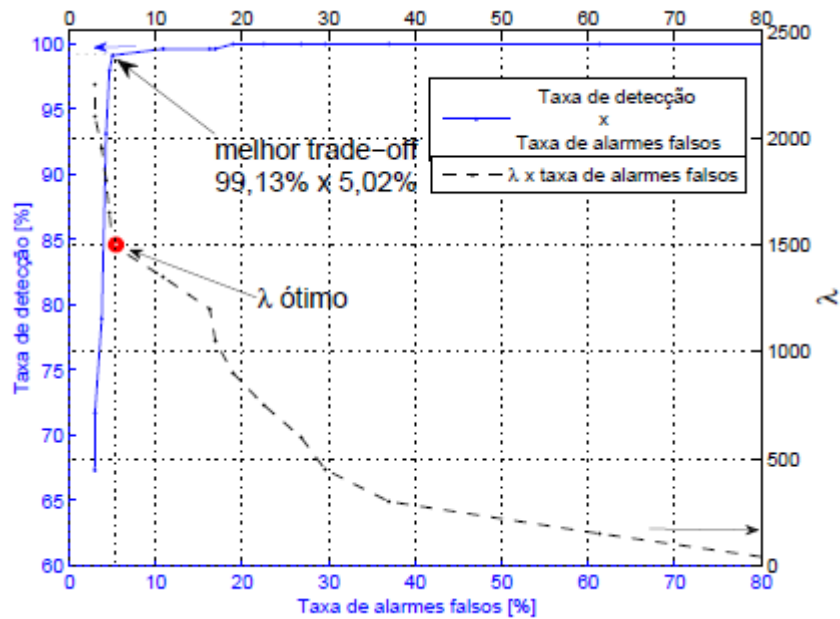


Figura 6.4 Alarmes gerados pelo SDA desenvolvido no período de 20 – 26/04/2009, objeto *iplnReceives* e λ



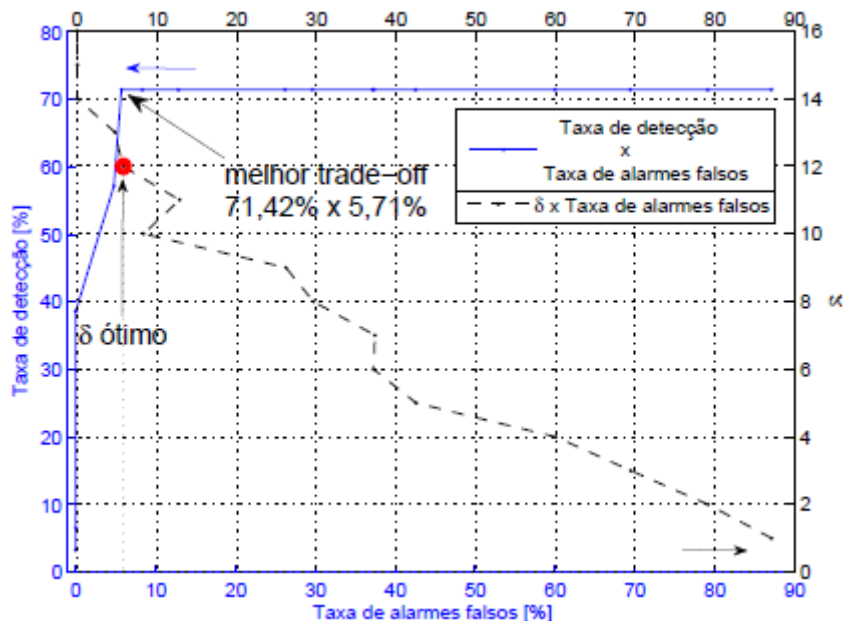
Com o objetivo de avaliar a precisão do SDA, foram executados experimentos com diferentes valores de *threshold* λ , a fim de se encontrar o melhor *trade-off* entre taxa de detecção e alarmes falsos. A figura 6.5 descreve o desempenho do algoritmo PSO-CIs em termos do *trade-off*, dado λ como parâmetro variável. Os valores de λ representados através da curva RoC [41] estão no intervalo de [1350; 1800]. Os resultados obtidos após 100 realizações do algoritmo PSO-CIs para cada valor de λ , confirmaram que o método é aplicável na detecção de anomalias, alcançando o melhor *trade-off* entre taxa de detecção e alarmes falsos: 99, 13% e 5, 02% para $\lambda = 1500$.

Figura 6.5 Taxa de detecção \times taxa de alarmes falsos e $\lambda \times$ taxa de alarmes falsos, gerados pelo SDA desenvolvido, no cenário 2.



A fim de comparar a precisão do SDA desenvolvido que é baseado em um algoritmo heurístico, o PSO-CIs, foi implementado um SDA baseado em um algoritmo determinístico, proposto por Zarpelão em [24]. O método é baseado no mecanismo de histerese que utiliza um parâmetro denominado δ , definido pelo usuário, que determina o tamanho do intervalo em que amostras fora do padrão caracterizam uma anomalia. Os experimentos realizados para avaliar esse algoritmo, foram executados utilizando-se o mesmo conjunto de dados, variando-se os valores de δ com o objetivo de determinar qual o valor ótimo de delta que resulta no melhor *trade-off* entre taxa de detecção e falsos alarmes. A figura 6.6 apresenta os resultados obtidos por esse modelo, e demonstra que o modelo determinístico foi capaz de alcançar uma taxa de detecção de 71,42% contra 5,71% alarmes falsos, para $\delta = 12.60$

Figura 6.6 Taxa de detecção \times taxa de alarmes falsos e $\delta \times$ taxa de alarmes falsos para SDA baseado no algoritmo determinístico, cenário 2.



A comparação dos resultados obtidos através dos experimentos com ambos os modelos, heurístico e determinístico, utilizando-se os dados do cenário 2, demonstrou que o modelo heurístico baseado no algoritmo PSO-CIs, obteve um ganho de desempenho de 27,71% em relação à taxa de detecção, e uma redução de 0,69% na taxa de alarmes falsos, sobre o algoritmo determinístico. Esses resultados indicam que o SDA desenvolvido é mais adequado para o ambiente de rede utilizado nesse experimento, e apresenta um grande potencial para a detecção de anomalias.

6.3 CENÁRIO 3

Com o objetivo de avaliar o SDA desenvolvido considerando diferentes características de anomalias, foi implementada uma definição parametrizada de anomalias de volume. São utilizados dois parâmetros nessa definição, α , que está relacionado com a amplitude de anomalia e γ que representa sua duração. Foi introduzido outro conceito denominado intervalo de alerta (histerese), que funciona da seguinte forma: se durante o monitoramento, uma amostra de tráfego exceder ou ficar abaixo do DSNS em $\alpha\%$, um intervalo de alerta é aberto. Se dentro do intervalo de alerta existirem γ amostras que excedam ou fiquem abaixo do DSNS em $\alpha\%$, este intervalo é classificado como anômalo, e é disparado um alarme para alertar ao operador

da rede. Para avaliar essa implementação foram realizados dois experimentos: o primeiro utilizando apenas um dia de tráfego, onde ocorreu uma anomalia de longa duração, e o segundo, utilizando cinco dias úteis.

Nesse cenário foi utilizado o tráfego de rede coletado da UEL, referente ao dia 08/02/2010 do objeto *iflnOctets*, pertencente ao grupo *interface* da MIB, do servidor web. O objeto *iflnOctets* determina o número total de octetos recebidos pela interface de rede. O objetivo é verificar se o algoritmo é capaz de identificar e disparar alarmes, para a anomalia ocorrida no período das 17 – 22 horas, como visto na figura 6.7. A figura 6.8 apresenta os alarmes disparados pelo SDA no período de teste.

Figura 6.7 Tráfego e DSNS do dia 08/02/2010 do objeto *iflnOctets* do servidor web da UEL.

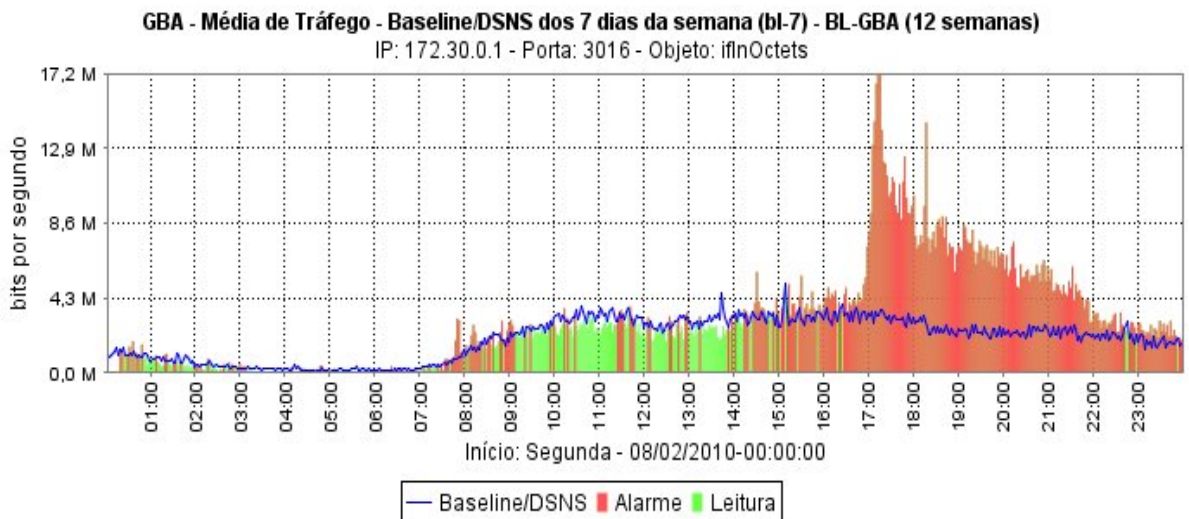
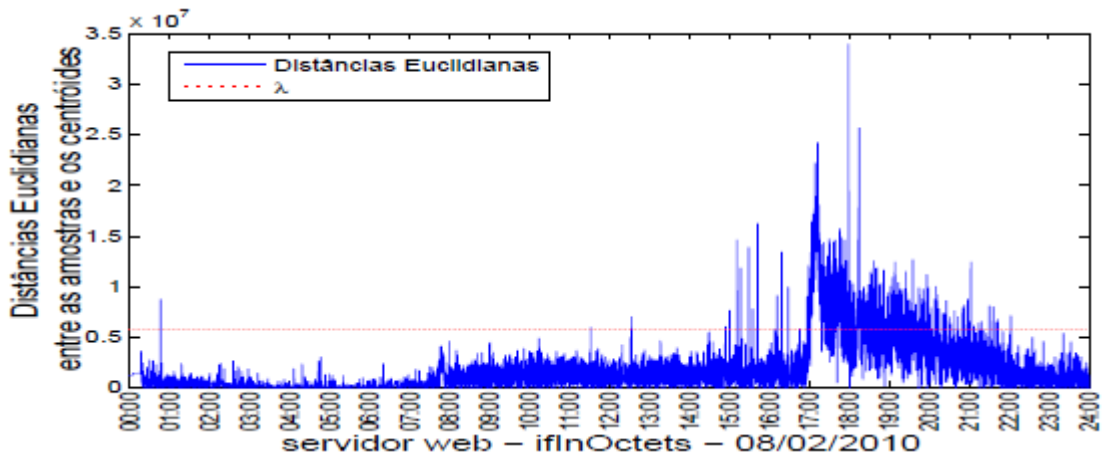
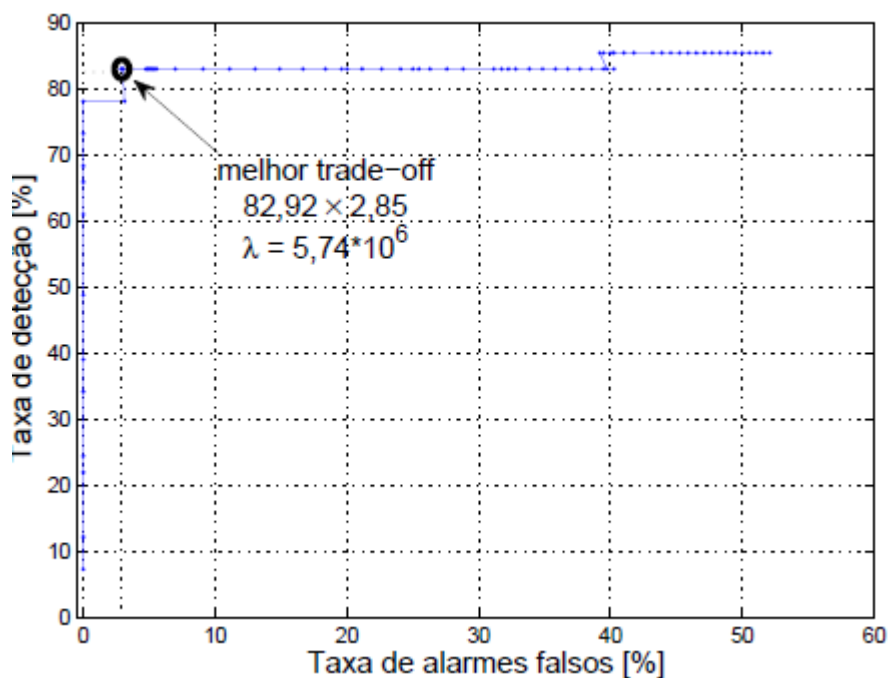


Figura 6.8 Alarmes disparados pelo SDA para o objeto *iflnOctets*, no dia 08/02/2010.



Foi definida uma classe de anomalias, denominada classe 1. Essa classe representa as anomalias de longa duração que excedem ou ficam abaixo do DSNS em pelo menos 60%. Os parâmetros utilizados na definição dessa classe são $\alpha\% = 60$ e $\gamma = 25$.

Figura 6.9 Resultados obtidos pelo SDA no dia 08/02/2010 no objeto *iflnOctets* do servidor web da UEL.



A fim de encontrar o valor ótimo de λ que resulte no melhor *trade-off* taxa de detecção \times taxa de falsos alarmes, o algoritmo de detecção foi executado 100 realizações para cada diferente valor de λ . A figura 6.9 apresenta o desempenho do sistema desenvolvido em termos deste *trade-off*, com λ variando no intervalo de $[1 \dots 10^7]$. O melhor valor encontrado, utilizando-se os parâmetros definidos na classe 1, foi $\lambda = 5,74 \times 10^6$. Os resultados confirmam que o SDA desenvolvido é aplicável para a detecção de anomalias em ambientes reais, utilizando-se as informações contidas nos objetos da MIB coletados através do protocolo SNMP. O melhor resultado, obtido foi de 82,92% para taxa de detecção e 2,85% para taxa de falsos alarmes, para $\lambda = 5,74 \times 10^6$.

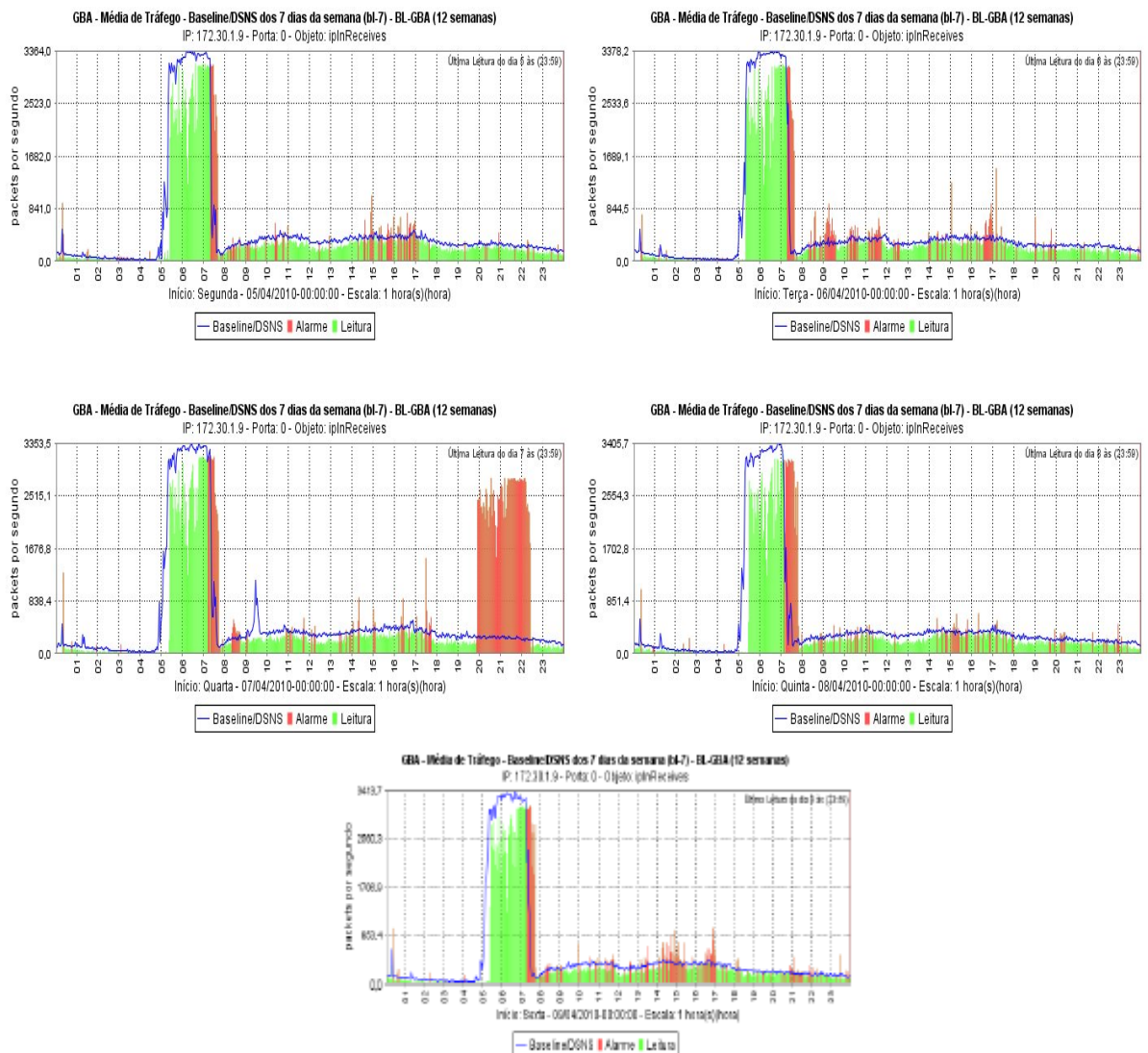
A figura 6.8 apresenta os alarmes gerados nesse período, levando em consideração os parâmetros definidos na classe 1, para $\lambda = 5,74 \times 10^6$. O eixo y representa a distância Euclidiana entre as amostras e os respectivos centróides, o eixo x

representa o horário do dia e a linha vermelha pontilhada representa o *threshold* λ . O melhor *trade-off* obtido foi de 82, 92% \times 2, 85% para taxa de detecção e taxa de falsos alarmes respectivamente⁶³

6.4 CENÁRIO 4

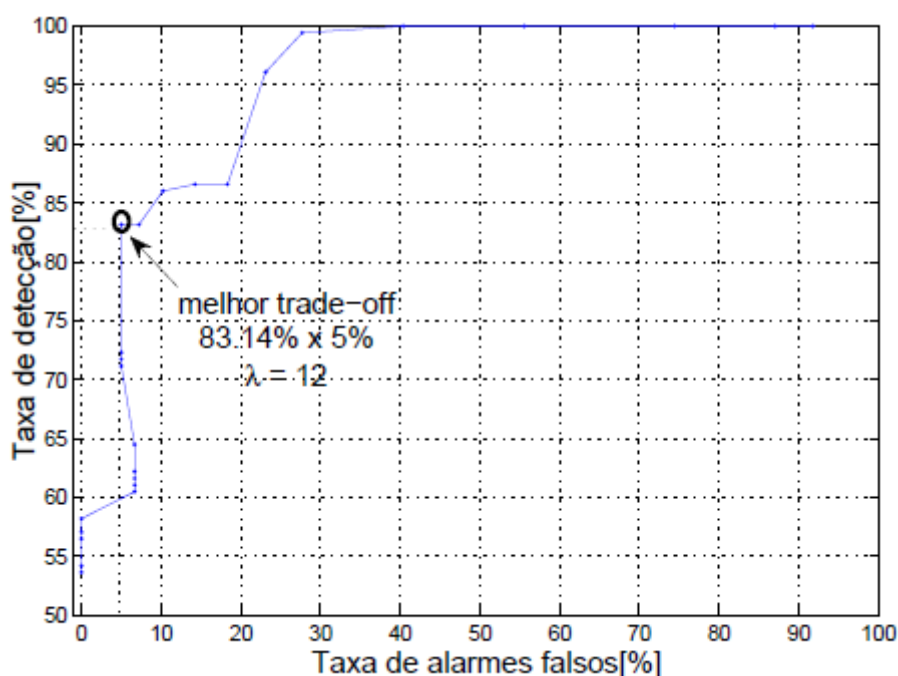
Neste cenário foi utilizado o tráfego de rede coletado no período de 05 – 09/04/2010, que representa os dias úteis da semana, como pode ser visto na figura 6.10. Os dados foram coletados no servidor web da UEL, do objeto *ipInReceives*, que contabiliza o número total de datagramas recebidos pela interface de rede, incluindo aqueles com erro.

Figura 6.10 Tráfego e DSNS dos dias 05 – 09/04/2010, do objeto *ipInReceives* do servidor web da UEL.



As métricas utilizadas para avaliação são a taxa de detecção eq. (6.1) e taxa de alarmes falsos eq. (6.2). Foram definidas duas classes de anomalias a serem detectadas. A classe 1, relacionada com anomalias de longa duração que se distanciam do DSNS (acima ou abaixo) em pelo menos 60%, e com duração de pelo menos 250 segundos (25 amostras). Os parâmetros utilizados para definir essa classe foram $\alpha = 60\%$ e $\gamma = 25$. A classe 2, refere-se às anomalias com menor duração, porém que excedam o DSNS em pelo menos 90%. Os parâmetros para essa classe são $\alpha = 90\%$ e $\gamma = 12$.

Figura 6.11 Desempenho do SDA utilizando-se a classe 1, para definição de anomalias, cenário 3b.



Com o objetivo de encontrar o valor ótimo de λ que resulte no melhor *trade-off* entre taxa de detecção \times taxa de falsos alarmes, o algoritmo de detecção foi executado 100 realizações para cada diferente valor de λ , para ambas as classes de anomalias. A figura 6.11 demonstra a precisão do sistema de detecção para a classe 1, com λ variando entre $[1 \dots 100]$, enquanto a figura 6.12 apresenta o desempenho para a classe 2. Os melhores valores obtidos através dos experimentos, foram $\lambda = 12$ para a classe 1 e $\lambda = 9$ para a classe 2. Os respectivos *trade-offs* foram $83,14\% \times 5\%$ para a classe 1, e $78\% \times 6,22\%$ para a classe 2.

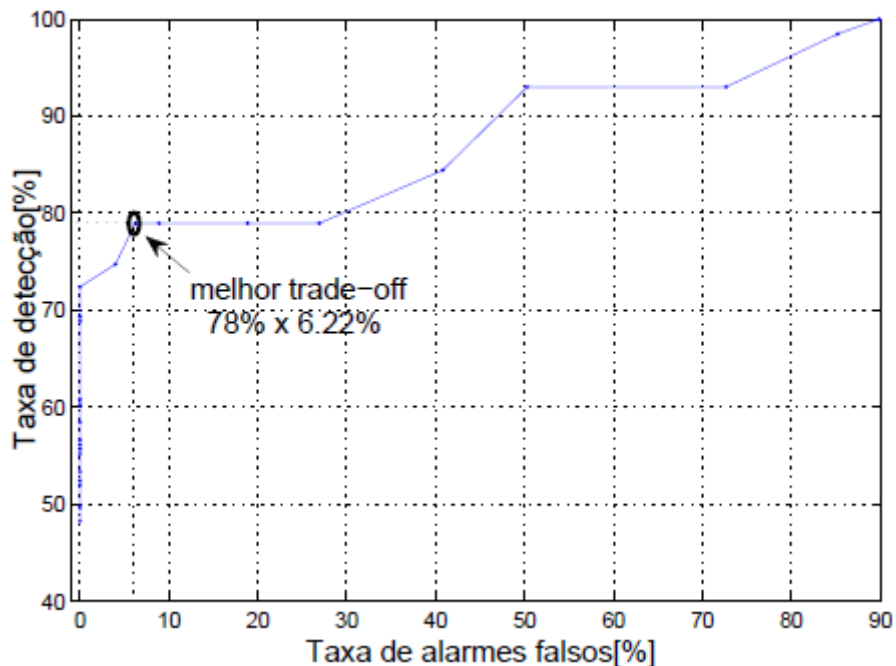
O SDA desenvolvido mostrou-se eficiente nesse experimento, atingindo resultados superiores a 78% para taxa de detecção, e taxa de alarmes falsos

abaixo de 3%. Além disso a implementação da parametrização do algoritmo de detecção, com a introdução de α e γ , possibilita uma precisão maior na detecção das anomalias, uma vez que as anomalias podem ser detectadas levando-se em consideração o conhecimento que o administrador de rede possui. A figura 6.13 apresenta os alarmes gerados para o período de 05 – 09/04/2010 da classe 1. Nesse cenário, as distâncias Euclidianas são normalizadas pela seguinte equação:

$$\tilde{\mathbf{D}} = \frac{\mathbf{D}}{\text{var}(\mathbf{D})} \quad (6.3)$$

onde D consiste das distâncias Euclidianas normalizadas, \mathbf{D} representa as distâncias Euclidianas, e $\text{var}(\mathbf{D})$ a variância (σ^2) de \mathbf{D} .

Figura 6.12 Desempenho do SDA utilizando-se a classe 2, para definição de anomalias, cenário 3b.



6.5 CENÁRIO 5

A fim de avaliar a precisão e desempenho do SDA desenvolvido que é baseado no algoritmo heurístico PSO-CIs, foi implementado um SDA baseado na análise de componentes principais (PCA), a fim de comparar os resultados obtidos por ambas as abordagens. Foram realizados três experimentos, todos utilizando tráfego coletado no período de 05/04/2010 a 11/04/2010 no servidor web da UEL: 1) Utiliza dados de apenas um objeto para o SDA proposto, *ipInReceives*; 2) Utiliza dados de quatro objetos para o SDA-PSO-CIs, *ipInReceives*, *ipInDelivers*, *tcpInSegs* e *ifInOctets*,

levando-se em consideração o grafo de dependência entre objetos da MIB, proposto em [42]; e 3) Utiliza dados de quatro objetos para o SDA-PCA, *ipInReceives*, *ipInDelivers*, *tcpInSegs* e *ifInOctets*. A figura 6.14 gerada pela ferramenta GBA, apresenta o tráfego e o DNS referente ao período utilizado nos testes.

Primeiramente foram realizados os experimentos para o SDA proposto, utilizando-se apenas um objeto da MIB. Foram definidos três parâmetros distintos para caracterizar as anomalias: $\gamma = 5$, $\gamma = 15$ e $\gamma = 25$. Os resultados obtidos são apresentados na66

Figura 6.13 Alarmes gerados para o período de 05 – 09/04/2010, do objeto *ipInReceives* no servidor web da UEL, cenário 3b.67

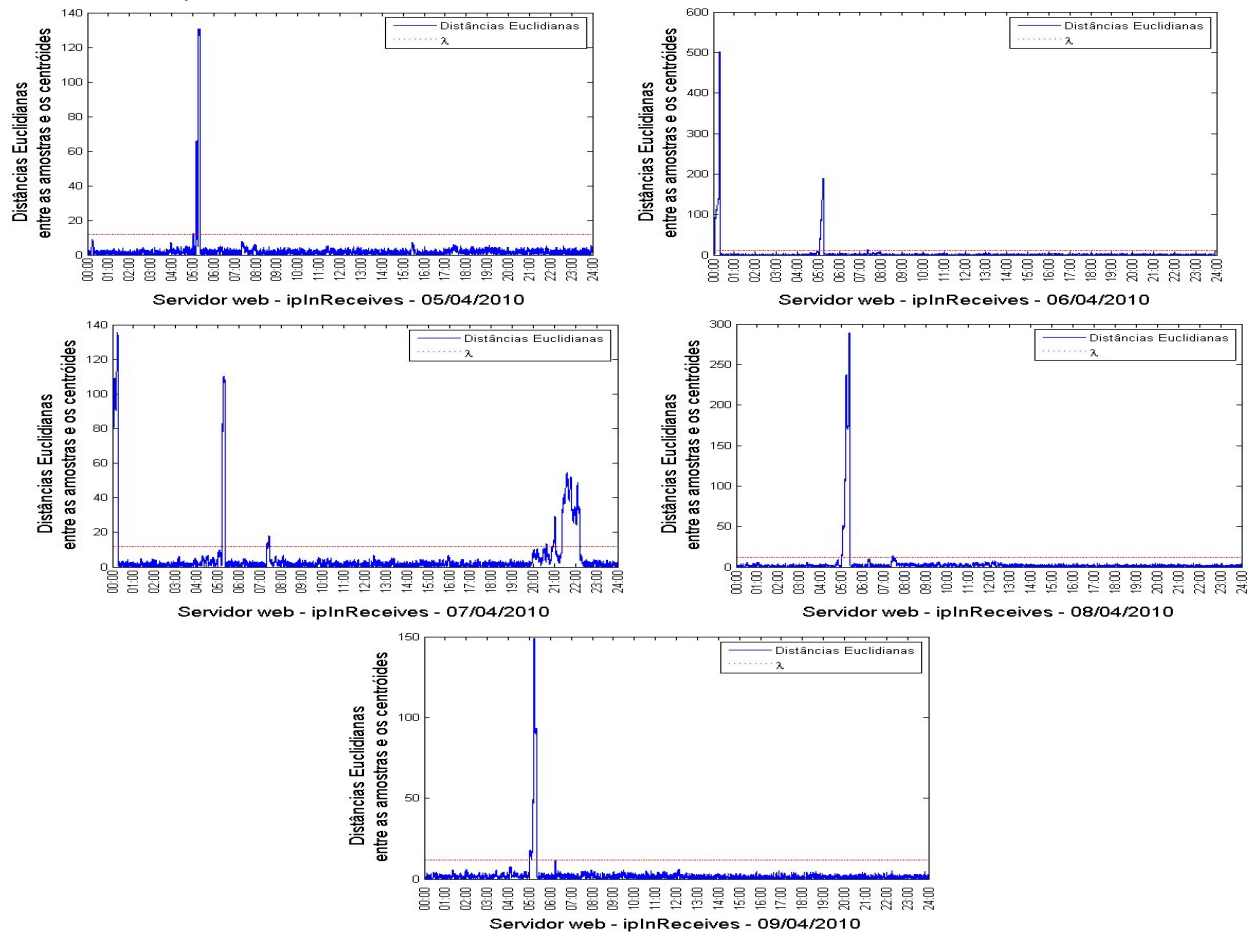
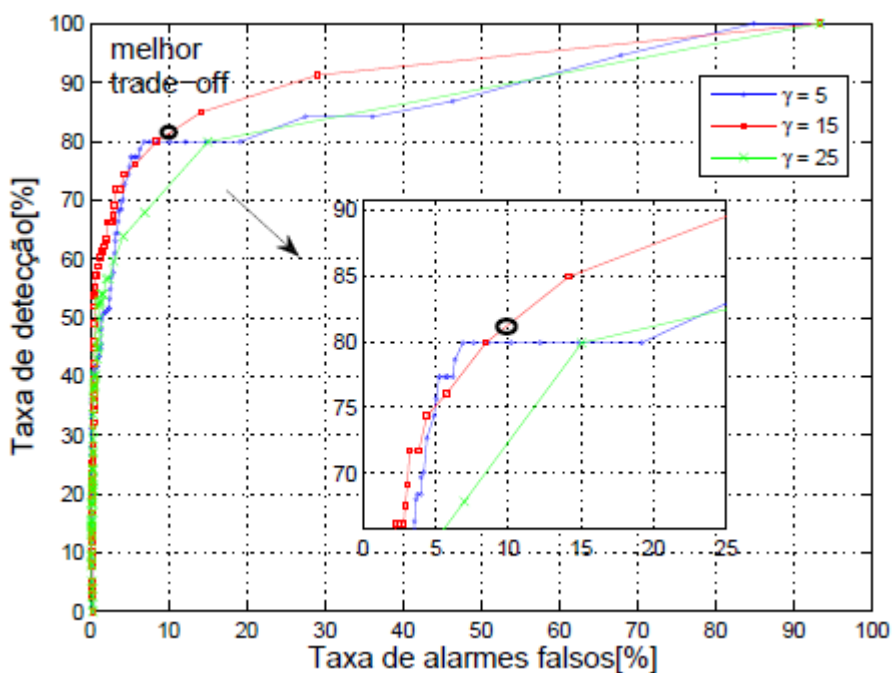


Figura 6.14 Gráficos da semana de tráfego do período de 05/04/2010 a 11/04/2010, objeto *ipInReceives*, cenário 468



figura 6.15. O melhor *trade-off* obtido dentre os três valores de γ foi de 82, 92% para taxa de detecção \times 9% para taxa de alarmes falso, para $\gamma = 15$. A Figura 6.16 apresenta os alarmes gerados para o objeto *ipInReceives*.

Figura 6.15 Taxa de detecção \times Taxa de alarmes falsos, para $\gamma = 5$, $\gamma = 15$ e $\gamma = 25$. Objeto *ipInReceives*. SDA PSO-CIs



Para $\gamma = 5$ e $\gamma = 15$, os melhores resultados obtidos foram 80% \times 7, 14% e 79, 80%, \times 15% respectivamente. Para $\gamma = 15$ o algoritmo apresentou uma performance superior, e ainda conseguiu atingir uma taxa de detecção de 91% com menos de 30% de taxa de alarmes falsos. Observou-se nesse experimento, que a detecção de anomalias, apresentou os melhores resultados para valores menores de γ , $\gamma = 5$ e $\gamma = 15$. Isso ocorre porque o intervalo utilizado para analisar o tráfego é de 300 segundos ($\gamma = 30$), de modo que quanto mais próximo o valor de γ for de 300 segundos, maior será a taxa de alarmes falsos.

No segundo experimento, foram utilizados quatro objetos simultâneos no SDA baseado no algoritmo PSO-CIs. A figura 6.17 apresenta os resultados obtidos para a semana de teste considerando $\gamma = 5$, $\gamma = 15$ e $\gamma = 25$. Através desse experimento, observa-se que os resultados para todos os valores de γ , ficaram mais semelhantes, uma vez que uma maior quantidade de dados foi utilizado. O melhor resultado obtido foi de 78% \times 9% para $\gamma = 15$. Apesar desse experimento não ter

apresentado melhora nos resultados, a correlação dos objetos da MIB é um caminho que deve ser investigado mais profundamente, em busca de melhores *trade-offs*. A figura 6.18 apresenta os alarmes disparados para a semana de teste, utilizando quatro objetos simultâneos

Figura 6.16 Alarmes disparados pelos ADS baseado no PSO-CIs no período de 05/04/2010 a 11/04/2010, objeto *ipInReceives*

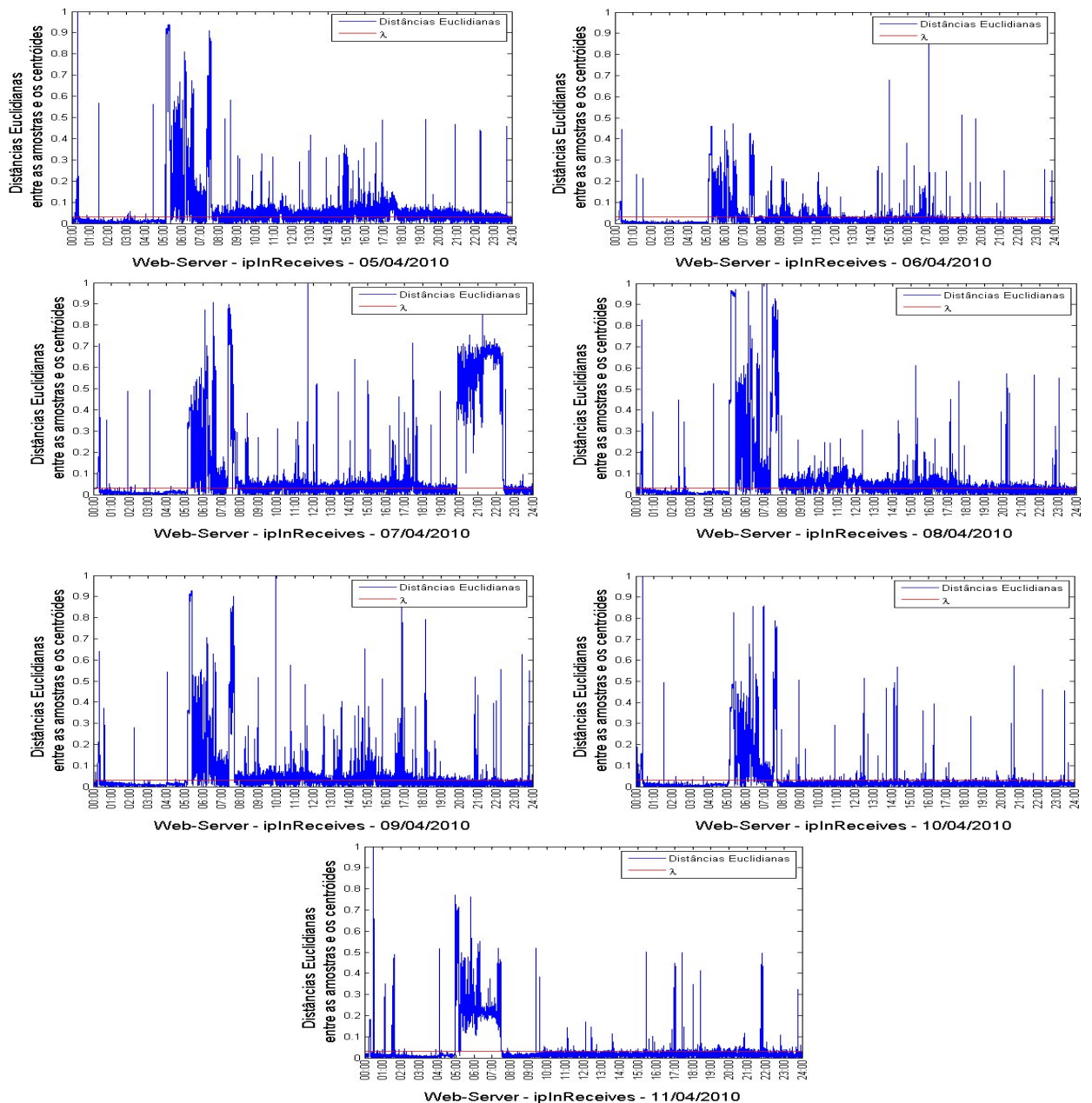
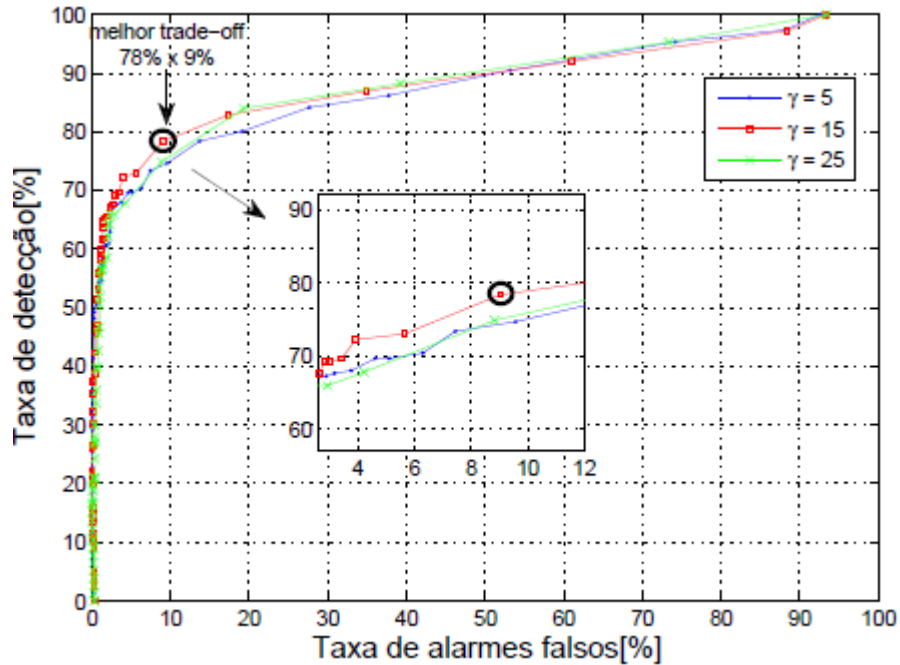


Figura 6.17 Taxa de detecção × Taxa de alarmes falsos, para $\gamma = 5$, $\gamma = 15$ e $\gamma = 25$. Quatro objetos SNMP simultâneos. SDA baseado no algoritmo PSO-CIs



O terceiro experimento foi realizado a fim avaliar o desempenho do SDA baseado no PCA. Assim como no experimento anterior, foram utilizados quatro objetos da MIB para a realização da etapa de treinamento do modelo normal, *ifInOctets*, *ipInReceives*, *ipInDelivers* e *tcpInSegs*. Foi utilizado o DSNS de cada objeto nessa etapa, enquanto na etapa de detecção, utilizou-se o movimento dos objetos.

Figura 6.18 Alarmes gerados pelo SDA baseado no algoritmo PSO-CIs, no período de 05/04/2010 a 11/04/2010, para os objetos *iplnReceives*, *iplnDelivers*, *tcpInSegs* e *ifInOctets*.

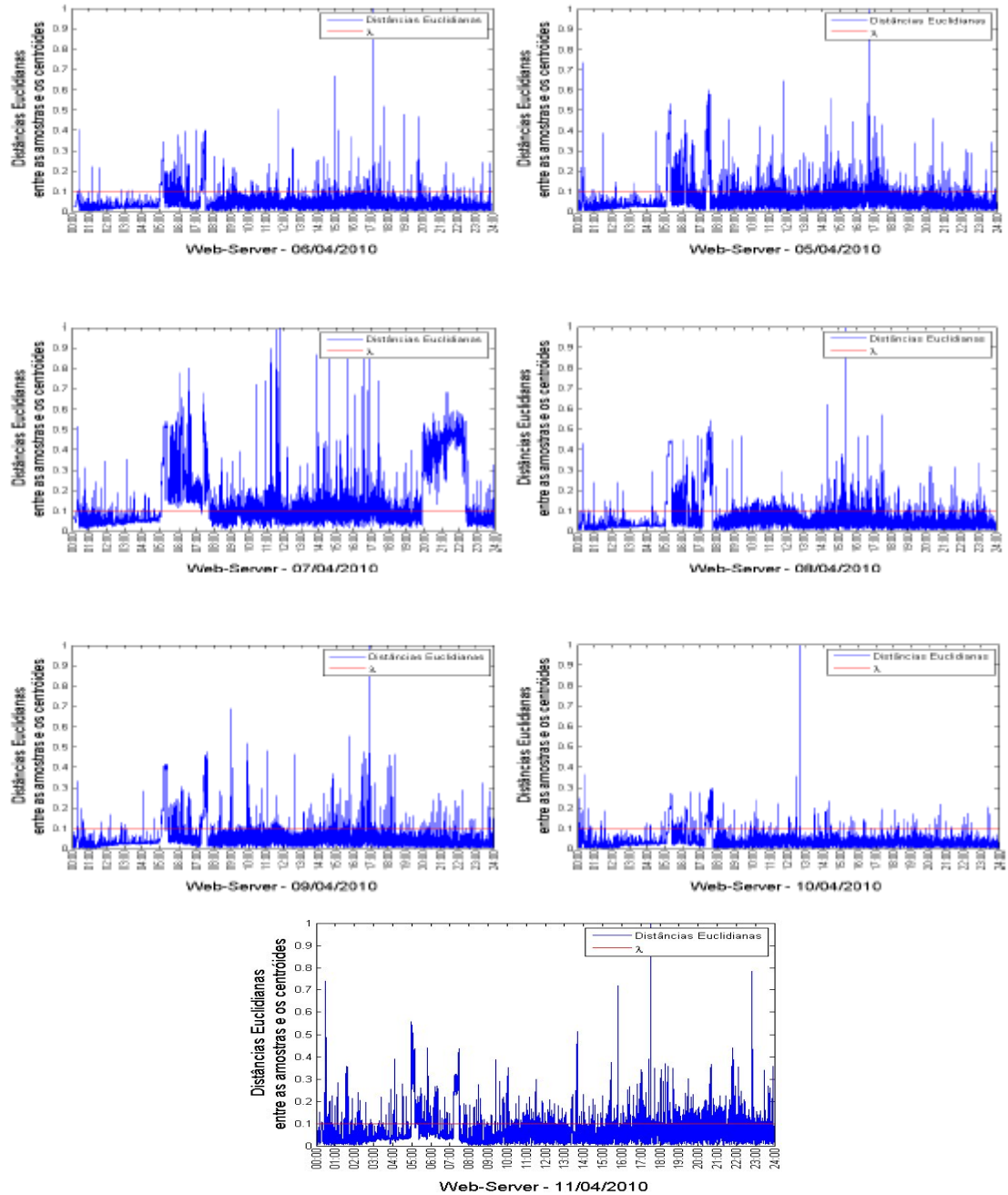
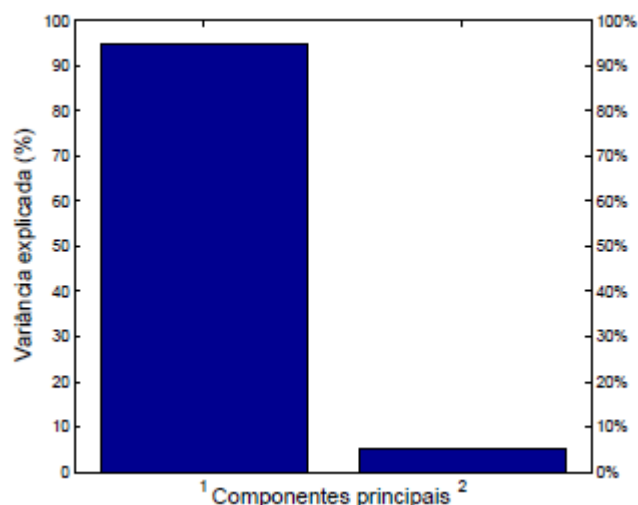


Figura 6.19 Variância percentual explicada pelos componentes principais

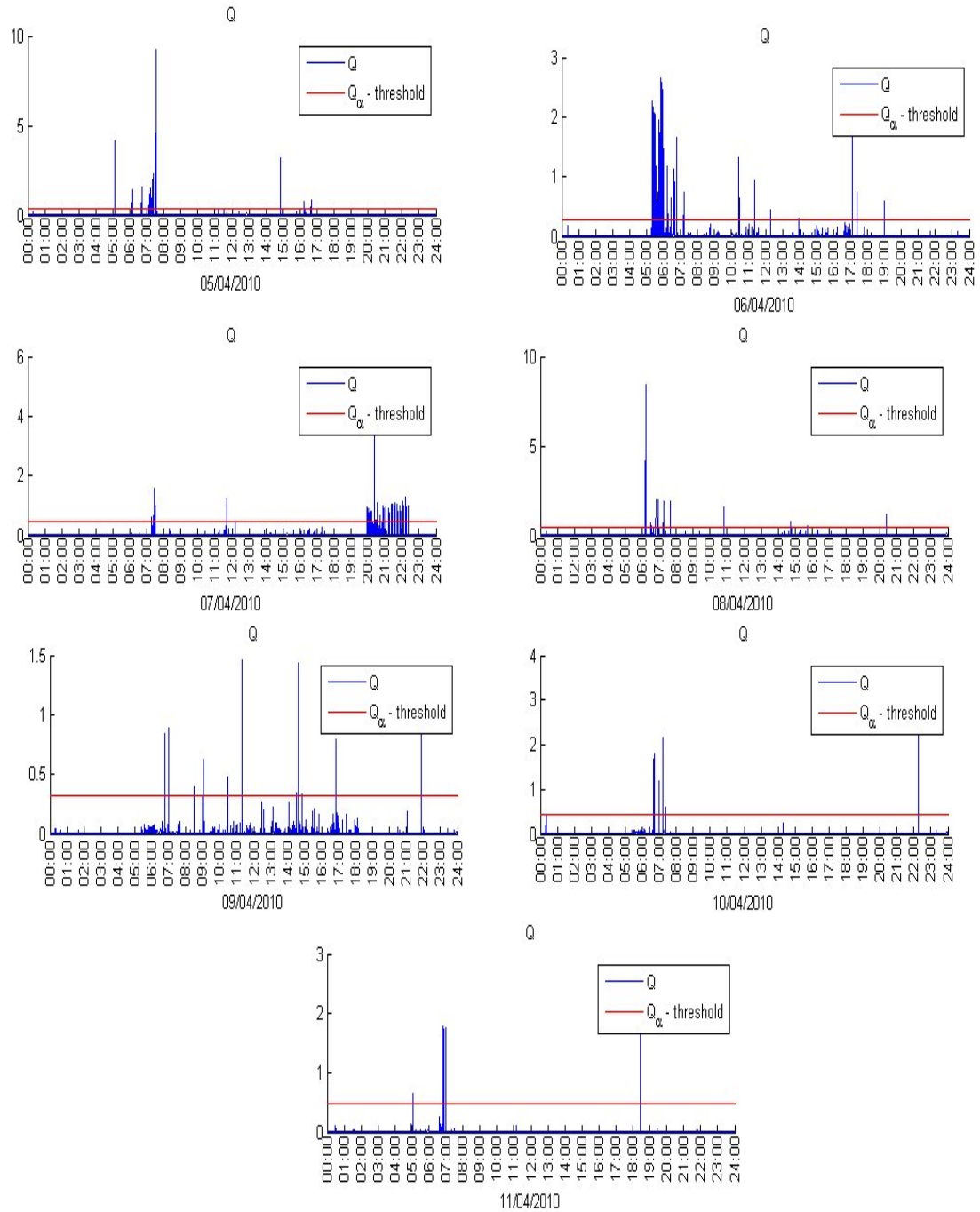


A figura 6.19 demonstra que mais de 90% da variação apresentada no conjunto de dados de treinamento, composto pelo DSNS dos quatro objetos, pode ser representada por apenas um componente principal. Embora o segundo componente apresente cerca de 5% da variabilidade, a abordagem utilizada neste trabalho, o CPV 5.7 considera somente os componentes que representem 90% da variabilidade. Deste modo, o conjunto de dados de treinamento é projetado sobre apenas o primeiro componente principal. A detecção de anomalias é obtida através do cálculo de Q , que utiliza o conjunto de dados composto pelo movimento dos quatro objetos. A performance obtida pelo modelo baseado no PCA foi de 70, 60% \times 1, 19% para taxa de detecção e alarmes falsos respectivamente. A figura 6.20 apresenta os alarmes gerados pelo modelo, em cada um dos dias da semana de teste.

Os resultados obtidos através dos experimentos neste cenário, demonstraram que o SDA desenvolvido é mais eficiente do que a abordagem utilizando o PCA, atingindo um *trade-off* de 82, 92% \times 9% para um objeto da MIB, e 78% \times 9% para quatro objetos simultâneos, contra 70, 60% \times 1, 19% obtidos pela abordagem PCA, também utilizando quatro objetos. Outro fator relevante, e que confirma nossa abordagem como superior para os cenários utilizados, é a análise de complexidade, vista na Seção 5.1. Contudo, é importante ressaltar que não foi realizado um estudo sobre a otimização dos parâmetros do PCA, o que deverá ser realizado de forma cautelosa nos trabalhos futuros. Apesar dos resultados obtidos pelo PCA terem sido ligeiramente inferiores e de sua maior complexidade, a técnica é promissora pois realiza a redução da dimensionalidade de grandes conjuntos de dados, tendo vasta aplicação no

cenário de redes atual. O PCA pode ainda ser combinado com outras técnicas, como o PSO ou o K-means, dando origem a uma variedade de novos algoritmos.

Figura 6.20 Alarmes gerados pelo SDA baseado no PCA.74



7 CONCLUSÕES

Neste trabalho foi abordado o problema de detecção de anomalias em tráfego de redes de comunicação, através do monitoramento de objetos da *Management Information Base* (MIB). O domínio de aplicação utilizado foi a rede da Universidade Estadual de Londrina (UEL), que atualmente possui mais de 5000 ativos, dentre computadores e equipamentos, que geram uma grande quantidade de dados que precisam ser monitorados em tempo real com o objetivo de manter a disponibilidade e a qualidade dos sistemas de ensino e pesquisa da universidade. Grande parte das soluções de monitoramento existente atualmente apóia-se em operadores de rede que são responsáveis pela configuração manual de *thresholds*, e pelo seu monitoramento, o que consiste de uma atitude pouco eficiente e que não explora a natureza intrínseca da correlação entre os dados. Portanto, a automação dos sistemas de monitoramento de redes consiste de uma ampla área de desenvolvimento e pesquisas.

Visando otimizar o monitoramento e a qualidade de redes de comunicação, este trabalho inspirou-se nos resultados promissores obtidos pelas pesquisas de Proença e Zarpelão [7, 24, 28, 42], que consideram a utilização dos dados presentes nos objetos da MIB para a análise do tráfego em busca de anomalias, tendo como base a utilização do DSNS [7], que conta com uma vasta base de dados histórica para a geração de perfis da rede da Universidade Estadual de Londrina. Esses trabalhos são baseados na aplicação de modelos estatísticos e determinísticos, que consistem de técnicas clássicas na literatura as quais vêm sendo amplamente utilizadas durante os últimos anos na área detecção de anomalias.

Com base na revisão bibliográfica realizada, nos trabalhos recentes encontrados na literatura, e.g., [1, 2, 19, 43], tem sido observado uma grande quantidade de pesquisas envolvendo a aplicação de técnicas heurísticas como o PSO e outras mais antigas, que agora têm sido exploradas como ferramentas para a detecção de anomalias. A alta dimensionalidade dos dados decorrente da complexidade dos sistemas e das redes de comunicação, em conjunto com a diversidade de fontes de dados e as freqüentes mudanças no comportamento padrão da rede, fazem com que os modelos heurísticos sejam mais eficientes em relação aos modelos determinísticos, devido à sua natureza em resolver problemas⁷⁵ complexos e de alta dimensionalidade, com menor complexidade computacional. É com base nesse estudo que desenvolvemos um sistema de detecção de anomalias (SDA), que utiliza uma

abordagem heurística combinada com a utilização do DSNS para a detecção de anomalias.

7.1 CONTRIBUIÇÕES

Com base nos aspectos levantados durante a revisão da literatura, foi desenvolvido um sistema de detecção de anomalias que combina a utilização de três algoritmos. Seu objetivo é a detecção de anomalias de volume, através do monitoramento de objetos da MIB via protocolo SNMP. O sistema desenvolvido, é baseado no algoritmo PSO-CIs que consiste da combinação dos algoritmos *K-means* e Particle Swarm Optimization (PSO), que são aplicados sobre o movimento do tráfego e o DSNS. A detecção de anomalias se dá através da análise da distância entre os dados clusterizados e os centróides dos clusters. A fim de avaliar o sistema desenvolvido, foram realizados estudos de complexidade, otimização dos parâmetros e comparação com outros modelos.

Foram realizados diversos experimentos, sob a perspectiva de diferentes cenários, com objetivo de avaliar o desempenho e a precisão do SDA desenvolvido. O sistema foi testado através da aplicação de diferentes objetos SNMP, tais como o *ipInReceives*, *ipInDelivers*, *ifInOctets* e *tcpInSegs*. Foram avaliados os resultados obtidos através da aplicação do sistema para objetos únicos, bem como para objetos monitorados simultaneamente. Os dados utilizados nos experimentos foram coletados no ambiente de rede da Universidade Estadual de Londrina, em diferentes períodos de tempo.

Os resultados obtidos através dos experimentos demonstraram que o DAS desenvolvido é capaz de detectar anomalias de volume de forma eficaz, alcançando resultados bastante satisfatórios e com uma baixa complexidade se comparado com outros modelos. Os testes resultaram em taxas de detecção superiores a 80% enquanto as taxas de alarmes falsos não ultrapassaram a faixa dos 10%. Os resultados foram comparados com os obtidos pela aplicação de um algoritmo determinístico, e também com os da abordagem utilizando o PCA, como apresentado no capítulo 6. Os resultados foram bastante satisfatórios uma vez que nosso sistema, baseado no algoritmo PSO-CIs, se mostrou mais eficiente obtendo desempenho superior em até 20% para taxa de detecção. Os resultados obtidos através da aplicação do SDA baseado na abordagem de PCA [38][44][45], também mostrou-se bastante promissora, obtendo

ótimos resultados. Contudo, a comparação dos resultados entre o SDA baseado no PSO-CIs e a abordagem PCA, demonstraram um ganho de quase 10% na taxa de detecção,76 para o algoritmo PSO-CIs, levando em consideração que não foi realizado uma análise de otimização dos parâmetros da abordagem PCA, o que deverá ser realizado futuramente.

Os trabalhos futuros incluem o aperfeiçoamento da análise simultânea dos objetos da MIB, a fim de reduzir o ruído apresentado nas distâncias Euclidianas, e com isso obter uma redução na taxa de alarmes falsos. Também pretende-se, incluir um módulo para análise e classificação das anomalias, e também o desenvolvimento de um modelo para análise do tráfego baseado em pacotes IP, o que possibilitaria uma investigação mais profunda a respeito das anomalias de rede. A utilização do PCA para redução de dimensionalidade dos dados, também pode vir a ser combinada com o SDA desenvolvido, quando o número de objetos analisados simultaneamente crescer demasiadamente.

REFERÊNCIAS

- 1 PATCHA, A.; PARK, J.-M. An overview of anomaly detection techniques: Existing solutions and latest technological trends. *Comput. Netw.*, Elsevier North-Holland, Inc., New York, NY, USA, v. 51, n. 12, p. 3448–3470, 2007. ISSN 1389-1286.
- 2 CHANDOLA, V.; BANERJEE, A.; KUMAR, V. Anomaly detection: A survey. *ACM Comput. Surv.*, ACM, New York, NY, USA, v. 41, n. 3, p. 1–58, 2009. ISSN 0360-0300.
- 3 HODGE, V.; AUSTIN, J. A survey of outlier detection methodologies. *Artif. Intell. Rev.*, v. 22, n. 2, p. 85–126, 2004.
- 4 YANKOV, D.; KEOGH, E.; REBBAPRAGADA, U. Disk aware discord discovery: finding unusual time series in terabyte sized datasets. *Knowl. Inf. Syst.*, v. 17, n. 2, p. 241–262, 2008.
- 5 KUANG, L.; ZULKERNINE, M. An anomaly intrusion detection method using the csi-knn algorithm. In: *Proceedings of the 2008 ACM symposium on Applied computing*. New York, NY, USA: ACM, 2008. (SAC '08), p. 921–926. ISBN 978-1-59593-753-7. Disponível em: <<http://doi.acm.org/10.1145/1363686.1363897>>.
- 6 MCCLOGHRIE, K.; ROSE, M. T. *Management Information Base for Network Management of TCP/IP-based internets:MIB-II*. United States: RFC Editor, 1991.
- 7 PROENÇA, M. et al. The hurst parameter for digital signature of network segment. In: SOUZA, J. de; DINI, P.; LORENZ, P. (Ed.). *Telecommunications and Networking – ICT 2004*. Springer Berlin / Heidelberg, 2004, (Lecture Notes in Computer Science, v. 3124). p. 772–781. Disponível em: <http://dx.doi.org/10.1007/978-3-540-27824-5_103>.
- 8 XIAO, L.; SHAO, Z.; LIU, G. K-means algorithm based on particle swarm optimization algorithm for anomaly intrusion detection. In: *WCICA 2006 . The Sixth World Congress on Intelligent Control and Automation*. [S.l.: s.n.], 2006. p. 5854 – 5858.
- 9 KENNEDY, J.; EBERHART, R. C. *Swarm intelligence*. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc., 2001. ISBN 1-55860-595-9.
- 10 LIM, S. Y.; JONES, A. Network anomaly detection system: The state of art of network behaviour analysis. In: *ICHIT '08: Proceedings of the 2008 International Conference on Convergence and Hybrid Information Technology*. Washington, DC, USA: IEEE Computer Society, 2008. p. 459–465. ISBN 978-0-7695-3328-5.
- 11 TAVALLAEE, M.; STAKHANOVA, N.; GHORBANI, A. A. Toward credible evaluation of anomaly-based intrusion-detection methods. *Systems, Man, and Cybernetics, Part C: Applications and Reviews, IEEE Transactions on*, PP, n. 99, p. 1 –9, 2010. ISSN 1094-6977.
- 12 GADDAM, S. R.; PHOHA, V. V.; BALAGANI, K. S. K-means+id3: A novel method for supervised anomaly detection by cascading k-means clustering and id3 decision

- tree learning methods. *Knowledge and Data Engineering, IEEE Transactions on*, p. 345–354, 2007.78
- 13 ZHANG, C.; ZHANG, G.; SUN, S. A mixed unsupervised clustering-based intrusion detection model. In: . Los Alamitos, CA, USA: IEEE Computer Society, 2009. v. 0, p. 426–428. ISBN 978-0-7695-3899-0.
 - 14 JIANLIANG, M.; HAIKUN, S.; LING, B. The application on intrusion detection based on k-means cluster algorithm. In: *Proceedings of the 2009 International Forum on Information Technology and Applications - Volume 01*. Washington, DC, USA: IEEE Computer Society, 2009. p. 150–152. ISBN 978-0-7695-3600-2. Disponível em: <<http://portal.acm.org/citation.cfm?id=1606748.1606787>>.
 - 15 ENSAFI, R.; DEGHANZADEH, S.; T, M.-R. A. Optimizing fuzzy k-means for network anomaly detection using pso. In: *Proceedings of the 2008 IEEE/ACS International Conference on Computer Systems and Applications*. Washington, DC, USA: IEEE Computer Society, 2008. (AICCSA '08), p. 686–693. ISBN 978-1-4244-1967-8. Disponível em: <<http://dx.doi.org/10.1109/AICCSA.2008.4493603>>.
 - 16 LIU, L. li; LIU, Y. Mqpso based on wavelet neural network for network anomaly detection. In: *Wireless Communications, Networking and Mobile Computing, 2009. WiCom '09. 5th International Conference on*. [S.l.: s.n.], 2009. p. 1–5.
 - 17 MA, R. et al. Network anomaly detection using rbf neural network with hybrid qpso. In: *IEEE International Conference on Networking, Sensing and Control*. [S.l.: s.n.], 2008.
 - 18 GAO, X.; OVASKA, S.; WANG, X. Particle swarm optimization of detectors in negative selection algorithm. In: *Systems, Man and Cybernetics, 2007. ISIC. IEEE International Conference on*. [S.l.: s.n.], 2007. p. 1236–1242.
 - 19 JIANZHEN, W.; JINRONG, S. Research on the application of particle swarm optimization algorithm in anomaly detection. In: *Computer Science and Education (ICCSE), 2010 5th International Conference on*. [S.l.: s.n.], 2010. p. 474–476.
 - 20 BRAUCKHOFF, D.; SALAMATIAN, K.; MAY, M. Applying pca for traffic anomaly detection: Problems and solutions. In: *INFOCOM 2009, IEEE*. [S.l.: s.n.], 2009. p. 2866–2870. ISSN 0743-166X.
 - 21 CALLEGARI, C. et al. A novel multi time-scales pca-based anomaly detection system. In: *Performance Evaluation of Computer and Telecommunication Systems (SPECTS), 2010 International Symposium on*. [S.l.: s.n.], 2010. p. 156–162.
 - 22 CASAS, P. et al. Volume anomaly detection in data networks: An optimal detection algorithm vs. the pca approach. In: . Berlin, Heidelberg: Springer-Verlag, 2009. p. 96–113. ISBN 978-3-642-04575-2. Disponível em: <<http://portal.acm.org/citation.cfm?id=1612488.1612497>>.
 - 23 LEE, D. C. et al. Fast traffic anomalies detection using snmp mib correlation analysis. In: *Proceedings of the 11th international conference on Advanced Communication Technology - Volume 1*. Piscataway, NJ, USA: IEEE Press, 2009.

- (ICACT'09), p. 166–170. ISBN 978-8- 9551-9138-7. Disponível em: <<http://portal.acm.org/citation.cfm?id=1701955.1701986>>.
- 24 ZARPELÃO, B. B. et al. Parameterized anomaly detection system with automatic configuration. In: *Proceedings of the 28th IEEE conference on Global telecommunications*. Piscataway, NJ, USA: IEEE Press, 2009. (GLOBECOM'09), p. 2224–2229. ISBN 978-1-4244- 4147-1. Disponível em: <<http://portal.acm.org/citation.cfm?id=1811681.1811749>>.
- 25 THE third international knowledge discovery and data mining tools competition data set KDD99-Cup. Disponível em: <<http://kdd.ics.uci.edu/databases/kddcup99/kddcup99.html>>. Acesso em: 25 nov. 2011.
- 26 LINCOLN Laboratory, MIT. DARPA intrusion detection data sets, 2009. Disponível em: <<http://kdd.ics.uci.edu/databases/kddcup99/kddcup99.html>>. Acesso em: 25 nov. 2011.
- 27 ABILENE observatory data collections. Disponível em: <www.internet2.edu/observatory>. Acesso em: 25 nov. 2011.
- 28 PROENÇA, M. et al. Baseline to help with network management. In: ASCENSO, J. et al. (Ed.). *e-Business and Telecommunication Networks*. Springer Netherlands, 2006. p. 158–166. ISBN 978-1-4020-4761-9. Disponível em: <http://dx.doi.org/10.1007/1-4020-4761-4_12>.
- 29 Proença Jr., M. L.; ROCHOL, J. Uma ferramenta para auxílio no gerenciamento de redes com backbone atm. In: *V CONGRESO ARGENTINO DE CIENCIAS DE LA COMPUTACIÓN - CACIC'99*. [S.l.: s.n.], 1999.
- 30 CASE, J. D. et al. *Simple Network Management Protocol (SNMP)*. United States: RFC Editor, 1990.
- 31 CASE, J. D. et al. *Introduction to Community-based SNMPv2*. United States: RFC Editor, 1996.
- 32 CASE, J. D. et al. *Introduction to Version 3 of the Internet-standard Network Management Framework*. United States: RFC Editor, 1999.
- 33 FIROUZI, B.; NIKNAM, T.; NAYERIPOUR, M. A new evolutionary algorithm for cluster analysis. In: *International Journal of Computer Science*. [S.l.: s.n.], 2009.
- 34 EBERHART, R.; KENNEDY, J. A new optimizer using particle swarm theory. In: *Micro Machine and Human Science, 1995. MHS '95., Proceedings of the Sixth International Symposium on*. [S.l.: s.n.], 1995. p. 39–43.
- 35 NEDJAH, N.; MOURELLE, L. M. *Swarm Intelligent Systems*. Springer-Verlag Berlin Heidelberg: Springer, 2006.
- 36 CHATTERJEE, A.; SIARRY, P. Nonlinear inertia weight variation for dynamic adaptation in particle swarm optimization. *Comput. Oper. Res.*, Elsevier Science Ltd., Oxford, UK, UK, v. 33, p. 859–871, March 2006. ISSN 0305-0548. Disponível em: <<http://portal.acm.org/citation.cfm?id=1115087.1115140>>.

- 37 EBERHART; SHI, Y. Particle swarm optimization: developments, applications and resources. In: *Evolutionary Computation, 2001. Proceedings of the 2001 Congress on*. [S.l.: s.n.], 2001.
- 38 LAKHINA, A.; CROVELLA, M.; DIOT, C. Diagnosing network-wide traffic anomalies. In: *Proceedings of the 2004 conference on Applications, technologies, architectures, and protocols for computer communications*. New York, NY, USA: ACM, 2004. (SIGCOMM '04), p. 219–230. ISBN 1-58113-862-8. Disponível em: <<http://doi.acm.org/10.1145/1015467.1015492>>.80
- 39 JACKSON, J. E. *A user's guide to principal components*. 1. ed. [S.l.]: Wiley-Interscience, 1991. ISBN 0471622672.
- 40 ZUMOFFEN, D.; BASUALDO, M. From large chemical plant data to fault diagnosis integrated to decentralized fault-tolerant control: Pulp mill process application. *Industrial and Engineering Chemistry Research*, v. 47, n. 4, p. 1201–1220, 2008.
- 41 AXELSSON, S. The base-rate fallacy and the difficulty of intrusion detection. *ACM Trans. Inf. Syst. Secur.*, ACM, New York, NY, USA, v. 3, p. 186–205, August 2000. ISSN 1094-9224. Disponível em: <<http://doi.acm.org/10.1145/357830.357849>>.
- 42 ZARPELÃO, B. B.; MENDES, L. D. S.; PROENÇA JR., M. L. Anomaly detection aiming pro-active management of computer network based on digital signature of network segment. *J. Netw. Syst. Manage.*, Plenum Press, New York, NY, USA, v. 15, p. 267–283, June 2007. ISSN 1064-7570. Disponível em: <<http://portal.acm.org/citation.cfm?id=1272213.1272221>>.
- 43 PUKKAWANNA, S.; FUKUDA, K. Combining sketch and wavelet models for anomaly detection. In: *Intelligent Computer Communication and Processing (ICCP), 2010 IEEE International Conference on*. [S.l.: s.n.], 2010. p. 313 –319.
- 44 ALVAREZ, D. G. Fault detection using principal component analysis (pca) in a wastewater treatment plant (wwtp). In: *Proceedings of 62-th International Student's Scientific Conference*. [S.l.: s.n.], 2009.
- 45 RINGBERG, H. et al. Sensitivity of pca for traffic anomaly detection. In: *Proceedings of the 2007 ACM SIGMETRICS international conference on Measurement and modeling of computer systems*. New York, NY, USA: ACM, 2007. (SIGMETRICS '07), p. 109–120. ISBN 978-1-59593-639-4. Disponível em: <<http://doi.acm.org/10.1145/1254882.1254895>>81

TRABALHOS PUBLICADOS PELO AUTOR

- 1 Networking Anomaly Detection using DSNS and Particle Swarm Optimization with Re- Clustering. Moises F. Lima, Lucas D. H. Sampaio, Bruno B. Zarpelão, Joel J. P. C. Rodrigues, Taufik Abrão, Mario Lemes Proença Jr; IEEE GLOBAL COMMUNICATIONS CONFERENCE (IEEE GLOBECOM 2010), December 2010, Miami, USA.

- 2 Anomaly detection using baseline and K-means clustering. Moises F. Lima, Bruno B. Zarpelão, Lucas D. H. Sampaio, Joel J. P. C. Rodrigues, Taufik Abrão, Mario Lemes Proença Jr; 18th International Conference on Software, Telecommunications and Computer Networks (SoftCOM 2010), IEEE Communication Society (ComSoc), 2010, Split, Hvar, Korcula setembro, 2010.
- 3 A heuristic-based network anomaly detection system using Digital Signature of Network Segment. Moisés F. Lima, Lucas D. H. Sampaio, Joel J. P. C. Rodrigues, Taufik Abrão, Mario Lemes Proença Jr.; Journal of Network and Systems Management, Springer, 2011. (submetido para avaliação)