

Compartilhamento de Dados de Tráfego de Rede Utilizando Privacidade Diferencial

Felipe C. Monteiro, Felipe T. Brito, Iago C. Chaves, Javam C. Machado

Departamento de Computação – Universidade Federal do Ceará (UFC)
60440-900 – Campus do Pici - Bloco 910 – Fortaleza – CE – Brasil

{felipe.monteiro, felipe.timbo, iago.chaves, javam.machado}@lsbd.ufc.br

Abstract. *Network traffic data is useful for a variety of applications. Entities that collect this type of data, such as Internet Service Providers (ISPs), generally share their network traffic information with external entities. However, this sharing can lead to privacy violations for the individuals contained in this data. This work proposes a new approach to sharing network traffic data using differential privacy, a method that aims to add noise to the original data. Experimental results show that the proposed approach introduces less noise into the data when compared to other techniques that also adopt differential privacy.*

Resumo. *Dados de tráfego de redes são úteis para uma variedade de aplicações. Geralmente as entidades que coletam esse tipo de dado, por exemplo provedores de internet (ISPs), compartilham suas informações de tráfego de rede com entidades externas. Contudo, esse compartilhamento pode levar a violações de privacidade dos indivíduos contidos nesses dados. Este trabalho propõe uma nova abordagem para compartilhamento de dados de tráfego de rede utilizando privacidade diferencial, um método que tem como objetivo adicionar ruído sobre os dados originais. Resultados experimentais mostram que a abordagem proposta introduz menos ruído nos dados quando comparada a outras técnicas que também adotam privacidade diferencial.*

1. Introdução

Com a expansão de serviços e plataformas *web* na era digital, surge uma crescente demanda por coleta e monitoramento de dados de tráfego de redes [Joshi e Hadi 2015]. Esses dados são geralmente coletados, gerenciados e monitorados por provedores de serviço de *Internet (Internet Service Provider - ISP)*, e se tornam essenciais para uma variedade de aplicações, tais como detecção de pontos de acessos, conhecimento do perfil de utilização da rede, identificação de comportamento anômalo de tráfego e alocação adequada de recursos [Shafiq et al. 2016].

No entanto, para que essas análises possam ser realizadas, ISPs geralmente precisam compartilhar ou comercializar seus dados de tráfego de rede com entidades externas e tal fato pode levar a violações de privacidade dos indivíduos contidos nesses dados [Mirim 2018]. Consequentemente, sanções podem ser aplicadas ao ISP devido a falta de conformidade às leis de privacidade vigentes no país onde o dado é coletado, como a Lei Geral de Proteção de Dados Pessoais no Brasil [LGPD 2019], Regulamento Geral de Proteção de Dados na Europa [GDPR 2018] e a Comissão Federal de Comunicação nos Estados Unidos [Comissão Federal de Comunicação 2018].

Uma das abordagens mais promissoras para proteger a privacidade dos indivíduos é anonimizar os dados antes de compartilhá-los com terceiros [Brito e Machado 2017]. Dessa forma, os dados deixam de ser pessoais e as leis de privacidade passam a não mais se aplicar aos mesmos [LGPD 2019, GDPR 2018]. Nesse contexto, várias técnicas de anonimização de dados foram propostas nas últimas décadas, por exemplo *k-anonymity*, *l-diversity*, *δ -presence* e *t-closeness* [Brito e Machado 2017]. No entanto, todas essas abordagens assumem que um adversário (usuário malicioso) possui conhecimento limitado sobre os dados, o que não é verdadeiro em situações do mundo real.

Para exemplificar, considere o seguinte cenário de contagem de serviços de *streaming* coletados por um ISP em uma determinada região do país: suponha um total de 100 usuários de serviços de *streaming*, dos quais 60 utilizam *Netflix*, 30 utilizam *AmazonVideo*, 9 utilizam *AppleTv* e apenas um indivíduo utiliza o serviço *Looke*. Considere que o ISP, detentor desses dados, queira comercializá-los com a empresa *Looke*, a qual deseja descobrir quantos usuários estão utilizando seus serviços nessa região. Esse compartilhamento de dados pode ser útil para eventuais direcionamentos de marketing da empresa *Looke*, com o intuito de aumentar sua carteira de clientes na região. Caso o ISP comercialize ou compartilhe os dados de tráfego da sua rede com a essa entidade externa, informações pessoais de seus clientes podem estar presentes no conjunto de dados compartilhado. No exemplo acima, caso um adversário, munido de informações externas, conheça o usuário que utiliza o serviço de streaming *Looke* na região, a privacidade deste usuário é violada, e assim, descumpre-se também as leis de privacidade vigentes.

A privacidade diferencial [Dwork 2006] é uma técnica que visa proteger a privacidade dos indivíduos enquanto permite que os dados sejam compartilhados e utilizados para fins de análise. Nos últimos anos, essa técnica tem se tornado o padrão para compartilhamento de dados de maneira privada [Sangeetha e Sudha Sadasivam 2019]. Ela consiste em adicionar um nível controlado de ruído aos dados originais de forma a impedir a identificação de indivíduos específicos no conjunto de dados compartilhado. Por meio dessa técnica, um adversário não deve ser capaz de aprender nada sobre um indivíduo específico que ele já não poderia ter aprendido anteriormente sem acesso aos dados [Dwork 2006]. Logo, a privacidade diferencial assume que o conhecimento adversário sobre os dados é ilimitado. Essa suposição é feita para garantir que mesmo um adversário que possua acesso a muitas fontes externas de informação e possa correlacionar diferentes bases de dados, não consiga usar essas informações para identificar indivíduos específicos sobre os dados compartilhados. É importante mencionar que, após a inserção de ruído, a utilidade dos dados para eventuais análises pode ser comprometida. Dessa forma, torna-se desafiador propor soluções que eficientemente adicionem ruído a dados de tráfego de rede de tal sorte que a privacidade dos indivíduos é protegida, enquanto a utilidade dos dados é mantida.

O restante desse trabalho está dividido da seguinte forma: Na Seção 2, discutem-se os trabalhos relacionados. A seção 3 aborda os conceitos sobre privacidade diferencial e suas propriedades. Já a Seção 4 detalha a abordagem proposta. Na Seção 5, avaliam-se os resultados obtidos neste trabalho. Por fim, a Seção 6 apresenta as conclusões obtidas e direciona os trabalhos futuros.

2. Trabalhos Relacionados

Alguns trabalhos existentes na literatura propõem diferentes técnicas de preservação de privacidade para dados de tráfego de rede. Os autores em Lu et al. 2018 implementam um método para preservar a privacidade de informações de tráfego de rede enquanto permitem a classificação dessas informações para análise de segurança. O método proposto utiliza a técnica de perturbação para adicionar ruído aleatório às informações de tráfego de rede antes da classificação. Os resultados dos experimentos mostram que o método é capaz de preservar a privacidade das informações de tráfego de rede enquanto fornece resultados precisos de classificação. Contudo, o trabalho não adota garantias formais com base na privacidade diferencial, e assim oferecem uma noção mais fraca de privacidade.

Já no trabalho de Koufogiannis e Pappas 2017, uma abordagem para proteger a privacidade de dados sensíveis em redes distribuídas é apresentada. A proposta utiliza a técnica de difusão para proteger a privacidade dos dados, adicionando ruído aleatório aos dados originais antes de transmiti-los pela rede. Por outro lado, o trabalho de Castro Vidal et al. 2020 propõe uma estratégia diferencialmente privada para estimar frequências de valores no contexto de dados de casas inteligentes, focando em fluxo contínuo de dados. Já o trabalho de Zhang et al. 2022 apresenta um método para proteger a privacidade do usuário durante a transmissão de dados em tempo real, distorcendo os dados para que possam ser úteis para análises de terceiros. A Tabela 1 compara os trabalhos relacionados.

Tabela 1. Tabela comparativa de trabalhos relacionados.

Trabalho	Estratégia de Privacidade	Aplicação
[Koufogiannis e Pappas 2017]	Privacidade Diferencial	Serviços
[Lu et al. 2018]	Perturbação	Protocolos
[de Castro Vidal et al. 2020]	Privacidade Diferencial	Serviços
[Zhang et al. 2022]	Privacidade Diferencial	Serviços
Este Trabalho	Privacidade Diferencial	Portas, Protocolos e Serviços

Apesar dos estudos realizados nos trabalhos mencionados anteriormente, nenhum deles foca no compartilhamento de tráfego de rede de maneira privada. Dessa forma, este trabalho propõe uma nova técnica para permitir aos provedores de internet e entidades que coletam dados de tráfego de rede compartilharem, ou comercializarem, seus dados por meio da privacidade diferencial. Em particular, desenvolvemos uma técnica que permite compartilhar três informações importantes sobre tráfego de rede de maneira privada: contagem de portas, contagem de protocolos e contagem de serviços.

3. Privacidade Diferencial

Privacidade diferencial (PD) é um modelo matemático que possibilita análises estatísticas sobre um conjunto de dados, sem comprometer a privacidade dos indivíduos contidos neles [Dwork 2006]. Ela dificulta a reidentificação de indivíduos a partir de ataques de ligação, em que o atacante possui um conhecimento prévio por meio de fontes externas. A privacidade diferencial protege a privacidade dos indivíduos baseada no conceito de indistinguibilidade de conjunto de dados. Considere uma função de consulta de contagem $q : \mathcal{D} \rightarrow \mathbb{Z}$ onde \mathcal{D} indica o conjunto de possibilidades de todos os conjuntos de dados. A função de consulta de contagem q é aplicada em um conjunto de dados e retorna um número inteiro. A privacidade diferencial se baseia no conceito de conjunto de dados

vizinhos. Dois conjuntos de dados $\mathcal{D}_1 \in \mathcal{D}$ e $\mathcal{D}_2 \in \mathcal{D}$ são ditos *conjuntos de dados vizinhos* se eles diferirem em um registro, isto é, $|\mathcal{D}_1 - \mathcal{D}_2| = 1$ [Dwork 2006].

Privacidade diferencial é satisfeita por um algoritmo aleatório, comumente chamado de *mecanismo* \mathcal{M} , capaz de adicionar um ruído apropriado para gerar uma resposta à consulta de contagem q realizada pelo usuário. Esse mecanismo é controlado por um parâmetro denominado *orçamento de privacidade* ε (ou do inglês - *privacy budget*) [Dwork 2006]. Um menor ε corresponde diretamente a uma maior garantia da preservação de privacidade, ou seja, uma maior quantidade de ruído adicionado ao dado original. Em geral, para se definir um orçamento de privacidade adequado para uma aplicação, especialistas devem realizar um amplo estudo para garantir que a privacidade dos indivíduos seja suficientemente protegida, mantendo bons níveis de precisão nas informações compartilhadas, isto é, uma boa utilidade dos dados [Bureau 2021]. Neste trabalho, utilizamos uma variação do orçamento de privacidade de 0.1 a 1.0 (Seção 4), valores comumente reportados na literatura, independente da aplicação.

Formalmente, a privacidade diferencial é definida por:

Definição 1 (*Privacidade Diferencial [Dwork 2006]*) Um mecanismo aleatório \mathcal{M} satisfaz Privacidade Diferencial, se para todos os conjuntos de dados vizinhos \mathcal{D}_1 e \mathcal{D}_2 que diferem em pelo menos um elemento e para qualquer possibilidade de saída O de \mathcal{M} ,

$$Pr[\mathcal{M}(\mathcal{D}_1) = O] \leq \exp(\varepsilon) Pr[\mathcal{M}(\mathcal{D}_2) = O],$$

onde $Pr[\cdot]$ indica a probabilidade de um dado evento.

Os mecanismos mais utilizados para consultas numéricas são o de Laplace [Dwork 2006] e o Geométrico [Ghosh et al. 2009]. Neste trabalho, utilizamos o mecanismo Geométrico, uma variante discreta do mecanismo de Laplace, que é usualmente adotado para consultas com respostas inteiras, no caso, consultas de contagem. A distribuição geométrica simétrica, com média 0 e parâmetro $\alpha \in [0, 1]$, é a distribuição de probabilidade tal que, para todos os inteiros x , uma variável aleatória X tem função de massa de probabilidade:

$$Pr[X = x] = \frac{1 - \alpha}{1 + \alpha} \alpha^{|x|} \quad (1)$$

Teorema 1 (*Mecanismo Geométrico [Ghosh et al. 2009]*) O mecanismo Geométrico \mathcal{M} que adiciona ruído independente a partir da distribuição geométrica simétrica, com $\alpha = \exp(-\varepsilon/\Delta q)$, satisfaz Privacidade Diferencial.

Mecanismos atuam adicionando uma certa quantidade de ruído aleatório às respostas das consultas. O ruído adicionado, além de ser dependente do *orçamento de privacidade*, é também dependente da sensibilidade global da consulta realizada.

Definição 2 (*Sensibilidade Global*) A sensibilidade global de uma consulta q é definida por:

$$\Delta q = \max_{\mathcal{D}_1, \mathcal{D}_2} \| q(\mathcal{D}_1) - q(\mathcal{D}_2) \|_1,$$

para todo $\mathcal{D}_1 \in \mathcal{D}$ e $\mathcal{D}_2 \in \mathcal{D}$ [Dwork 2006].

A sensibilidade global mede o maior impacto em relação à presença ou ausência de um registro em todos os possíveis pares de conjuntos de dados vizinhos. Quando q é uma consulta (função) de contagem, a sensibilidade global $\Delta q = 1$, uma vez que a adição

ou remoção de qualquer registro em um conjunto de dados pertencente a \mathcal{D} impacta em no máximo 1 sobre qualquer consulta de contagem. A privacidade diferencial também apresenta propriedades úteis, como o pós-processamento e composição sequencial.

Definição 3 (*Pós-Processamento*) *Considere \mathcal{M} qualquer mecanismo aleatório tal qual $\mathcal{M}(q)$ é diferencialmente privado. Para qualquer função f , $f(\mathcal{M}(q))$ também satisfaz Privacidade Diferencial [Dwork 2006].*

A propriedade de Pós-Processamento define que qualquer função aplicada sobre uma saída de um mecanismo diferencialmente privado também satisfaz a privacidade diferencial [Dwork 2006].

Definição 4 (*Composição Sequencial*) *Para cada mecanismo \mathcal{M}_i que provê ε_i -privacidade diferencial, uma sequência de mecanismos diferencialmente privados \mathcal{M}_i , provê $\sum_i \varepsilon_i$ -privacidade diferencial [McSherry 2009].*

A composição sequencial é aplicada quando executamos uma série de mecanismos sequencialmente no mesmo conjunto de dados. Isso implica que o orçamento de privacidade ε usado em cada consulta sobre o mesmo conjunto de dados precisa ser dividido [McSherry 2009].

Uma outra maneira de garantir a privacidade diferencial, ao invés da adição de ruído sobre as respostas originais, é através da geração de dados sintéticos por meio do PrivBayes [Zhang et al. 2017]. Essa técnica utiliza redes Bayesianas para modelar as distribuições de probabilidade das variáveis presentes no conjunto de dados original. Esse processo é realizado utilizando a privacidade diferencial. Com base nessa modelagem, o método gera dados sintéticos que preservam as estatísticas e as relações entre as variáveis contidas nos dados originais, garantindo que as propriedades estatísticas do conjunto de dados original sejam mantidas nos dados gerados sinteticamente.

4. Abordagem Diferencialmente Privada para Dados de Tráfego de Rede

Com o objetivo de proteger as informações sensíveis dos usuários de acessos indesejados, este trabalho propõe uma nova técnica para compartilhar contagens relacionadas a dados de tráfego de rede utilizando privacidade diferencial. Contagens são um objeto natural de estudo para análise de dados no contexto de preservação da privacidade [Machanavajjhala et al. 2017]. Elas permitem que usuários computem o número de indivíduos em um conjunto de dados que satisfazem predicados. São exemplos de consultas de contagem: histogramas, intervalos (*range queries*), funções de distribuição cumulativas, entre outras. Consequentemente, ao publicar consultas de contagens de maneira privada, entidades que coletam dados podem realizar uma série de análises importantes sobre os dados sem comprometer a privacidade dos indivíduos pertencentes a eles.

Conforme mencionado anteriormente, este trabalho visa publicar (compartilhar) o resultado de três contagens distintas de maneira privada: contagem de portas de destino, de protocolos utilizados e de serviços dos fluxos de rede. Em particular, a contagem de portas de destino fornece informações cruciais sobre a capacidade e o desempenho de uma rede. Por meio dessa contagem é possível determinar quantos dispositivos podem estar conectados à rede, bem como quantas conexões de rede estão em uso [Tanenbaum e Wetherall 2021]. Por outro lado, ao contar a quantidade de protocolos TCP e UDP em

uma rede, os administradores de rede podem monitorar o tráfego e identificar possíveis problemas de desempenho. Por exemplo, se houver uma grande quantidade de tráfego TCP, isso pode indicar uma grande quantidade de transferência de dados entre dispositivos na rede, o que pode levar a problemas de congestionamento. Além disso, a contagem de protocolos TCP e UDP também pode auxiliar na identificação de possíveis ameaças de segurança, como ataques de negação de serviço (DoS) ou tentativas de invasão, muitas vezes manifestados como tráfego incomum ou excessivo em um determinado protocolo [Tanenbaum e Wetherall 2021]. Por fim, a contagem de serviços provenientes dos fluxos de rede é particularmente importante para auxiliar administradores a rastrear o uso da rede por diferentes departamentos ou usuários, o que pode ser útil para fins de faturamento ou para identificar possíveis violações de políticas de uso da rede. Dessa forma, é possível responder a perguntas, de maneira privada, do tipo: “qual serviço de *streaming* é mais utilizado?”.

A abordagem mais usual para contornar esse tipo de problema é aplicar diretamente a privacidade diferencial através do mecanismo Geométrico sobre cada uma das três consultas de contagem. Em outras palavras, para cada contagem de portas de destino, protocolos e serviços, um ruído é adicionado com base na distribuição geométrica simétrica (Teorema 1). Contudo, como três contagens distintas necessitam ser obtidas sobre o mesmo conjunto de dados, o orçamento de privacidade ϵ deve ser dividido por três (Definição 4). Isso faz com que o mecanismo Geométrico produza saídas com ruídos maiores e, conseqüentemente, contagens com maiores erros. Esses resultados são apresentados na Seção 4.

Para lidar com o problema de compartilhamento de dados de tráfego de rede de maneira privada, nós propomos uma abordagem baseada em três etapas: (1) pré-processamento do conjunto de dados original com base em informações públicas; (2) aplicação do mecanismo Geométrico sobre os dados pré-processados; e (3) agrupamento das contagens a partir dos dados ruidosos. A Figura 1 exemplifica a abordagem proposta a partir de dados originais de tráfego de rede até o compartilhamento ruidoso das contagens. Vale ressaltar que, neste trabalho um fluxo de rede, isto é, o dado original é composto por IP de destino (IP Dest.), Porta de destino (Porta Dest.), Protocolo e Serviço utilizados.

4.1. Pré-processamento dos Dados

O objetivo do pré-processamento é agrupar os dados de tal sorte que a adição de ruído, via privacidade diferencial, seja executada apenas uma vez, ao invés de três vezes seguindo a abordagem mais usual. Essa estratégia evita que o orçamento de privacidade seja dividido e, conseqüentemente, produzindo contagens mais próximas das originais.

Inicialmente, a partir do conjunto de dados original, o IP de destino é removido, por ser um identificador explícito, e os dados originais são agrupados por triplas (porta, protocolo e serviço). Essa etapa de agrupamento é realizada com base nos registros de dados de tráfego de rede disponibilizados publicamente pela Internet Assigned Numbers Authority [IANA 2023], organização responsável por coordenar a alocação global de recursos relacionados à Internet. Em outras palavras, a IANA mantém um registro de portas de protocolo e serviços atribuídos, que inclui informações sobre os protocolos associados a essas portas e os serviços que eles fornecem. Assim, é possível agrupar os dados com base nessas informações sem ferir a privacidade dos indivíduos, visto que essas informações são de domínio público.

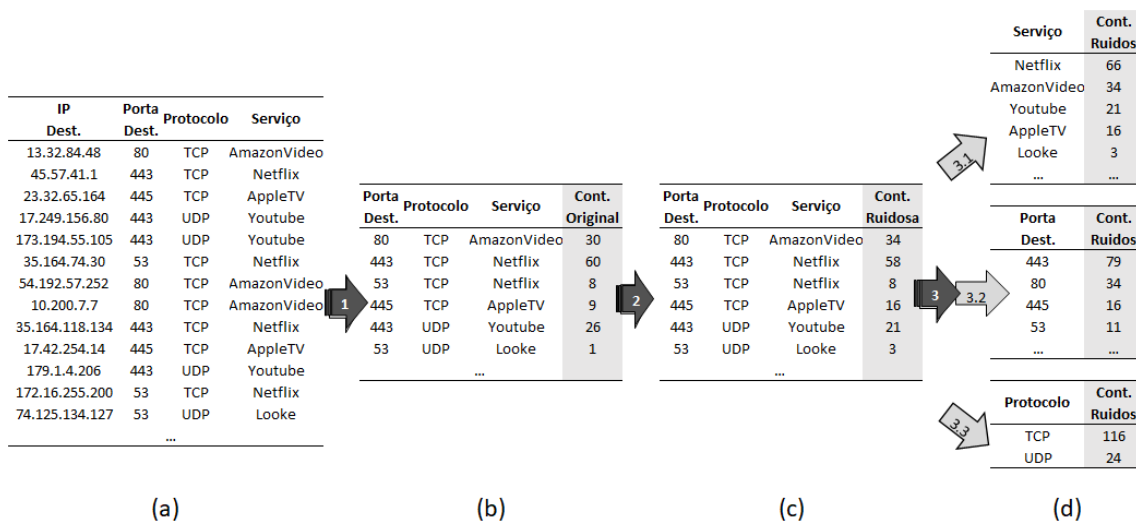


Figura 1. Exemplo de conjunto de dados de tráfego de rede compartilhados por meio da nossa abordagem. (a) Conjunto de dados original. (b) Conjunto de dados pré-processados e suas respectivas contagens originais. (c) Conjunto de dados agregado e sua contagem ruidosa. (d) Contagem ruidosa de serviços, portas e protocolos.

Para cada tripla (porta, protocolo e serviço) agrupada com base nos padrões da IANA, suas respectivas contagens são calculadas no conjunto de dados original. Essas contagens são informações privadas e, uma vez compartilhadas sem garantias formais de privacidade, podem levar à re-identificação de indivíduos. A Figura 1b apresenta um exemplo do conjunto de dados original pré-processado com suas respectivas contagens.

4.2. Aplicação do Mecanismo Geométrico

Nesta etapa, o mecanismo Geométrico é aplicado sobre as contagens de cada tripla pré-processada na fase anterior. Esse processo envolve a introdução de ruído aleatório na contagem original calibrada pelo orçamento de privacidade ϵ e pela sensibilidade global Δq . Dessa forma, a privacidade dos dados é formalmente garantida. A Figura 1c mostra os dados agrupados anteriormente, agora com suas respectivas contagens ruidosas. No exemplo em questão, para a tripla (80, TCP, AmazonVideo) um ruído de valor 4 foi adicionado à contagem original. Já para a tripla (443, TCP, Netflix) um ruído de valor -2 foi acrescentado ao dado original. Já para a tripla (53, TCP, Netflix) um ruído de valor 0 foi introduzido, isto é, o dado permaneceu o mesmo. Pelo fato do ruído adicionado aos dados ser aleatório, um usuário malicioso que tem acesso aos dados compartilhados não possui conhecimento de quais contagens foram alteradas e quais permaneceram as mesmas, e nem a quantidade de ruído introduzido em cada contagem.

4.3. Pós-processamento

Finalmente, após adição de ruído aos dados, uma etapa de pós-processamento é conduzida a fim de obter as contagens agregadas por portas de destino, protocolos e serviços, isto é, o objetivo final da nossa abordagem. A Figura 1d apresenta os dados de contagem de serviços, portas destino e protocolos agregados pela soma de suas contagens ruidosas. Vale ressaltar que, ao compartilhar uma determinada contagem, seja de qualquer tipo, a nossa abordagem pode retornar um valor negativo devido à adição de ruído aleatório

pelo mecanismo Geométrico ser tanto positiva quanto negativa (Teorema 1). Nosso pós-processamento arredonda os valores agregados de contagem negativa para o menor valor inteiro positivo possível, isto é, 1 (um). É importante mencionar também que a etapa de pós-processamento não viola as garantias formais da privacidade diferencial, visto que qualquer função aplicada sobre um conjunto de dados já diferencialmente privado também satisfaz a privacidade diferencial (Definição 3).

5. Resultados Experimentais

Experimentos foram conduzidos em ambiente com sistema operacional Linux Ubuntu 22.04 LTS, Processador Intel i7 2.7 GHz e 16GB de memória RAM. Comparamos nossa abordagem com três concorrentes existentes na literatura: Mecanismo Geométrico [Ghosh et al. 2009], Mecanismo Log-Laplace [Le Ny e Pappas 2013] e Privbayes [Zhang et al. 2017]. A abordagem proposta neste trabalho e os concorrentes foram implementados na linguagem de programação Python. Adicionalmente, os experimentos foram executados 10 vezes e a média destes resultados foram reportadas.

Utilizamos dois conjuntos de dados reais: o primeiro é um conjunto de dados privado coletado em um laboratório de pesquisa localizado no Brasil, denominado “*conjunto de dados local*”. Esse conjunto de dados possui 33.845 fluxos de rede. Os valores mínimos e máximos de contagem para este conjunto de dados e seus principais atributos estatísticos são descritos na Tabela 2. O segundo conjunto de dados é usualmente adotado na literatura para fins científicos, denominado “*IP Network Traffic Flows*”¹. Esse conjunto possui 1.046.015 fluxos de rede. Dados estatísticos sobre esses dados são igualmente descritos na Tabela 2.

Tabela 2. Tabela de contagem – Conjuntos de Dados.

	<i>Local</i>			<i>IP Network Traffic Flows</i>		
	Portas	Protocolos	Serviços	Portas	Protocolos	Serviços
<i>Contagem</i>	29	2	25	8,465	2	120
<i>Min</i>	1	5.802	1	1	443.273	1
<i>25%</i>	2	11.362	2	1	483.140	21.5
<i>50%</i>	20	16.922	20	1	523.007	241.5
<i>75%</i>	140	22.482	140	2	562.874	2406
<i>Max</i>	22.844	28.043	22.844	423.039	602.742	230.847

Neste trabalho, examinamos tanto o erro relativo médio quanto os top- k serviços e portas mais frequentes para todos as técnicas analisadas.

5.1. Avaliação do Ruído Introduzido

A avaliação do ruído introduzido é importante para mensurar a utilidade de dados após inserção de ruído. O erro relativo médio (MRE) é uma métrica comum utilizada para avaliar a utilidade de dados ruidosos. Dessa forma, quanto menor o erro relativo médio – MRE, mais preciso é o mecanismo diferencialmente privado na preservação da utilidade dos dados.

Essa métrica é definida por:

$$MRE(y, \hat{y}) = \frac{1}{N} \sum_{i=0}^{N-1} \frac{|y_i - \hat{y}_i|}{|y_i|} \quad (2)$$

¹<https://www.kaggle.com/datasets/jsrojas/ip-network-traffic-flows-labeled-with-87-apps>

Coletamos o erro das contagens de portas, protocolos e serviços com 3 valores diferentes de orçamentos de privacidade ϵ : 0.1, 0.5 e 1.0. Os resultados de MRE com seus respectivos intervalos de confiança (95%) para 10 execuções são reportados na Figura 2 para o “conjunto de dados local” e Figura 3 para o “IP Network Traffic Flows”.

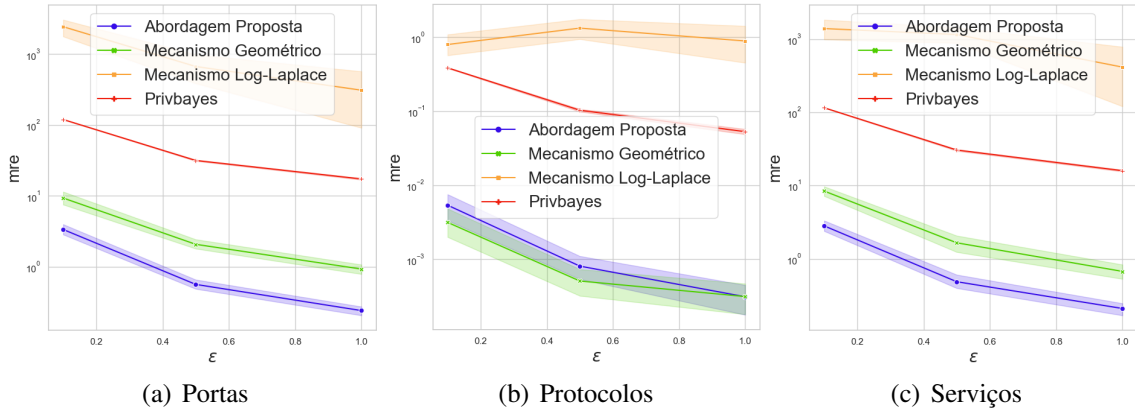


Figura 2. Erro Relativo Médio do conjunto de dados Local.

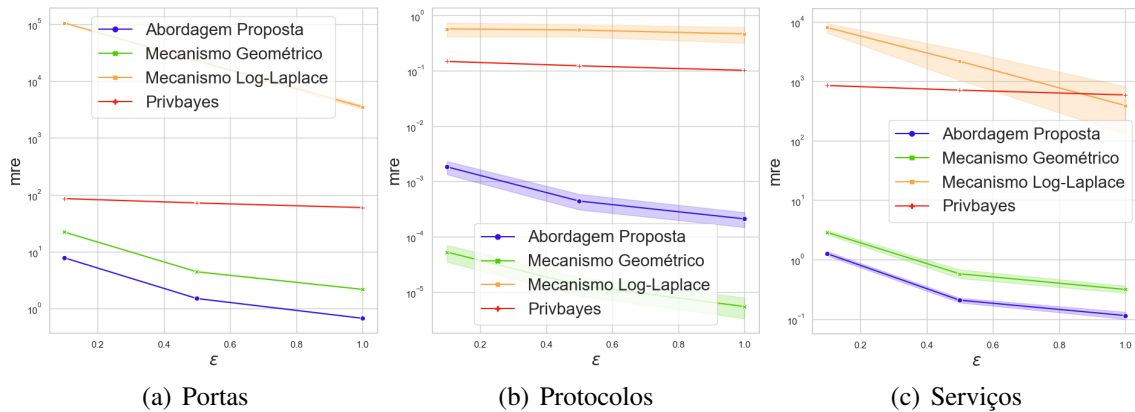


Figura 3. Erro Relativo Médio do conjunto de dados IP Network Traffic Flows.

Conforme esperado, à medida que o orçamento de privacidade ϵ aumenta, o erro relativo diminui. Nota-se que a abordagem proposta apresenta melhores resultados quando comparada à aplicação direta da privacidade diferencial via Mecanismo Geométrico, ao mecanismo Log-Laplace e ao PrivBayes para as contagens de portas e serviços. Isto ocorre devido à abordagem proposta economizar orçamento de privacidade ϵ em todo o processo, não necessitando de divisão. Já para a contagem de protocolos (Figura 2b) nossa abordagem apresenta resultados um pouco inferiores a aplicação direta do Mecanismo Geométrico. Apesar disso, o erro médio relativo introduzido é bem baixo, isto é, na ordem de 10^{-3} quando $\epsilon = 0.5$, por exemplo. Vale ressaltar que existem apenas dois protocolos (TCP e UDP) compartilhados juntamente com suas respectivas contagens.

Os resultados no segundo conjunto de dados, “IP Network Traffic Flows” – Figura 3, apresentam um comportamento similar. Especificamente, quando $\epsilon = 0.5$, os valores para contagem de portas, protocolos e serviços são respectivamente 1.50, 0.0004 e 0.20. Quando comparamos com a aplicação direta da privacidade diferencial, obtêm-se respectivamente, quando o $\epsilon = 0.5$, valores 4.44, 1.33 e 0.57. Na aplicação do mecanismo Log-Laplace, os valores do erro absoluto médio são: 22838, 0.54 e 2160. Considerando

a aplicação do Privbayes, os valores são: 72, 0.12 e 707. A contagem de protocolos (Figura 3b) utilizando a abordagem proposta também apresenta resultados moderadamente superiores quando comparada ao mecanismo Geométrico. Novamente, mesmo com esses resultados, o erro médio relativo introduzido é bem baixo, isto é, na ordem de 10^{-3} quando $\varepsilon = 0.5$, por exemplo.

5.2. Avaliação dos Top- k Serviços e Portas

Contagens diferencialmente privadas podem ser utilizadas também para identificar os serviços ou portas mais frequentes, medindo os top- k registros mais utilizados em ambos os conjuntos de dados e avaliando a similaridade de Jaccard. Esses resultados são exibidos nas Figuras 4 e 5. Convém salientar que, para protocolos, os resultados referentes a top- k não foram reportados, visto que ambos os conjuntos de dados possuem apenas 2 tipos de protocolos: TCP e UDP. A similaridade de Jaccard é dada pela seguinte fórmula:

$$Jaccard(y, \hat{y}) = \frac{|y \cap \hat{y}|}{|y \cup \hat{y}|} \quad (3)$$

onde y e \hat{y} são os conjuntos de top- k original e ruidoso, respectivamente. Valores próximos de 1 sugerem que os dois conjuntos são muito similares, enquanto que valores próximos de 0 indicam que os conjuntos y e \hat{y} são mais disjuntos. Nestes resultados, fixamos o orçamento de privacidade em 0.1.

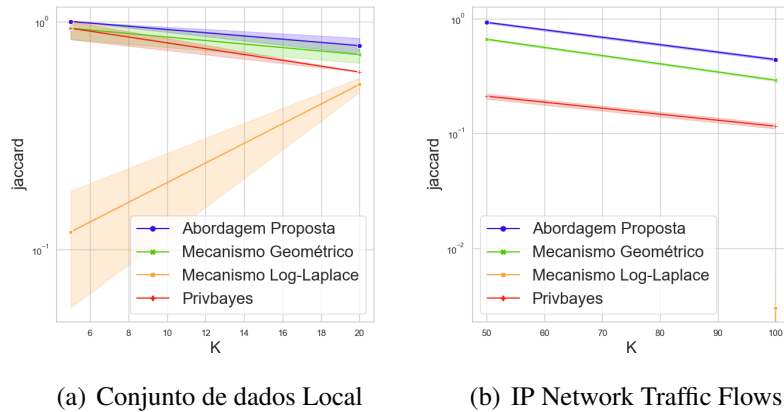


Figura 4. Similaridade de Jaccard para contagem de portas.

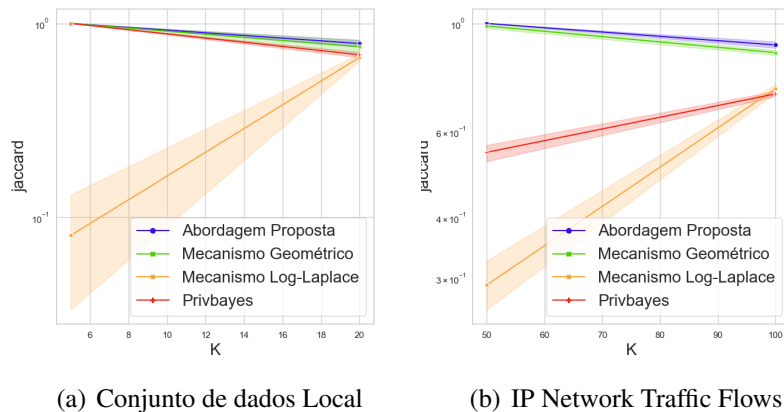


Figura 5. Similaridade de Jaccard para contagem de serviços.

Nota-se que a abordagem proposta neste trabalho obtém melhores resultados em comparação a outras técnicas e mecanismos diferencialmente privados. Dessa forma, nossa abordagem garante melhor utilidade dos dados mesmo quando são avaliados os top- k serviços e portas, para ambos os conjuntos de dados analisados.

6. Conclusão

Neste trabalho propomos uma abordagem diferencialmente privada para publicar dados de tráfego de rede com garantias formais de privacidade. Inicialmente, agrupamos os dados em triplas contendo portas, protocolos e serviços, para que a adição de ruído fosse executada apenas uma vez, ao invés de três vezes utilizando a abordagem mais usual de privacidade diferencial. Em seguida, o mecanismo Geométrico foi aplicado sobre as contagens de cada tripla pré-processada, a fim de introduzir ruído aleatório na contagem original. Por fim, obtemos as contagens agregadas por portas de destino, protocolos e serviços via pós-processamento de dados, isto é, o objetivo final da nossa abordagem. Resultados experimentais mostraram que a abordagem proposta adicionou menos ruído aos dados quando comparada a outras três técnicas existentes na literatura, proporcionando maior utilidade aos dados e, ao mesmo tempo, garantindo formalmente a privacidade dos tráfegos de rede. Como trabalhos futuros, pretendemos investigar informações privadas em outros domínios de redes, tais como duração do fluxo de tráfego, tamanho do fluxo e comprimento do fluxo de rede.

Agradecimentos

Este trabalho foi parcialmente financiado pela Lenovo, como parte de seu investimento em P&D pela lei de informática. Os autores agradecem ao CNPq (316729/2021-3), à CAPES (89723/2018-01, 88882.454584/2019-01) e ao LSB/D/UFC pelo financiamento parcial deste trabalho.

Referências

- Brito, F. T. and Machado, J. C. (2017). Preservação de privacidade de dados: Fundamentos, técnicas e aplicações. *Jornadas de atualização em informática*, pages 91–130.
- Bureau, U. C. (2021). Census bureau sets key parameters to protect privacy in 2020 census results. <https://www.census.gov/newsroom/press-releases/2021/2020-census-key-parameters.html>. Acesso 03 de fevereiro de 2023.
- Comissão Federal de Comunicação (2018). Privacidade do cliente. <https://www.fcc.gov/general/customer-privacy>. Acesso 17 de fevereiro de 2023.
- de Castro Vidal, I., da Costa Mendonça, A. L., Rousseau, F., and Machado, J. C. (2020). Protecting: An application of local differential privacy for iot at the edge in smart home scenarios. In *Anais do XXXVIII Simpósio Brasileiro de Redes de Computadores e Sistemas Distribuídos*, pages 547–560. SBC.
- Dwork, C. (2006). Differential privacy. In *International Colloquium on Automata, Languages, and Programming*, pages 1–12. Springer.
- GDPR (2018). General data protection regulation. <https://gdpr-info.eu/>. Acesso 17 de fevereiro de 2023.

- Ghosh, A., Roughgarden, T., and Sundararajan, M. (2009). Universally utility-maximizing privacy mechanisms. In *Proceedings of the forty-first annual ACM symposium on Theory of computing*, pages 351–360.
- IANA (2023). Internet Assigned Numbers Authority. <https://www.iana.org/>. Acesso em: 24 de março de 2023.
- Joshi, M. and Hadi, T. H. (2015). A review of network traffic analysis and prediction techniques. *arXiv preprint arXiv:1507.05722*.
- Koufogiannis, F. and Pappas, G. J. (2017). Diffusing private data over networks. *IEEE Transactions on Control of Network Systems*, 5(3):1027–1037.
- Le Ny, J. and Pappas, G. J. (2013). Privacy-preserving release of aggregate dynamic models. In *Proceedings of the 2nd ACM international conference on High confidence networked systems*, pages 49–56.
- LGPD (2019). Lei geral de proteção de dados pessoais. http://www.planalto.gov.br/ccivil_03/_Ato2019-2022/2019/Lei/L13853.htm. Acesso em: 17 de fevereiro de 2023.
- Lu, Y., Tian, H., Shen, H., and Xu, D. (2018). Privacy preserving classification based on perturbation for network traffic. In *Parallel and Distributed Computing, Applications and Technologies: 19th International Conference, PDCAT 2018, Jeju Island, South Korea, August 20-22, 2018, Revised Selected Papers 19*, pages 121–132. Springer.
- Machanavajjhala, A., He, X., and Hay, M. (2017). Differential privacy in the wild: A tutorial on current practices & open challenges. In *Proceedings of the 2017 ACM International Conference on Management of Data*, pages 1727–1730.
- McSherry, F. D. (2009). Privacy integrated queries: an extensible platform for privacy-preserving data analysis. In *Proceedings of the 2009 ACM SIGMOD International Conference on Management of data*, pages 19–30.
- Mimir (2018). Collection of User Data by ISPs and Telecom Providers, and Sharing with Third Parties. <https://www.ivpn.net/blog/collection-of-user-data-by-isps-and-telecom-providers-and-sharing-with-third-parties>. Acesso em: 28 de março de 2023.
- Sangeetha, S. and Sudha Sadasivam, G. (2019). Privacy of big data: a review. *Handbook of big data and iot security*, pages 5–23.
- Shafiq, M., Yu, X., Laghari, A. A., Yao, L., Karn, N. K., and Abdessamia, F. (2016). Network traffic classification techniques and comparative analysis using machine learning algorithms. In *2016 2nd IEEE International Conference on Computer and Communications (ICCC)*, pages 2451–2455. IEEE.
- Tanenbaum, A. S. and Wetherall, D. (2021). *Network Computer 6th Edition*. Pearson.
- Zhang, J. et al. (2017). Privbays: Private data release via bayesian networks. *ACM Transactions on Database Systems (TODS)*, 42(4):1–41.
- Zhang, X., Hamm, J., Reiter, M. K., and Zhang, Y. (2022). Defeating traffic analysis via differential privacy: a case study on streaming traffic. *International Journal of Information Security*, 21(3):689–706.