

ARTIGO COMPLETO/FULL PAPER

Otimização da Coleta de Dados BGP: Seleção de Coletores Chave

Optimizing BGP Data Collection: Selection of Key Collectors

Pedro Maciel Carneiro Ferreira • ✉ pedroferreira@furg.br
Universidade Federal do Rio Grande (FURG)

Pedro De Botelho Marcos • ✉ pbmarcos@furg.br
Universidade Federal do Rio Grande (FURG)

RESUMO. Este trabalho investiga redundâncias na coleta de dados BGP e propõe um algoritmo para sugerir subconjuntos eficazes de coletores. A análise revela que diferentes coletores oferecem melhores resultados para distintas métricas: alguns são mais adequados para ASes, enquanto outros se destacam na coleta de prefixos, enlaces, espaço de endereçamento ou comunidades BGP. Utilizando dados dos projetos Route Views, RIPE RIS e PCH, a abordagem otimiza recursos computacionais sem comprometer a visibilidade da rede. Além disso, foi avaliada a consistência dos coletores mais eficazes ao longo do tempo, verificando se mantêm seu desempenho em diferentes períodos.

ABSTRACT. This paper investigates redundancies in BGP data collection and proposes an algorithm to suggest effective subsets of collectors. The analysis shows that different collectors perform better for specific metrics: some are more suitable for ASes, while others excel in collecting prefixes, links, address space or BGP communities. Using data from the Route Views, RIPE RIS, and PCH projects, the approach optimizes computational resources without compromising network visibility. Additionally, the consistency of the most effective collectors over time was evaluated to verify whether they maintain their performance across different periods.

PALAVRAS-CHAVE: BGP • Internet • Prefixos • Roteamento • Sistemas Autônomos

KEYWORDS: BGP • Internet • Prefixes • Routing • Autonomous Systems

1 Introdução

A Internet é composta por milhares de Sistemas Autônomos (ASes) interconectados que permitem a troca de dados globalmente [1]. Esses ASes se conectam por enlaces que definem as rotas, garantindo que pacotes alcancem seus destinos.

Para monitorar e compreender a infraestrutura e o comportamento da Internet, plataformas de coleta de dados do *Border Gateway Protocol* (BGP), como Route Views (RV) [2], RIPE Routing Information Service (RIPE RIS) [3] e Packet Clearing House (PCH) [4], mantêm sessões de *peering* com roteadores. Esses roteadores compartilham voluntariamente (total ou parcialmente) suas tabelas de roteamento, permitindo que as plataformas colem dados essenciais, como prefixos anunciados, topologias de ASes e conectividade entre diferentes partes da rede. Essas informações são fundamentais para análises científicas e operacionais da infraestrutura da Internet.

No entanto, o protocolo BGP propaga amplamente mensagens de conectividade, o que leva múltiplos coletores a registrarem informações semelhantes, gerando uma grande redundância de dados. Essa redundância resulta em volumes massivos de dados, aumentando

os custos de processamento e armazenamento sem, necessariamente, agregar novas informações relevantes. Além disso, a concentração de dados duplicados pode gerar lacunas de visibilidade, já que informações específicas ou relevantes podem ser obscurecidas em meio à repetição [5].

Este trabalho propõe um algoritmo de seleção otimizada de coletores, sugerindo subconjuntos eficazes para diferentes tipos de análise. A abordagem visa alcançar alta cobertura de ASes, enlaces, prefixos, espaço de endereçamento e comunidades BGP, sem a necessidade de processar todos os dados coletados, minimizando redundâncias e otimizando o uso dos recursos computacionais. Com essa otimização, analistas e operadores de rede podem concentrar-se nas fontes de dados mais relevantes, evitando o desperdício de esforço com dados excessivamente redundantes. Adicionalmente, foi adotada a normalização do espaço de endereçamento em blocos /24 para IPv4 e /48 para IPv6.

A eficácia do algoritmo foi testada ao longo do tempo, com o objetivo de verificar se os coletores selecionados mantêm sua relevância e consistência em diferentes períodos.

Este estudo também complementa iniciativas re-

centes, como o GILL, que busca mitigar a redundância futura por meio de uma nova plataforma de coleta [6]. Enquanto o GILL propõe novas estratégias para evitar redundância em plataformas futuras, a pesquisa aqui apresentada oferece uma solução prática para lidar com a redundância histórica existente, permitindo que analistas aproveitem ao máximo as plataformas atuais, como Route Views, RIPE RIS e PCH.

2 Contextualização

As plataformas RV, RIPE RIS e PCH desempenham um papel central na coleta de dados BGP, permitindo análises da estrutura e comportamento da Internet. Elas mantêm sessões de peering com roteadores que compartilham voluntariamente suas tabelas de roteamento.

Os dados coletados incluem três atributos principais [7]: (a) timestamp da atualização, (b) prefixo IP (IPv4/IPv6) anunciado, e (c) AS Path, sequência de ASes até o prefixo. Além disso, apenas RV e RIPE RIS registram (d) comunidades BGP, metadados opcionais que indicam preferências de roteamento, cuja ausência no PCH limita certas análises.

Esses dados são utilizados em diversos cenários, como a identificação de mudanças de rotas, detecção de incidentes, e mapeamento das relações entre ASes [8]. No entanto, como as plataformas frequentemente registram dados redundantes, os custos de armazenamento e processamento aumentam sem agregar valor proporcional.

Com o acúmulo de dados históricos, torna-se fundamental desenvolver abordagens eficientes de análise. Este trabalho propõe um algoritmo para selecionar coletores eficazes, garantindo ampla visibilidade da rede enquanto minimiza o volume de dados processado.

3 Metodologia de Coleta e Processamento

Os dados deste estudo foram coletados das plataformas RV [2], RIPE RIS [3] e PCH [4], utilizando o bgpscaner, uma ferramenta de parsing de MRT do extinto projeto Isolario [9]. A análise principal baseou-se em uma *snapshot* de 25 de abril de 2024, abrangendo 352 coletores e fornecendo uma visão atualizada da infraestrutura de roteamento da Internet. Além disso, conduziu-se uma análise temporal para verificar a consistência dos coletores selecionados nas datas de 1º de janeiro de 2023 e 1º de janeiro de 2022.

Foram extraídos e analisados os seguintes recursos para IPv4 e IPv6:

- Prefixos: Blocos de endereços IP anunciados.
- Enlaces: Conexões observadas entre sistemas autônomos.

- ASes: Sistemas Autônomos identificados nas rotas.
- Comunidades BGP Standard: Metadados de roteamento tradicionais, contendo valores de 32 bits.
- Comunidades BGP Large: Comunidades estendidas que utilizam 48 bits, oferecendo maior capacidade de descrição.
- Espaço de Endereçamento: Blocos de endereços IP normalizados em /24 para IPv4 e /48 para IPv6.

Cabe ressaltar que, embora RV e RIPE RIS registrem as comunidades BGP, o PCH não oferece esse tipo de dado. Essa distinção foi considerada na análise e interpretação dos metadados, especialmente para entender o impacto das diferentes comunidades na visibilidade da rede. Nenhum dado referente a comunidades BGP extended foi encontrado no conjunto de coletores.

3.1 Normalização e Filtragem do Espaço de Endereçamento

O espaço de endereçamento é fundamental para avaliar a cobertura dos blocos de forma consistente. Foi aplicada uma normalização: para IPv4, prefixos menores que /8 e maiores que /24 foram descartados, com os restantes normalizados para /24. Para IPv6, prefixos fora do intervalo /16 a /48 foram removidos e os demais normalizados para /48.

Essa higienização eliminou 115.545 prefixos IPv4 e 14.166 IPv6, correspondendo a 9,25% e 5,55% do total original, respectivamente, permitindo uma comparação uniforme entre diferentes coletores.

3.2 Resumo do Volume de Dados

Após a normalização e filtragem, a base de dados final consistiu nos seguintes recursos:

- IPv4: 1.248.651 prefixos, 614.617 enlaces, 77709 ASes e 11.888.342 blocos de endereçamento /24.
- IPv6: 255.246 prefixos, 403.910 enlaces, 41526 ASes e 14.616.109.062 blocos de endereçamento /48.
- Comunidades Standard IPv4: 88383 instâncias.
- Comunidades Standard IPv6: 45387 instâncias.
- Comunidades Large IPv4: 42620 instâncias.
- Comunidades Large IPv6: 43125 instâncias.

Esses recursos estão distribuídos entre um total de 352 coletores. A análise apresentada na próxima seção avalia como esses recursos se sobrepõem entre coletores e como a redundância observada afeta a eficiência das análises.

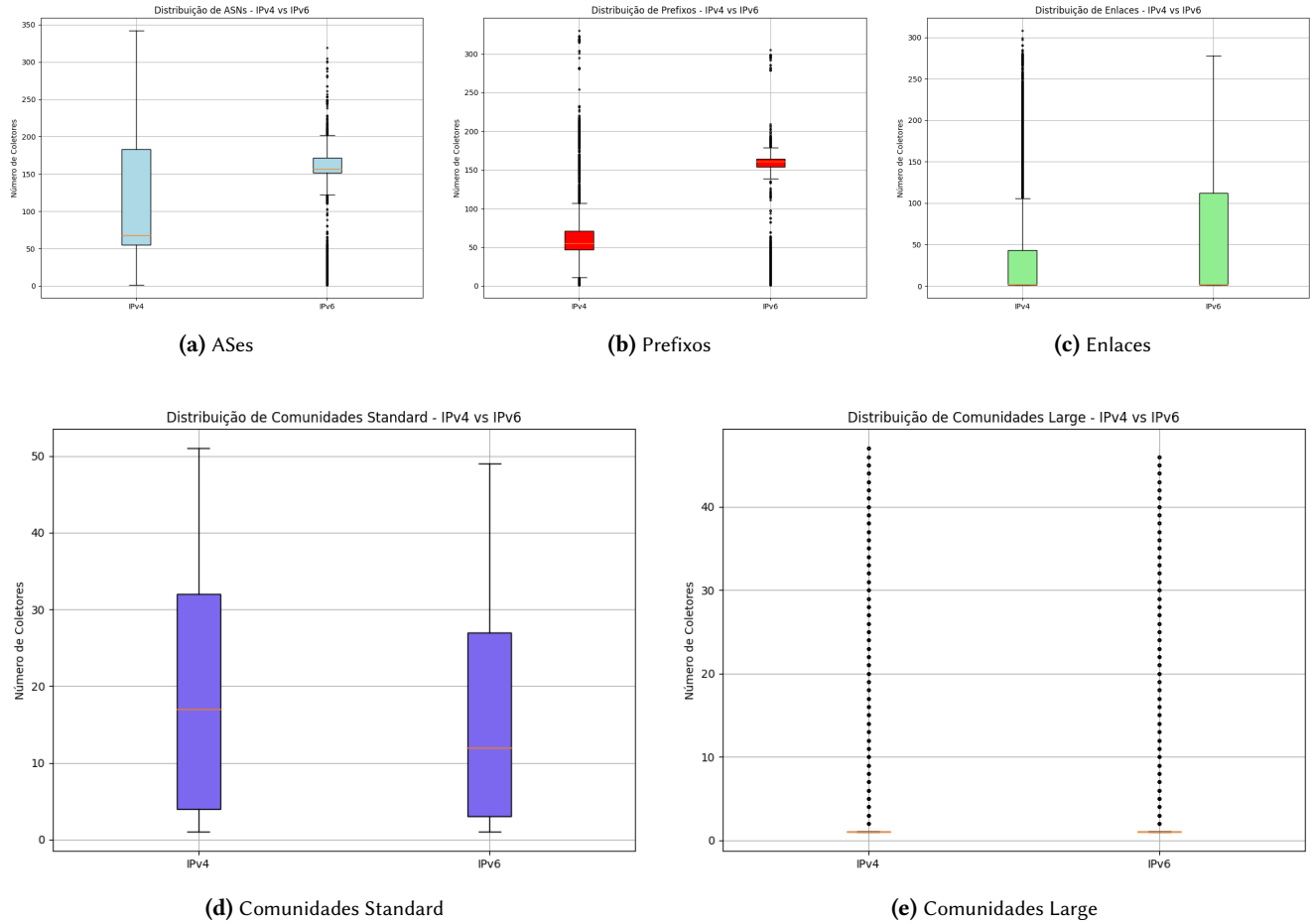


Figura 1. Boxplots mostrando a distribuição dos diferentes recursos entre coletores.

4 Análise da Redundância dos Dados BGP

Nesta seção, é analisada a redundância dos dados coletados pelas plataformas Route Views, RIS e PCH, mostrando como diferentes recursos se distribuem entre os coletores. A análise indica que muitos desses recursos aparecem repetidamente em diversos coletores, gerando redundância significativa e aumentando o custo de processamento sem agregar novas informações relevantes.

Os boxplots apresentados na Figura 1 demonstram a variabilidade dos diferentes recursos entre os coletores. Eles ajudam a visualizar como os recursos, como ASes, prefixos e comunidades BGP, estão distribuídos e identificam quais métricas são mais concentradas ou dispersas entre os coletores.

Os resultados dos boxplots destacam que alguns recursos são altamente concentrados em poucos coletores, enquanto outros são mais distribuídos. Por exemplo:

- **ASes e prefixos:** São observados em muitos coletores, sugerindo alta redundância. Isso indica que poucos coletores podem capturar a maioria desses recursos.

- **Comunidades BGP large:** São mais raras, com cerca de 75% das instâncias aparecendo em apenas um coletor. Isso enfatiza a importância de selecionar cuidadosamente os coletores para garantir a captura desses recursos menos redundantes.
- **Enlaces e comunidades standard:** Apresentam uma distribuição intermediária, exigindo uma quantidade moderada de coletores para assegurar cobertura adequada.

Além disso, a análise revela que os recursos em IPv6 tendem a ser observados em mais coletores comparados ao IPv4. Esse comportamento pode ser atribuído à maior representatividade dos coletores do PCH em IPv6, proporcionando uma visibilidade mais abrangente desse protocolo.

Essa distribuição reforça a necessidade de uma abordagem otimizada na seleção de coletores para minimizar o custo de processamento e maximizar a visibilidade, especialmente para recursos menos comuns e menos redundantes.

4.1 Redundância do Espaço de Endereçamento

A análise revela que um único coletor é suficiente para capturar quase toda a cobertura dos blocos IPv4 e IPv6. O coletor RRC25 do RIPE RIS, por exemplo, observa 99,64% do espaço IPv4 e 99,56% do IPv6.

Dado que um coletor já fornece praticamente toda a cobertura necessária, gráficos de CDF ou boxplots não foram aplicados. Embora um algoritmo de seleção pudesse ser utilizado, não foi aplicado neste caso, dado o alto nível de cobertura alcançado com apenas um coletor.

5 Problema de Cobertura de Conjuntos

O Set Cover Problem (SCP) consiste em, dado um conjunto universal U e uma coleção de subconjuntos S_1, S_2, \dots, S_n , onde $S_i \subseteq U$, selecionar o menor número de subconjuntos cuja união cubra todos os elementos de U . É necessário garantir que todos os recursos de U estejam presentes nos subconjuntos escolhidos, minimizando o total de subconjuntos selecionados.

No contexto deste trabalho, cada coletor é representado por um S_i com os recursos que ele observa, e U é o total de recursos de interesse.

5.1 Complexidade do Problema

O SCP é NP-completo, e encontrar uma solução ótima exigiria testar todas as combinações de subconjuntos, o que é inviável em instâncias suficientemente grandes. Assim, uma aproximação eficiente é essencial para lidar com centenas de coletores [10].

5.2 Solução Aproximada com Algoritmo Guloso

Dada a complexidade do SCP, uma abordagem eficiente é o algoritmo guloso. Esse método seleciona iterativamente o subconjunto que cobre mais elementos ainda não cobertos, repetindo o processo até que todos os elementos estejam incluídos. Embora não seja necessariamente ótimo, o algoritmo oferece uma aproximação: o custo final será no máximo $H(n)$ vezes o custo ótimo, onde $H(n)$ é o n -ésimo número harmônico [10, 11].

No contexto deste trabalho, a abordagem gulosa é útil para priorizar coletores que maximizam a cobertura enquanto minimizam redundâncias e esforço computacional.

6 Resultados e Avaliação

Nesta seção são apresentados os resultados da aplicação do algoritmo guloso para a seleção de coletores eficientes. A análise examina a relação entre a quantidade de coletores utilizados e a cobertura alcançada, além de avaliar o impacto no volume de dados processados para otimizar os recursos computacionais.

6.1 Cobertura dos Recursos e Tamanho dos Coletores

A seguir são apresentadas as tabelas com os resultados obtidos para IPv4 e IPv6. O conjunto completo de dados inclui 352 coletores, com um volume total de 113 GB. As tabelas mostram como diferentes métricas alcançam 95% de cobertura com um número menor de coletores e dados processados, destacando o equilíbrio entre maximizar a visibilidade e reduzir a redundância.

| IPv4 | Percentual | Coletores Utilizados | GB Processados |
|-------------------------|------------|----------------------|----------------|
| ASes | 98,31% | 1 | 7,75 |
| Enlaces | 95,13% | 32 | 82,8 |
| Prefixos | 95,00% | 7 | 25,4 |
| Espaço de Endereçamento | 99,64% | 1 | 7,75 |
| Comunidades Standard | 95,05% | 12 | 52,70 |
| Comunidades Large | 96,15% | 3 | 14,40 |

Tabela 1. Cobertura alcançada utilizando os coletores selecionados para IPv4.

Os resultados para IPv4, evidenciados na Tabela 1 mostram que ASes e espaço de endereçamento são rapidamente capturados com apenas um coletor, atingindo 98,31% e 99,64% de cobertura, respectivamente. Já outras métricas, como enlaces e comunidades BGP, exigem mais coletores para alcançar uma cobertura completa. Por exemplo, para enlaces, são necessários 32 coletores para atingir 95,13% de visibilidade. Isso indica que, enquanto alguns recursos estão concentrados em poucos coletores, outros são mais dispersos, demandando um esforço maior de coleta para garantir uma visão abrangente da rede.

| IPv6 | Percentual | Coletores Utilizados | GB Processados |
|-------------------------|------------|----------------------|----------------|
| ASes | 96,00% | 1 | 0,05 |
| Enlaces | 95,00% | 32 | 75,02 |
| Prefixos | 95,46% | 6 | 21,00 |
| Espaço de Endereçamento | 99,56% | 1 | 7,75 |
| Comunidades Standard | 95,43% | 15 | 57,20 |
| Comunidades Large | 95,07% | 3 | 18,10 |

Tabela 2. Cobertura alcançada utilizando os coletores selecionados para IPv6.

Os resultados para IPv6 seguem um padrão semelhante ao IPv4, com ASes e espaço de endereçamento rapidamente capturados por um único coletor, indicando alta concentração. Por outro lado, métricas como enlaces e comunidades BGP exigem mais coletores para alcançar 95% de visibilidade, refletindo sua menor redundância. A necessidade de 32 coletores para enlaces e 15 para comunidades standard sugere que esses recursos são mais raros, já que uma parte significativa deles é observada em poucos coletores, exigindo uma coleta mais abrangente para garantir uma cobertura ampla.

6.2 Impacto do Algoritmo Guloso

A aplicação do algoritmo guloso permitiu selecionar um subconjunto eficiente de coletores, garantindo alta cobertura para várias métricas com um volume controlado de dados processados. A análise revelou que menos de 10 de um total de 352 coletores são suficientes para cobrir 95% de ASes, prefixos e do espaço de endereçamento, enquanto enlaces e comunidades BGP standard exigem mais coletores.

Esses resultados confirmam a importância de uma seleção otimizada, assegurando a captura dos recursos essenciais e evitando o processamento excessivo de dados redundantes.

6.3 Consistência ao Longo do Tempo

A consistência dos coletores selecionados pelo algoritmo guloso foi avaliada ao longo do tempo. As Tabelas 3 e 4 apresentam a cobertura alcançada em 2023 e 2022, utilizando os mesmos coletores definidos para garantir pelo menos 95% de visibilidade em 2024.

| IPv4 | 2024 | 2023 | 2022 |
|-------------------------|--------|--------|--------|
| ASes | 98,31% | 98,57% | 96,75% |
| Enlaces | 95,13% | 92,78% | 89,95% |
| Prefixos | 95,00% | 93,72% | 90,21% |
| Espaço de Endereçamento | 99,64% | 99,58% | 99,70% |
| Comunidades Standard | 95,05% | 94,32% | 93,24% |
| Comunidades Large | 95,15% | 81,15% | 87,51% |

Tabela 3. Cobertura alcançada ao longo do tempo utilizando os coletores selecionados para IPv4.

| IPv6 | 2024 | 2023 | 2022 |
|-------------------------|--------|--------|--------|
| ASes | 96,00% | 81,11% | 81,93% |
| Enlaces | 95,00% | 82,09% | 81,60% |
| Prefixos | 95,46% | 92,59% | 92,41% |
| Espaço de Endereçamento | 99,56% | 99,86% | 99,73% |
| Comunidades Standard | 95,43% | 94,86% | 94,64% |
| Comunidades Large | 95,07% | 76,41% | 69,30% |

Tabela 4. Cobertura alcançada ao longo do tempo utilizando os coletores selecionados para IPv6.

Os resultados das Tabelas 3 e 4 mostram que a cobertura dos coletores selecionados para 2024 mantém-se alta nos anos anteriores, embora com pequenas variações.

Para IPv4, ASes, prefixos e espaço de endereçamento apresentaram estabilidade, enquanto as comunidades large mostraram maior volatilidade, caindo para 81,15% em 2023.

Em IPv6, a maior instabilidade observada pode ser associada ao crescimento da adoção do protocolo, o que torna a visibilidade menos previsível em anos distintos. A variabilidade nas comunidades large sugere a necessidade de reavaliação constante para garantir alta cobertura ao longo do tempo.

7 Trabalhos Futuros

Embora o algoritmo guloso tenha sido eficiente, alternativas mais flexíveis podem ser exploradas. Uma proposta é adotar um sistema de pontuação, atribuindo pesos às métricas, como ASes e prefixos, para priorizar coletores com maior impacto. Esse sistema permitiria uma seleção dinâmica, adaptando-se a mudanças na topologia em tempo real. A ideia é melhorar a eficiência do monitoramento, evitando redundância e assegurando que as informações mais relevantes sejam capturadas, mantendo a visibilidade da rede sem sobrecarregar o processamento.

8 Conclusão

Este trabalho analisou redundâncias na coleta de dados BGP e aplicou um algoritmo guloso para selecionar subconjuntos eficientes de coletores. A abordagem identificou que diferentes coletores se destacam em métricas específicas, como ASes e prefixos, otimizando a cobertura dos recursos e reduzindo o processamento redundante.

A análise temporal indicou que os coletores mantêm um desempenho consistente ao longo do tempo, reforçando a necessidade de revisões periódicas. Assim, a solução apresentada permite otimizar a análise dos dados nas plataformas Route Views, RIS e PCH, sem comprometer a visibilidade da rede.

Referências

1 Gao, L. On inferring autonomous system relationships in the Internet. *IEEE/ACM Transactions on Networking*, v. 9, n. 6, p. 733–745, 2001. DOI: [10.1109/90.974527](https://doi.org/10.1109/90.974527).

2 University of Oregon. *Route Views Project*. 2024. Disponível em: <http://www.routeviews.org/>.

3 RIPE NCC. *RIPE Routing Information Service (RIS)*. 2024. Disponível em: <https://ris.ripe.net/docs/mrt/>.

4 Packet Clearing House. *Packet Clearing House (PCH)*. 2024. Disponível em: https://www.pch.net/resources/Routing_Data/.

5 Milolidakis, A. et al. On the Effectiveness of BGP Hijackers That Evade Public Route Collectors. *IEEE Access*, v. 11, p. 31092–31124, 2023. DOI: [10.1109/ACCESS.2023.3261128](https://doi.org/10.1109/ACCESS.2023.3261128).

- 6 Alfroy, T. *et al.* The Next Generation of BGP Data Collection Platforms. *In: ACM SIGCOMM '24: Proceedings of the ACM SIGCOMM 2024 Conference*. 2024. P. 794–812. DOI: [10.1145/3651890.3672251](https://doi.org/10.1145/3651890.3672251).
- 7 Orsini, C. *et al.* BGPStream: A Software Framework for Live and Historical BGP Data Analysis. *IMC '16: Proceedings of the 2016 Internet Measurement Conference*, v. 11, p. 429–444, 2016. DOI: [10.1145/2987443.2987482](https://doi.org/10.1145/2987443.2987482).
- 8 Luckie, M. *et al.* AS relationships, customer cones, and validation. *IMC '13: Proceedings of the 2013 conference on Internet measurement conference*, p. 243–256, 2013. DOI: [10.1145/2504730.2504735](https://doi.org/10.1145/2504730.2504735).
- 9 Cogotti, L. *bgpscanner: MRT parser in C for high-speed BGP data analysis*. 2019. Disponível em: <https://gitlab.com/Isolario/bgpscanner>.
- 10 Vazirani, V. V. *Approximation Algorithms*. Irvine, CA: Springer, 2001. ISBN 978-3-540-65367-7.
- 11 Chandu, D. P. Improved Greedy Algorithm for Set Covering Problem. *SSRG International Journal of Computer Science and Engineering*, abs/1506.04220, 2015. Disponível em: <https://api.semanticscholar.org/CorpusID:28267446>.