

Metodologia para publicação de dados com preservação da privacidade para órgãos públicos: de maneira proativa e solicitada

Bruno R. S. Moraes¹, Josenildo C. Silva², Ariel S. Teles³, Antonio A. B. Júnior¹,
Francisco J. S. Silva¹, Luciano R. Coutinho¹

¹Lab. de Sis. Distrib. Inteligentes (LSDi), Universidade Federal do Maranhão – UFMA
CEP 65.085-580 – Av. dos Portugueses, 1966 – Campus do Bacanga
São Luís – MA – Brazil

²Instituto Federal do Maranhão – IFMA
CEP 65.030-005 – Av. Getúlio Vargas, 04 – Monte Castelo – São Luís – MA – Brazil

³Instituto Federal do Maranhão – IFMA
CEP 65.570-000 – Rua José de Alencar – Bairro Comprida – Araiões – MA – Brazil

{brs.moraes, antonio.batista, francisco.silva, luciano.rc}@ufma.br,
{jcsilva, ariel.teles}@ifma.edu.br

Abstract. *The publication of government data enables transparency and scientific advancement. In Brazil, this publication is regulated by the LAI, may be mandatory or requested, and must comply with the LGPD. Data privacy is the Manager's responsibility, and suppressing explicit identifiers is insufficient to ensure privacy. This paper proposes a methodology that encompasses both forms of publication. In the mandatory, the Manager's ensures the protection of the record without data suppression. In the requested, the Miner's is included in the data anonymization process. A case study was conducted with public data, where it was possible to uniquely select 7,357 records. Applying the methodology, it was possible to create indistinguishable groups of size 10.*

Resumo. *A publicação de dados governamentais possibilita transparência e avanço científico. No Brasil, essa publicação é regulamentada pela LAI, pode ser obrigatória ou solicitada e deve estar em conformidade com a LGPD. A privacidade dos dados é de responsabilidade do Gestor, e a supressão de identificadores explícitos é insuficiente para garantir a privacidade. Este artigo propõe uma metodologia que abrange ambas as formas de publicação. Na obrigatória, o Gestor garante a proteção do registro sem a supressão dos dados. Na solicitada, o Minerador é incluído no processo de anonimização dos dados. Foi realizado um estudo de caso com dados públicos, onde foi possível selecionar unicamente 7.357 registros. Aplicando a metodologia, foi possível criar grupos indistinguíveis de tamanho 10.*

1. Introdução

A publicação de dados desempenha um papel importante para a transparência governamental, disponibilizando para o cidadão informações referentes a gastos públicos,

assim como ajuda a comunidade científica a analisar tendências, identificação de padrões que auxiliam em tomadas de decisões, impulsionando avanços em diversas áreas [Carvalho et al. 2023, Direito and Barros 2025]. No caso específico do Brasil, todos os órgãos e entidades públicas são obrigados a disponibilizar informações para o acesso público. Esta obrigatoriedade encontra-se regulamentada pela Lei nº 12.527/2011, Lei de Acesso à Informação (LAI), [Brasil 2011]. Entretanto, a lei supracitada considera de forma particular os dados pessoais, em seu Artigo 31, seção V, o qual afirma que o tratamento de dados pessoais deve respeitar a “intimidade, vida privada, honra e imagem das pessoas, bem como às liberdades e garantias individuais”. Portanto, a publicação de dados pelo setor público não deve violar a privacidade dos cidadãos [Queiroz and Motta 2015, Tejedo-Romero et al. 2025]. Neste contexto, temos a Lei Geral de Proteção de Dados (LGPD) [Brasil 2018], que dispõe sobre o tratamento de dados pessoais no Brasil, definindo diretrizes para a proteção da privacidade dos indivíduos. Portanto, na publicação de dados governamentais, por um lado, temos a LAI que obriga a divulgação dos dados. Por outro, temos a LGPD que exige a anonimização dos dados pelos meios técnicos razoáveis e disponíveis na ocasião do tratamento dos dados. As leis supracitadas contribuem para maior transparência tanto na publicação quanto no uso de dados pessoais. Porém, não é objetivo das regulamentações definir em termos técnicos quais tecnologias devem ser utilizadas para prover a privacidade. Elas abordam de maneira geral como os gestores dos dados – pessoa ou entidade considerada confiável que é responsável pelos dados – devem tratar os dados pessoais, cabendo a eles a escolha dos métodos de tratamento.

Na área da computação, há uma subárea denominada de Publicação de Dados com Preservação de Privacidade (PDPP), que visa criar um conjunto de dados anônimo que proteja a privacidade enquanto mantém níveis ótimos de utilidade dos dados. Esse objetivo pode ser alcançado por meio de técnicas de preservação da privacidade, como anonimização de dados, generalização e perturbação. Tais técnicas visam combater as ameaças de vincular um indivíduo a um registro ou a um atributo sensível. Essas ameaças são chamadas de vinculação de registro e vinculação de atributo, respectivamente [Ge et al. 2024]. De maneira geral, a PDPP é dividida em duas abordagens principais: a primeira visa diminuir o espaço de busca do conjunto de dados original, sem adicionar ruído, por meio dos operadores de generalização e/ou supressão, tendo como principais modelos de privacidade o k -anonimato [Samarati 2001] e ℓ -diversidade [Machanavajjhala et al. 2007]. A segunda abordagem adiciona, por meio da adição de ruído, uma imprecisão nos dados, tendo como modelo de referência a privacidade diferencial [Dwork 2006].

Gerenciar o equilíbrio entre privacidade e utilidade é desafiador, devido aos seus princípios serem inversamente proporcionais [Brito and Machado 2017]. Os gestores podem utilizar técnicas de privacidade para compartilhar dados, de maneira a garantir um certo nível de privacidade e assim não violar as regulamentações de proteção de dados. Entretanto, essas técnicas impactam a utilidade dos dados, o que pode ser um problema para o Minerador, pessoa que tem acesso aos dados publicados ou os solicita, já que deseja extrair o máximo de informações possível, e essa perda de utilidade pode impedir o acesso a certos dados relevantes.

A LAI regulamenta a publicação proativa de dados de órgãos públicos, mas

também prevê que a publicação pode ser solicitada. Nesse contexto, o Gestor lida com ambos os tipos de publicação e, independentemente de qual seja, deve respeitar a privacidade dos dados, conforme regulamentado pela LGPD.

Diante desta problemática, o objetivo principal deste artigo é propor uma metodologia para publicação de dados com preservação da privacidade, onde esta publicação pode ocorrer de maneira proativa e solicitada. Na proativa, o Gestor generaliza os dados de maneira que não seja possível identificar unicamente um registro. Na solicitada, principal contribuição deste trabalho, é a inclusão do Minerador no processo de anonimização dos dados. Nesta publicação, o Gestor, de maneira interativa com o Minerador, apresenta os atributos que necessitam ser manipulados para alcançar um nível de privacidade desejada. O Minerador, por sua vez, especifica as manipulações que podem ser realizadas sobre os dados, tentando impactar o mínimo possível a utilidade dos dados. Desta forma, o equilíbrio entre utilidade e privacidade fica a cargo do processo interativo entre o Gestor e o Minerador. Na publicação de dados, encontramos dados de diferentes naturezas: áudio, vídeo, imagem, tabulares e geoespaciais [Eynden 2011]. Entretanto, na metodologia proposta, investigamos a preservação da privacidade na publicação de dados tabulares estruturados. Em relação à privacidade, a publicação deve possuir proteção tanto do registro quanto do atributo.

2. Conceitos Básicos

Nesta seção, serão apresentados os conhecimentos básicos e os problemas envolvidos na publicação de dados.

2.1. Dados Tabulares

Dados tabulares são representados por tabelas com estrutura de linha e coluna, onde cada linha corresponde a uma tupla do conjunto de dados e as colunas aos atributos das tuplas. Esta representação é denominada de microdados, quando estes representam dados pessoais, cada tupla corresponde a uma pessoa e os atributos às suas características [Fung et al. 2010, Machado et al. 2019]. Os microdados podem ser classificados em 04 tipos [Fung et al. 2010], que são: **Identificador Explícito (IE)**: são atributos que identificam unicamente um registro. Exemplo: CPF, CNH; **Quase-Identificador (QI)**: são atributos que não são IE, mas podem criar um IE, quando se correlacionam com outros QI's. Exemplo: Data de nascimento, CEP; **Atributo Sensível (AS)**: consistem em informações sensíveis específicas da pessoa. Exemplo: doença, salário e situação de incapacidade; **Atributo Não Sensível (ANS)**: contém todos os atributos que não se classificam como IE, QI e AS.

Importante lembrar que, ao disponibilizar os microdados, é obrigatória a supressão dos identificadores explícitos. Entretanto, remover somente os identificadores explícitos mostrou-se ineficaz para prover a privacidade [Affonso and Sant'Ana 2017]. Os demais microdados podem ser utilizados para reidentificar indivíduos, ligando ou combinando com outros dados ou observando características únicas encontradas nos dados divulgados [Sweeney 2002, Karagiannis et al. 2024, Coelho et al. 2025].

2.2. Abordagem não perturbativa

A abordagem não perturbativa visa limitar o espaço de busca sobre os dados originais, porém, mantendo a veracidade em relação às consultas que podem ser realiza-

das sobre os dados. Isso é feito por meio do uso dos operadores de privacidade: generalização e/ou supressão [Fung et al. 2010, Carvalho et al. 2023]. O trabalho de [De Capitani Di Vimercati et al. 2012, Machado and Neto 2021] define que o grau de proteção desfrutado na publicação de dados utilizando esta abordagem pode ser mensurado por um valor numérico. O nível de proteção do registro está diretamente relacionado ao valor de k , que representa o tamanho do menor grupo formado associado aos quase-identificadores, tendo como referência o modelo de privacidade k -anonimato, proposto por [Samarati 2001]. O nível de proteção do atributo está relacionado ao valor de ℓ , que representa a menor diversidade de atributos sensíveis relacionados aos grupos criados por meio da obtenção de k . O modelo de privacidade referência é o ℓ -diversidade, proposto por [Machanavajjhala et al. 2007].

2.3. Abordagem perturbativa

A abordagem perturbativa utiliza o operador de ruído, ela permite a publicação de todo o conjunto de microdados, onde esses dados possuem um certo nível de imprecisão em relação aos dados originais [Aggarwal et al. 2008, Carvalho et al. 2023]. Quando comparamos a abordagem perturbativa com a não perturbativa, percebemos que a abordagem perturbativa, apesar de tornar os dados menos precisos, mantém a estrutura dos dados, preservando certas propriedades estatísticas que são perdidas no processo de generalização e supressão da abordagem não perturbativa [Fung et al. 2010]. Neste trabalho, adotamos como operador de ruído a privacidade diferencial, proposta por [Dwork 2006]. A privacidade diferencial é uma medida rigorosa de privacidade para análise de dados, que oferece garantias probabilísticas robustas, garantindo que a presença ou ausência de qualquer indivíduo em um conjunto de dados não afete significativamente o resultado de uma análise ou consulta, permitindo o compartilhamento de um conjunto de dados sem revelar informações individuais de cada indivíduo [Dwork 2006, Alves et al. 2024a, Alves et al. 2024b].

3. Trabalhos Relacionados

A literatura apresenta trabalhos relacionados à publicação de dados com preservação da privacidade. Nestes trabalhos, são apresentadas tanto soluções com sistematizações de procedimentos que o Gestor pode adotar para evitar a violação da privacidade no compartilhamento de dados quanto interfaces que permitem ao Gestor interagir diretamente com os dados, a fim de balancear a utilidade e a privacidade.

Em relação as soluções de sistematização de procedimentos, [Baloukas et al. 2024] apresenta um *framework* para apoiar o compartilhamento de dados pessoais para pesquisa científica integrando anonimização, avaliação de risco e geração de acordos de licença. Tendo como objetivo simplificar o processo de compartilhamento de dados entre organizações, de maneira que este compartilhamento ocorra de forma segura e em conformidade com as regulamentações do uso de dados pessoais. [Queiroz and Motta 2015] verificou a forma de anonimização dos dados adotada no setor público brasileiro, analisando os microdados do Instituto Nacional de Estudos e Pesquisas Educacionais Anísio Teixeira (INEP) referentes ao ano de 2013, e constatou que a anonimização adotada apresentava fragilidades. Neste trabalho foi demonstrado que a aplicação de formas sistemáticas de anonimização reduz o risco de reidentificação de indivíduos. Os modelos de anonimização utilizados foram o k -anonimato [Samarati 2001],

para garantir a proteção do registro criando grupos indistinguíveis com base nos quase-identificadores, e o *Distinct ℓ -diversidade* [Machanavajjhala et al. 2007], utilizado para a proteção dos atributos sensíveis.

As soluções baseadas em interfaces possibilitam o Gestor manipular diretamente os dados a serem publicados, utilizando os operadores de privacidade: generalização, supressão e adição de ruído. [Kim and Kim 2024] apresenta uma metodologia dentro do sistema KMBIG para o compartilhamento de dados de saúde. No sistema, o responsável pela publicação interage com os dados por meio de uma interface que permite realizar operações de anonimização. Após a manipulação dos dados, é gerado um relatório que é avaliado por um comitê, de forma a definir se o nível de privacidade é aceitável para publicação. [Wang et al. 2018] aborda a anonimização de dados tabulares multiatributo, e propõe uma interface visual combinada com um pipeline de manipulação de dados. O sistema visual é composto por duas funcionalidades visuais: PER-Tree (*Privacy Exposure Risk Tree*) e UPD-Matrix (*Utility Preservation Degree Matrix*). A PER-Tree apresenta ao usuário os problemas de privacidade, permitindo que ele possa aplicar os operadores de privacidade para sanar essas fragilidades. A UPD-Matrix fornece um *feedback* de maneira visual dos impactos destas manipulações sobre a utilidade dos dados. A interface auxilia o usuário, responsável pela publicação, encontrar de maneira interativa um equilíbrio entre a proteção da privacidade e a utilidade dos dados. [Abu Attieh et al. 2024] descreve a criação e avaliação da ORCHESTRA *Pseudonymization Tool* (OPT), em redes de pesquisa biomédica. Oferecendo funcionalidades de pseudonimização e gerenciamento de dados de participantes e amostras biológicas. A pseudoanonimização ocorre através da substituição de dados identificáveis por pseudônimos, garantindo a proteção da privacidade, mas permitindo a vinculação de dados para análises científicas. A interação do usuário para a manipulação dos dados se dá via interface gráfica.

Podemos notar que os trabalhos de [Baloukas et al. 2024], [Abu Attieh et al. 2024] e [Queiroz and Motta 2015] apresentam uma perspectiva onde o Gestor deve considerar uma estrutura composta por etapas para publicar dados, de modo que não viole as regulamentações que tratam do uso e compartilhamento de dados pessoais. Enquanto os propostos por [Wang et al. 2018], [Abu Attieh et al. 2024] e [Kim and Kim 2024] apresentam uma perspectiva na qual o Gestor interage com os dados a serem publicados por meio de uma interface, manipulando-os com operadores de privacidade (generalização, supressão e perturbação).

Ambas as perspectivas objetivam garantir que os dados publicados não violem a privacidade, tentando impactar ao mínimo a utilidade dos dados. Entretanto, fica nítido que em ambas as perspectivas o Minerador fica alheio ao processo de anonimização de dados, não interagindo com o Gestor e algumas manipulações realizadas por ele podem inviabilizar a utilidade dos dados para os fins desejados pelo Minerador.

Mediante a problemática, este trabalho propõe uma metodologia para a publicação de dados, abrangendo tanto a forma proativa, quando o compartilhamento de dados obrigatório, quanto a solicitada. Na publicação solicitada, uma das principais contribuições deste trabalho, a anonimização dos dados ocorrerá de forma interativa entre o Gestor e o Minerador, buscando o equilíbrio entre utilidade e privacidade por meio dessa interação.

4. Metodologia

Diante desse desafio a metodologia proposta para publicar dados com preservação da privacidade opera de duas maneiras: *proativa* e *solicitada*. A maneira *proativa* ocorre quando órgãos e/ou entidades públicas são obrigados, por determinação legal, a disponibilizar informações. Na *solicitada*, órgãos e/ou entidades disponibilizam dados mediante solicitações específicas. A Figura 1 ilustra de maneira geral a metodologia de publicação de dados proativa e solicitada. Na publicação proativa, o Gestor, ao publicar os dados, não sabe quem irá acessá-los, o minerador é desconhecido, onde se faz necessário aplicar um maior nível de privacidade sobre os dados.

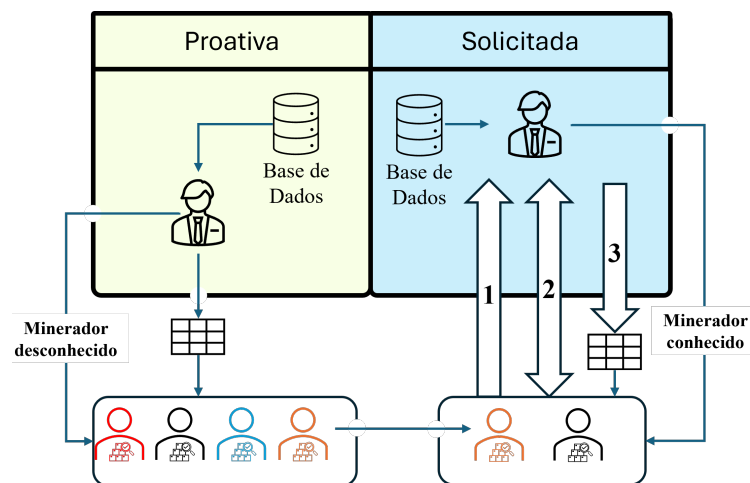


Figura 1. Visão geral da metodologia: proativa e solicitada

Na solicitada, o Minerador que teve acesso aos dados da publicação proativa e deseja uma melhor qualidade dos dados ou um Minerador que deseje acessar um determinado conjunto de dados realiza uma solicitação de publicação ao Gestor dos dados (1). Após, inicia-se um processo interativo entre o Gestor e o Minerador (2), até alcançar um equilíbrio entre a privacidade desejada pelo Gestor e o nível aceitável de utilidade para o Minerador. Quando este equilíbrio é alcançado, os dados são disponibilizados ao Minerador (3). Importante lembrar que o Minerador não tem acesso aos dados originais durante o processo interativo, somente especifica as manipulações que podem ser realizadas sobre os dados. Nas subseções 4.1 e 4.2 a seguir, será apresentada em mais detalhes a publicação proativa e solicitada.

4.1. Publicação Proativa (Publicação obrigatória)

Na publicação proativa, o Gestor tem a obrigação legal de disponibilizar dados. Essa exigência, imposta por regulamentações e/ou leis vigentes, determina que órgãos públicos tornem as informações sob sua guarda acessíveis, assegurando o direito fundamental ao acesso à informação. No entanto, é fundamental considerar outras normativas legais que regem a publicação de dados pessoais, as quais estabelecem regras e princípios para garantir a privacidade e a proteção dessas informações.

A Figura 2 ilustra a visão mais detalhada da publicação proativa com preservação da privacidade. Esse processo é composto por três etapas principais, divididas em seis passos.

– Etapa 1: composta de 04 passos - **Seleção dos dados, classificação dos metadados, exclusão do IE e a adição do Identificador de Rastreabilidade (IR).** Passo 1 - Seleção dos dados: este passo consiste no Gestor selecionar em sua base local o conjunto de dados solicitado. Passo 2 - Classificação dos metadados: neste passo o Gestor classifica as colunas do conjunto de dados em IE, QI, AS e ANS. Passo 3 - Excluir IE: neste passo o Gestor exclui as colunas classificadas no passo 2 como IE. Passo 4 - Adição do IR: neste passo é criado e adicionado um identificador de rastreabilidade para cada registro do conjunto de dados. O Identificador de Rastreabilidade (IR) consiste em um relacionamento criado entre os dados originais e os dados anonimizados. Essa relação é para possibilitar ao Gestor rastrear os dados que originaram os dados anonimizados, caso eles sejam solicitados. Assim como possibilita auditorias, caso se façam necessárias.

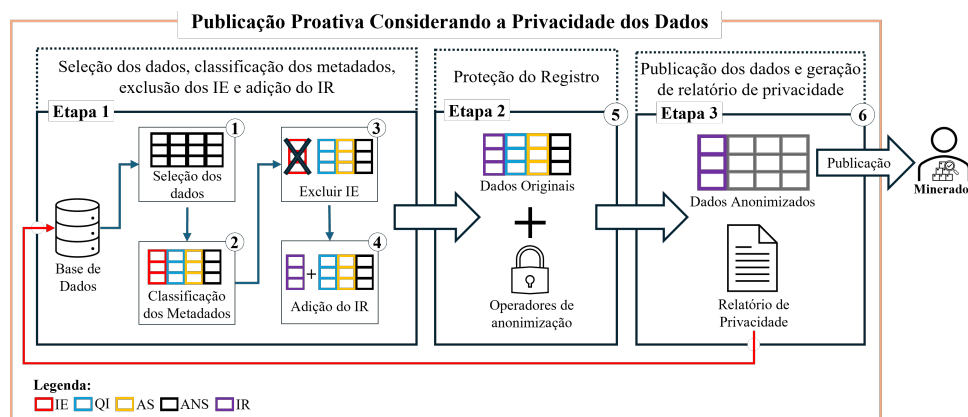


Figura 2. Etapas da publicação proativa

– Etapa 2: **Proteção do Registro** corresponde ao Passo 5. Nesta fase, aplicam-se operações de generalização ou mascaramento sobre os quase-identificadores com o objetivo de criar grupos indistinguíveis. Um grupo é considerado indistinguível quando um conjunto de quase-identificadores se torna idêntico para múltiplos registros no conjunto de dados, impossibilitando a identificação direta de um registro individual. O tamanho mínimo que todos os grupos devem atender é definido pelo Gestor, sendo obrigatório que este tamanho seja maior ou igual a 2. Após aplicar as manipulações de generalização os grupos com tamanho menores que o desejado pelo Gestor terão seus quase-identificadores mascarados.

– Etapa 3: **Publicação dos Dados e Geração de Relatório de Privacidade**, que corresponde ao Passo 6, envolve duas ações principais. Primeiramente, é gerado um relatório de privacidade que descreve as manipulações aplicadas ao conjunto de dados, sendo essas informações registradas em sua base de dados local. Em seguida, o conjunto de dados anonimizado é disponibilizado de forma pública.

A Figura 3 ilustra o comparativo entre a estrutura visual da tabela com os dados originais e os anonimizados de maneira proativa. De forma que **CPF** foi classificado como IE, sendo excluído; **CEP** e **Data de Nasc** classificados como QI, tiveram parte de suas informações mascaradas. Os QI's menores que o grupo desejado serão mascarados. O grupo indistinguível formado tem tamanho 2.

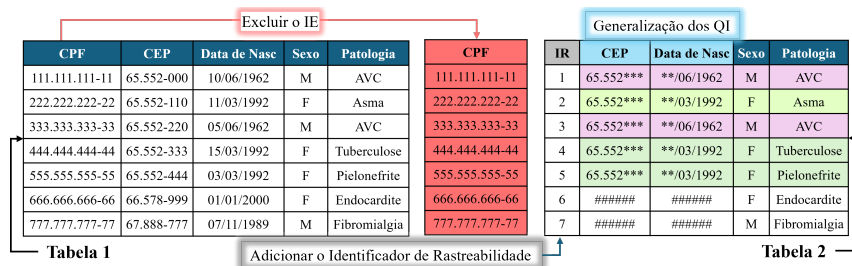


Figura 3. Ilustração visual entre a tabela original e a anonimizada

4.2. Publicação Solicitada (Processo Iterativo)

A Figura 4 ilustra as etapas do processo iterativo para publicação de dados com preservação da privacidade. Este processo é caracterizado como iterativo devido ao decorrer do processo de anonimização, o Minerador interage com o Gestor especificando as manipulações que podem ser realizadas sobre os dados e o Gestor, por sua vez, define o nível de privacidade que a publicação deve possuir. O equilíbrio entre privacidade e utilidade fica a cargo deste processo iterativo entre ambos. A publicação interativa é dividida em três etapas principais: proteção do registro, proteção do atributo e análise da utilidade.

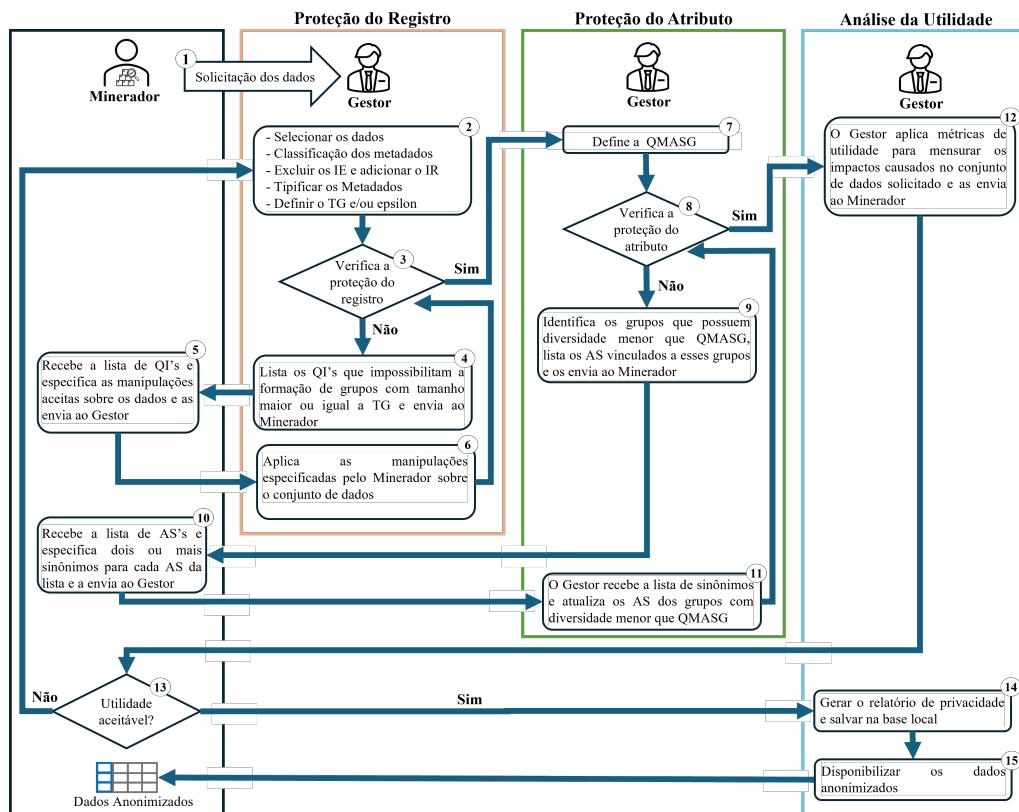


Figura 4. Publicação Interativa

Agora detalharemos os processos que constituem a publicação interativa ilustrada na Figura 4. Inicialmente é executado **processo 1**, onde o Minerador realiza a solicitação dos dados. Nesta solicitação consta: a *Identificação do Minerador*, a *Especificação dos Dados Solicitados* e a *Dimensão de Interesse*. Na *Identificação do Minerador*, ele informa seu tipo, que pode ser: pesquisador de instituição pública ou privada. Além disso,

deverá fornecer seu nome, endereço, número de identificação, tipo de identificação e anexar outras informações pertinentes. O tipo de identificação refere-se a um documento com validade legal e reconhecimento pelas autoridades competentes do território onde o minerador solicita os dados. Anexar outras informações pertinentes refere-se ao termo de responsabilidade, no qual o Minerador se compromete a não violar e/ou compartilhar os dados, bem como outros documentos que se fizerem necessários. Na *Especificação dos Dados Solicitados*, o Minerador informará um rótulo para o dado solicitado, juntamente com sua descrição. Exemplo: Data de nascimento – dado referente ao dia do nascimento da pessoa. No exemplo, o rótulo é “Data de nascimento” e a descrição é “dado referente ao dia do nascimento da pessoa”. Na *Dimensão de Interesse* o Minerador indica a dimensão de interesse em relação aos dados, podendo ser: Estrutura, Veracidade e Híbrida. Estrutura: garante a integridade de cada registro, não permitindo a operação de supressão. A privacidade é provida por meio da adição de ruído. Veracidade: emprega operadores de generalização e/ou supressão para reduzir o espaço de busca, mas assegurando que as informações obtidas sejam consistentes com as originais. Híbrida: pode utilizar os operadores de generalização, supressão e adição de ruído.

No **processo 2**, ao receber a solicitação dos dados, o Gestor executará os seguintes passos: *selecionar os dados*, *classificação dos metadados*, *excluir os IE e adição do IR*, *tipificação dos metadados* e *Definir o TG (Tamanho do Grupo) e/ou epsilon*. Em *Selecionar dos dados*: o Gestor seleciona o conjunto de dados solicitados pelo Minerador em sua base de dados local. Na *Classificação dos Metadados*: o Gestor classifica os metadados, representados pelos rótulos do conjunto de dados definidos em *Solicitação de Dados*, em IE, QI, AS e ANS. *Excluir os IE e adicionar o IR*: o Gestor suprime os IE's, cria e adiciona um IR para cada registro. Na *tipificação dos metadados*: o Gestor tipifica os QI e AS em numérico ou categórico. Aqueles classificados como numéricos devem ser especificados como contínuo (real) ou discreto (inteiro). *Definir o TG e/ou epsilon (ϵ)*: TG é o valor definido pelo Gestor para a melhor consulta associada aos QI's sobre o conjunto de dados solicitados, sendo esta consulta obrigatoriamente maior ou igual a 2. De modo que, os registros retornados pela consulta são indistinguíveis entre si. Por exemplo, se o valor de TG for 5, independentemente das combinações dos QI's, a consulta sempre retornará 5 ou mais registros, de forma que não seja possível diferenciá-los analisando os QI's. O TG é definido quando a dimensão de interesse for a veracidade. Quando a dimensão de interesse for a estrutura, será utilizado o operador de adição de ruído. Nesta proposta, utiliza-se a privacidade diferencial aplicada aos QI's, com o mecanismo de Laplace para metadados numéricos e o exponencial para os categóricos. Para esta dimensão, o Gestor definirá o valor do ϵ que será o nível de privacidade adicionado aos dados. Na híbrida, o Gestor combina princípios das dimensões de veracidade e estrutura, onde ele define o valor de TG e ϵ .

No **processo 3**, é verificada a proteção do registro. Este processo ocorre de forma automática, onde verifica se os QI's formam grupos maiores ou iguais ao TG. Caso sim, avança-se para o processo 7. Caso contrário, os dados serão manipulados para se alcançar a proteção do registro. As manipulações a serem realizadas sobre os dados variam de acordo com a dimensão de interesse escolhida pelo Minerador. A seguir, vamos detalhar como ocorre a manipulação dos dados de acordo com a dimensão de interesse.

– **Veracidade**: nesta dimensão é executado o fluxo ilustrado na Figura 4, inicia-se

com o **processo 4** onde o Gestor lista os QI's que impossibilitam a formação de grupos com tamanho maior ou igual a TG e envia ao Minerador. O Minerador, ao receber a lista de QI's, especifica as manipulações aceitas sobre os dados e as envia ao Gestor (**processo 5**). O Gestor, ao receber as especificações, aplica as manipulações sobre o conjunto de dados (**processo 6**) e verifica a proteção do registro (**processo 3**). Este processo iterativo continua até que os QI's formem grupos maiores ou iguais a TG. Um ponto a considerar neste processo iterativo é a possibilidade de o Minerador negociar com o Gestor o nível de privacidade (valor de TG) e o Gestor, por sua vez, solicitar algumas manipulações sobre os dados.

– **Estrutura:** o Gestor com base no valor de ϵ aplica os mecanismos de privacidade diferencial sobre os QI's. Em seguida, é avançado para o **processo 7**. A Figura 5 ilustra como ocorre o processo de adição de ruído via privacidade diferencial. Neste ponto, os QI's encontram-se tipificados. Para os atributos numéricos, é aplicado o mecanismo de laplace, onde é adicionado ao valor original um valor gerado pelo mecanismo. Este valor é obtido por meio da sensibilidade da função, que consiste no quociente do valor máximo onde a saída da função pode variar quando um único registro é adicionado ou removido do conjunto de dados, pelo ϵ .

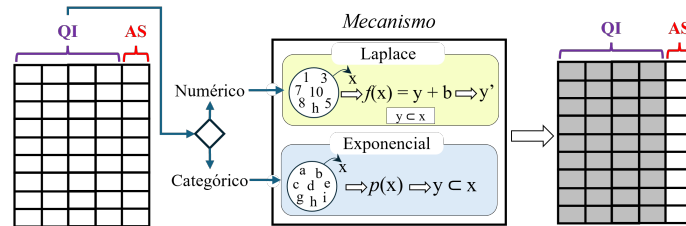


Figura 5. Adição de ruído nos quase-identificadores via privacidade diferencial

Para os atributos categóricos é aplicado o mecanismo exponencial, este mecanismo seleciona de maneira probabilística um valor dentro do conjunto de opções dos QI's para aquele atributo em específico. O cálculo da probabilidade é ponderada, levando em consideração a frequência do atributo e o valor de ϵ .

– **Híbrida:** ilustrada na Figura 6. Inicialmente, o Gestor seleciona os QI's que impossibilitam a formação de grupos maiores ou iguais a TG e os envia ao Minerador. O Minerador, ao receber os QI's, especifica as manipulações que ele aceita sobre os dados e as retorna ao Gestor.

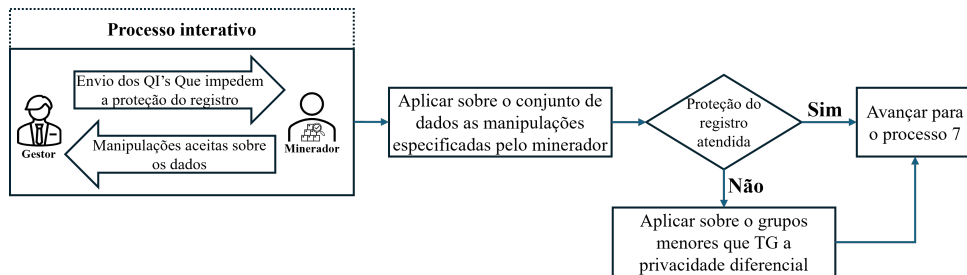


Figura 6. Ilustração do fluxo geral do processo híbrido

O Gestor aplica sobre os QI's as manipulações especificadas pelo Minerador. Em seguida, verifica se a proteção do registro foi alcançada. Se sim, avança-se para o **pro-**

processo 7. Caso contrário, é aplicada a privacidade diferencial sobre os QI's dos grupos menores do que TG e avança para o **processo 7**.

Neste ponto foi alcançada a proteção do registro. No **processo 7**, o Gestor inicialmente define o valor de QMASG (Quantidade Mínima de Atributos Sensíveis por Grupo), sendo que o menor valor que QMASG pode assumir é 2. No **processo 8**, é verificado de maneira automática se cada grupo formado pelos QI's possui uma diversidade de AS maior ou igual a QMASG. Se sim, avançamos para o **processo 12**. Caso contrário, os AS necessitarão ser diversificados. O processo de diversificação varia de acordo com a dimensão escolhida pelo Minerador.

– **Veracidade:** o Gestor identifica os grupos que possuem diversidade menor que a QMASG, e lista os AS vinculados a esses grupos e a envia ao Minerador para que ocorra o processo de diversificação (**processo 9**). O Minerador, ao receber a lista de AS inicia o processo de diversificação. Este processo consiste em atribuir a cada AS da lista dois ou mais sinônimos. Ao término, o Minerador terá associado a cada elemento da lista de AS enviadas pelo Gestor uma lista de sinônimos. Um sinônimo, neste contexto, é uma palavra semelhante ao AS. Essa lista é enviada ao Gestor para atualizar os AS dos grupos com diversidade menor que a QMASG (**processo 10**). O Gestor, ao receber a lista de sinônimos, atualiza os AS dos grupos menores que a QMASG (**processo 11**) e verifica a proteção do atributo (**processo 8**). Este processo iterativo repete-se até que os grupos formados pelos QI's possuam uma diversidade de AS maior ou igual a QMASG. Vale ressaltar que a QMASG pode ser negociada durante o processo iterativo.

– **Estrutura:** o Gestor com base no valor de ϵ aplica os mecanismos de privacidade diferencial sobre os AS dos grupos indistinguíveis, conforme ilustrado na Figura 7, com diversidade menor que a QMASG. Em seguida, é avançado para o **processo 12**.

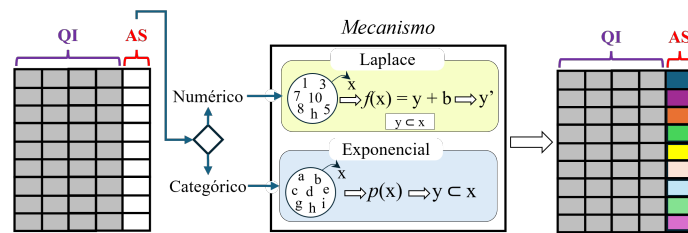


Figura 7. Diversificação dos atributos sensíveis via privacidade diferencial

– **Híbrida:** o Gestor envia uma lista de AS que impossibilitam a proteção do atributo. O Minerador diversifica os AS e devolve ao Gestor. O Gestor atualiza os AS dos grupos indistinguíveis menores que QMASG e verifica a proteção do atributo. Para os grupos indistinguíveis que possuírem a QMASG menores que o desejado após a atualização, será aplicada a privacidade diferencial, com base no valor de ϵ . Em seguida, avança-se para o **processo 12**. Este processo está ilustrado na Figura 8.

Neste ponto, foi alcançada a proteção do atributo. No **processo 12**, o Gestor utiliza métricas de utilidade para mensurar os impactos causados no conjunto de dados solicitado, resultantes das manipulações realizadas, e as envia ao Minerador.

No **processo 13**, o Minerador analisa as métricas e verifica se os dados mantiveram uma utilidade aceitável. Caso sim, o Minerador informa ao Gestor que deseja ter

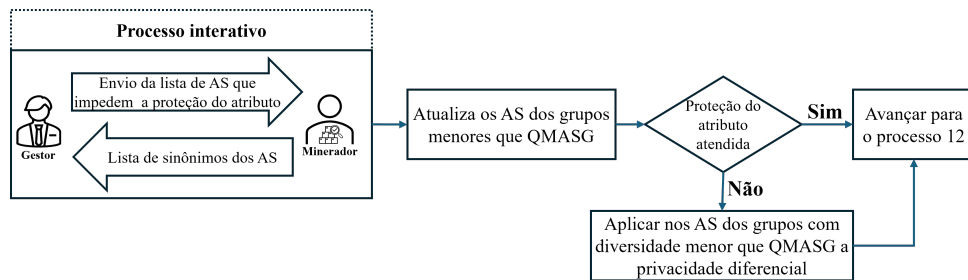


Figura 8. Processo de diversificação dos atributos, dimensão híbrida

acesso aos dados anonimizados, e avança-se para o **processo 14**. Caso contrário, inicia-se novamente o processo de proteção do registro.

No **processo 14**, o Gestor gera o relatório de privacidade que consiste na descrição dos procedimentos adotados no processo de anonimização de dados, assim como as informações de quem solicitou os dados.

No **processo 15**, consiste na disponibilização dos dados anonimizados ao Minerador.

5. Cenário de uso da metodologia utilizando dados públicos

Para demonstrar o uso da metodologia, utilizamos dados públicos referentes aos registros de vacinação COVID-19. Estes dados são publicados pelo Ministério da Saúde (SUS), por meio do portal OpenDataSUS, que disponibiliza informações que podem servir para subsidiar análises objetivas da situação sanitária, tomadas de decisão baseadas em evidências e elaboração de programas de ações de saúde. O conjunto de dados do COVID-19 é dividido por estado, podendo ser acessados por meio de uma API ou arquivo no formato CSV, disponíveis no portal¹. O motivo da escolha destes dados é o fato de sua publicação ser obrigatória. Analisamos o conjunto de dados publicados e identificamos que algumas colunas filtravam registros, agrupando-os em pequenos grupos ou até mesmo de maneira única. Diante dessa fragilidade, aplicamos a metodologia proposta para aumentar o nível de privacidade. O cenário de aplicação será apresentado para a publicação proativa e solicitada.

5.1. Verificação de fragilidades em relação a privacidade do registro

Inicialmente, foi selecionado um município do estado do Maranhão, o conjunto de dados possuía 11.543 registros, e realizou-se a classificação dos metadados, como IE (*paciente_id*), QI (*paciente_idade*, *paciente_dataNascimento*, *paciente_enumSexoBiologico*, *paciente_racaCor_valor*, *paciente_endereco_cep*, *vacina_dataAplicacao*) e AS (*vacina_grupoAtendimento_nome*). Em seguida, foram removidos os atributos classificados como IE. Após a remoção, verificou-se se os QI's permitiam selecionar pequenos grupos ou até mesmo um único registro. A utilização do QI *paciente_dataNascimento* (ano/mês/dia) resultou na seleção de 7.357 registros únicos, indicando que cada um deles possui uma data de nascimento distinta no conjunto de dados. Os QI's *paciente_endereco_cep*, *vacina_dataAplicacao* e *paciente_idade* selecionavam de maneira única 304, 69 e 3, respectivamente.

¹ <https://opendatasus.saude.gov.br/dataset/covid-19-vacinacao>

5.2. Publicação Proativa

Inicialmente, os dados foram selecionados, excluiu-se o IE *paciente_id*, sendo criado e adicionado a cada registro um IR. Após, realizaram-se operações de anonimização sobre os dados, que foram: para *paciente_idade*, foi adicionado um período de 10 em 10 anos; para *paciente_endereco_cep*, foram mascarados os últimos 05 caracteres; e para *vacina_dataAplicacao*, foram mascaradas as informações referentes ao dia e mês. O QI *paciente_dataNascimento* foi suprimido. Definiu-se que o nível de proteção do registro seria 10. Em seguida, verificou-se se os grupos formados pela generalização e mascaramento alcançaram o tamanho 10. Constatou-se que, com a combinação de todos os QI's, era possível selecionar de maneira única 468 registros. Assim, optou-se por mascarar os QI's dos registros que não formam grupos de tamanho 10. Para avaliação da utilidade dos dados, foram aplicados sobre os dados originais e anonimizados o modelo de classificação J48 [Quinlan 1993], para realizar a comparação da acurácia obtida. O modelo de classificação foi avaliado utilizando o Weka 3.8.6², com validação cruzada 10-fold.

5.3. Publicação Solicitada

Para o experimento da publicação solicitada, temos o seguinte cenário. O Minerador, ao acessar os dados publicados pelo Gestor na publicação proativa, deseja um maior detalhamento dos dados. Cada processo será apresentado de maneira explícita.

Processo 1: O Minerador solicita ao Gestor os dados originais referentes à publicação disponibilizada. Na solicitação, o Minerador informa seus dados pessoais, envia o termo de responsabilidade de não violar a privacidade dos dados e especifica que a dimensão de interesse é a veracidade.

Processo 2: O Gestor, ao receber a solicitação, por meio do IR, acessa os dados originais. A seguir, classifica e suprime os IE. Posteriormente, gera e adiciona aos dados novos IR e associa-os aos dados originais, tipifica os dados e define TG de 7.

Processo 3: O Gestor verifica se os dados originais formam grupos maiores ou iguais a TG. Ao realizar a verificação, constatou que *paciente_idade*, *paciente_dataNascimento* e *paciente_endereco_cep* formavam 13, 9263 e 452 grupos menores que TG. A combinação de todos os QI's formava 11.275 grupos menores que TG. Além de *paciente_idade* e *paciente_dataNascimento* selecionavam 03 e 7.357 registros de maneira única.

– Início do processo iterativo para alcançar a proteção do registro

Processo 4: O Gestor gera uma lista dos QI's que impossibilitam formar TG e envia ao Minerador.

Processo 5: O Minerador recebe a lista dos QI's que impossibilitam alcançar TG de 7. O Minerador especificou as seguintes operações: suprimir o dia de *paciente_dataNascimento*; definir intervalo de 5 anos para *paciente_idade*; suprimir os últimos 4 dígitos de *paciente_endereco_cep*; e suprimir o dia de *vacina_dataAplicacao*. Após, enviou-as ao Gestor.

²<https://ml.cms.waikato.ac.nz/weka/index.html>

Processo 6: (Gestor) Aplicou sobre os dados as especificações definidas pelo Minerador. Porém, ainda era possível selecionar de maneira única um registro pelo QI paciente_idade (**Processo 3**). Informando ao Minerador.

- **Minerador:** Poderia diminuir o valor de TG.

- **Gestor:** Sim, caso especifique novas manipulações, sugiro que a idade fique de 10 em 10 anos e onde tiver data verificar se pode ficar apenas o ano.

Processo 5: (Minerador) Especifica que as novas manipulações aceitas sobre os dados são: para paciente_idade, aplicar intervalo de 10 anos; para paciente_dataNascimento, manter somente o ano; para vacina_dataAplicacao, manter somente o ano; e para paciente_endereco_cep, suprimir os últimos 5 dígitos.

Processo 6: (Gestor) Definiu o novo TG de 5, aplicou as manipulações definidas pelo Minerador e verificou se os grupos formados são maiores ou iguais ao valor de TG (**Processo 3**). Mesmo aplicando as manipulações e diminuindo o TG para 5, a combinação de todos os QI's apresenta 2080 grupos menores que TG.

- **Minerador:** Podes baixar o TG para 3.

- **Gestor:** Sim. Entretanto, os grupos menores que 3 serão suprimidos. Aceita esta condição?

- **Minerador:** Sim. Condição aceita.

Processo 6: (Gestor) Atualizou o TG para 3 e excluiu 2126 registros. Alcançando a proteção do registro.

– Fim do processo iterativo para alcançar a proteção do registro

Processo 7: O Gestor verificou a diversidade de cada grupo formado, houve alguns grupos que apresentaram diversidade menor que 2. Porém, o AS associado foi: Faixa Etária, Pessoas de 12 a 17 anos, Pessoas de 18 a 64 anos, Pessoas de 5 a 11 anos, Pessoas de 65 a 69 anos, Pessoas de 70 a 74 anos, Pessoas de 75 a 79 anos e Pessoas de 80 anos ou mais. Não foram encontradas informações como Pneumopatias Crônicas Graves, Neoplasias, Doenças Cardiovasculares e Cerebrovasculares, que estão presentes no conjunto de dados. Como o Minerador está fornecendo garantias de que não irá violar os dados, não será necessário manipular os atributos. Avançando para o processo 12.

Processo 12: O Gestor aplica o modelo de classificação J48 e compara a acurácia obtida dos dados originais e anonimizados, obtendo uma diferença de 0,2256%.

Processo 13: O Minerador analisa os resultados enviados pelo Gestor, e os dados satisfazem suas necessidades.

Processo 14: O Gestor gera o relatório de privacidade, contendo todas as manipulações realizadas sobre os dados.

Processo 15: O Gestor disponibiliza os dados ao Minerador.

5.4. Resultados e discussão

Na **publicação proativa**, as operações de generalização e mascaramento possibilitaram publicar os dados com um grupo indistinguível de 10, com a supressão do *paciente_dataNascimento*. Isso implica que, para uma melhor consulta sobre o conjunto de

dados, haverá 10 registros associados, dificultando assim os ataques de reidentificação. Importante destacar que o IR possibilita o Gestor rastrear os dados originais, caso necessário. Na **publicação solicitada**, com a interação do Minerador, houve a preservação do QI *paciente_dataNascimento*. O processo interativo é importante, pois possibilita ao Minerador especificar manipulações que causem menos impacto nos dados e preservem mais a utilidade, para a finalidade que deseja. Este processo também abre possibilidades entre o Gestor e o Minerador de negociarem o nível de privacidade e utilidade. Como o Minerador, por meio de documentos formais, garante que não irá violar a privacidade, o Gestor pode flexibilizar o nível de privacidade para que os dados preservem mais utilidade. A Tabela 1 apresenta a acurácia obtida por meio do modelo de classificação J48. Percebe-se que a publicação solicitada apresentou uma acurácia próxima da obtida com os dados originais.

Tabela 1. Comparação da acurácia obtida do dado original e os obtidos pelo uso da metodologia, com uso do modelo de classificação J48.

Dados Original	Metodologia Proativa	Metodologia Solicitada
84.1753%	79.9653%	83.9497%

O resultado obtido com os dados da publicação proativa obteve uma acurácia menor, porém esta publicação possui um nível de privacidade maior, com grupo indistinguível de 10.

6. Conclusão

Neste trabalho apresentamos uma metodologia para publicação de dados públicos de maneira proativa ou solicitada com preservação da privacidade, para auxiliar o Gestor no processo de compartilhamento de dados. Na proativa, o Gestor adiciona privacidade aos dados formando grupos indistinguíveis. Na solicitada, a metodologia inclui o Minerador no processo de anonimização dos dados, onde pode especificar as manipulações que causem menor impacto nos dados. Este processo interativo entre Gestor e Minerador é o responsável pelo equilíbrio entre utilidade e privacidade. O experimento demonstrou que o processo interativo possibilita que o Minerador tenha dados mais úteis do que na publicação proativa. Independente da publicação, o identificador de rastreabilidade possibilita o acesso aos dados originais, para fins de solicitação e/ou auditorias. Para trabalhos futuros, aplicaremos a metodologia em empresas e hospitais para avaliar sua eficácia e levantar pontos para refinamento da metodologia.

Agradecimentos

Os autores agradecem ao Programa de Pós-Graduação em Ciência da Computação - Associação UFMA/UFPI. Este projeto conta com financiamento da FAPEMA/EMAP, processo número APP-09405/22, processo FAPESP número 23/00811-0, CNPq [processos 408548/2023-1, 308059/2022-0 e 441817/2023-8] e do Programa Porto do Futuro, Bolsa de Doutorado BD-08777/22.

Referências

Abu Attieh, H., Neves, D. T., Guedes, M., Mirandola, M., Dellacasa, C., Rossi, E., and Prasser, F. (2024). A scalable pseudonymization tool for rapid deployment in large

- biomedical research networks: Development and evaluation study. *JMIR Med Inform*, 12:e49646.
- Affonso, E. P. and Sant’Ana, R. C. G. (2017). PRESERVAÇÃO DA PRIVACIDADE NO ACESSO A DADOS POR MEIO DO MODELO K-ANONIMATO. *PontodeAcesso*, 11(1):20–41.
- Aggarwal, C. C., Yu, P. S., Elmagarmid, A. K., and Sheth, A. P., editors (2008). *Privacy-Preserving Data Mining: Models and Algorithms*, volume 34 of *Advances in Database Systems*. Springer US, Boston, MA.
- Alves, A. G. M., Pereira, F., Chaves, I., and Machado, J. (2024a). Privacidade diferencial em gradient boosting decision trees com técnicas de particionamento para dados categóricos. In *Anais do XXXIX Simpósio Brasileiro de Bancos de Dados*, pages 444–456, Porto Alegre, RS, Brasil. SBC.
- Alves, V., Costa, J., Gonzalez, L., Souza, A., and Villas, L. (2024b). Seleção de clientes adaptativa baseada em privacidade diferencial para aprendizado federado. In *Anais Estendidos do XLII Simpósio Brasileiro de Redes de Computadores e Sistemas Distribuídos*, pages 225–232, Porto Alegre, RS, Brasil. SBC.
- Baloukas, C., Papadopoulos, L., Demestichas, K., Weissenfeld, A., Schlarb, S., Aramburu, M., Redó, D., García, J., Gaines, S., Marquenie, T., Eren, E., and Erdogan Peter, I. (2024). A risk assessment and legal compliance framework for supporting personal data sharing with privacy preservation for scientific research. In *Proceedings of the 19th International Conference on Availability, Reliability and Security, ARES ’24*, New York, NY, USA. Association for Computing Machinery.
- Brasil (2011). Lei nº 12.527, de 18 de novembro de 2011.
- Brasil (2018). Lei nº 13.709, de 14 de agosto de 2018. *Presidência da República, Secretaria-Geral Subchefia para Assuntos Jurídicos*.
- Brito, F. and Machado, J. (2017). *Preservação de Privacidade de Dados: Fundamentos, Técnicas e Aplicações*. Sociedade Brasileira de Computação.
- Carvalho, T., Moniz, N., Faria, P., and Antunes, L. (2023). Survey on privacy-preserving techniques for microdata publication. *ACM Comput. Surv.*, 55(14s).
- Coelho, K., Okuyama, M., Nogueira, M., Vieira, A., Silva, E., and Nacif, J. (2025). Metodologia para avaliação da anonimização baseada em k-anonimato nos modelos de aprendizado de máquina. In *Anais do XLIII Simpósio Brasileiro de Redes de Computadores e Sistemas Distribuídos*, pages 742–755, Porto Alegre, RS, Brasil. SBC.
- De Capitani Di Vimercati, S., Foresti, S., Livraga, G., and Samarati, P. (2012). DATA PRIVACY: DEFINITIONS AND TECHNIQUES. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, 20(06):793–817.
- Direito, D. d. C. and Barros, V. L. M. (2025). Estado digital: dados e políticas públicas no brasil.
- Dwork, C. (2006). Differential privacy. In Bugliesi, M., Preneel, B., Sassone, V., and Wegener, I., editors, *Automata, Languages and Programming*, pages 1–12, Berlin, Heidelberg. Springer Berlin Heidelberg.

- Eynden, V. v. d. (2011). *Managing and sharing data: best practice for researchers*. UK Data Archive, Colchester, 3rd ed., fully rev edition. OCLC: 731028890.
- Fung, B. C., Wang, K., Fu, A. W.-C., and Yu, P. S. (2010). *Introduction to Privacy-Preserving Data Publishing: Concepts and Techniques*. Chapman & Hall/CRC, 1st edition.
- Ge, Y.-F., Wang, H., Cao, J., Zhang, Y., and Jiang, X. (2024). Privacy-preserving data publishing: an information-driven distributed genetic algorithm. *World Wide Web*, 27(1):1.
- Karagiannis, S., Ntantogian, C., Magkos, E., Tsohou, A., and Ribeiro, L. L. (2024). Mastering data privacy: leveraging k-anonymity for robust health data sharing. *International Journal of Information Security*, 23(3):2189–2201.
- Kim, I. and Kim, T. (2024). Kmbig: Safeguarding data sharing with advanced anonymization and risk management. In *2024 IEEE International Conference on Big Data and Smart Computing (BigComp)*, pages 439–444.
- Machado, J., Duarte Neto, E., and Bento Filho, M. (2019). Técnicas de Privacidade de Dados de Localização. In Cavalcanti, M. C. and Traina, A., editors, *Tópicos em Gerenciamento de Dados e Informações: Minicursos do SBBD 2019*, pages 8–37. SBC, 1 edition.
- Machado, J. C. and Neto, E. R. D. (2021). Privacidade de dados de localização: Modelos, técnicas e mecanismos. *Sociedade Brasileira de Computação*.
- Machanavajjhala, A., Kifer, D., Gehrke, J., and Venkitasubramaniam, M. (2007). *L* - diversity: Privacy beyond *k* -anonymity. *ACM Transactions on Knowledge Discovery from Data*, 1(1):3.
- Queiroz, M. and Motta, G. (2015). Privacidade e transparência no setor público: Um estudo de caso da publicação de microdados do inep. In *Anais do XV Simpósio Brasileiro de Segurança da Informação e de Sistemas Computacionais*, pages 362–365, Porto Alegre, RS, Brasil. SBC.
- Quinlan, R. (1993). *C4.5: Programs for Machine Learning*. Morgan Kaufmann Publishers, San Mateo, CA.
- Samarati, P. (2001). Protecting respondents’ identities in microdata release. *IEEE Transactions on Knowledge and Data Engineering (TKDE)*, 13(6):1010–1027.
- Sweeney, L. (2002). k-ANONYMITY: A MODEL FOR PROTECTING PRIVACY. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, 10(05):557–570.
- Tejedo-Romero, F., Ferraz Esteves Araujo, J. F., and Gonçalves Ribeiro, M. J. (2025). The usability of brazilian government open data portals: ensuring data quality. *Humanities and Social Sciences Communications*, 12(1):297.
- Wang, X., Chou, J.-K., Chen, W., Guan, H., Chen, W., Lao, T., and Ma, K.-L. (2018). A utility-aware visual approach for anonymizing multi-attribute tabular data. *IEEE Transactions on Visualization and Computer Graphics*, 24(1):351–360.