



**ESTUDO TÉCNICO SOBRE ANONIMIZAÇÃO DE DADOS NA  
LGPD: UMA VISÃO DE PROCESSO BASEADO EM RISCO E  
TÉCNICAS COMPUTACIONAIS**

**VERSÃO 1.0**

**BRASÍLIA/DF  
NOVEMBRO DE 2023**

## **Autoridade Nacional de Proteção de Dados**

### **Diretor-Presidente**

**Waldemar Gonçalves Ortunho Júnior**

### **Diretores**

**Arthur Pereira Sabbat**

**Joacil Basílio Rael**

**Miriam Wimmer**

### **Equipe de Elaboração**

**Marcelo Santiago Guedes – Coordenador-Geral de Tecnologia e Pesquisa (CGTP)**

**Diego Carvalho Machado – Especialista (CGTP)**

**Albert França Josué Costa – Especialista (CGTP)**

## SUMÁRIO

<b>1. INTRODUÇÃO .....</b>	<b>4</b>
<b>2. O PROCESSO DE ANONIMIZAÇÃO DE DADOS .....</b>	<b>5</b>
2.1. Utilidade do dado pessoal derivada da finalidade da operação de tratamento .....	5
2.2. Documentação do processo de anonimização .....	7
2.3. Gestão do risco de reidentificação .....	7
2.4. Limitações das técnicas de anonimização .....	12
<b>REFERÊNCIAS .....</b>	<b>18</b>
<b>APÊNDICES .....</b>	<b>20</b>
I. CADERNO DE TÉCNICAS DE ANONIMIZAÇÃO DE DADOS .....	20
II. GLOSSÁRIO .....	27

## 1. INTRODUÇÃO

Com o crescimento exponencial da quantidade de dados gerada e compartilhada diariamente, é essencial garantir que os dados pessoais dos titulares sejam devidamente protegidos contra o uso indevido ou não autorizado. A anonimização de dados é uma das possíveis formas que podem conduzir ao alcance dessa garantia. De forma geral, a anonimização consiste em um processo pelo qual os dados com capacidade de identificar um titular são transformados de maneira que a probabilidade de os associar, diretamente ou indiretamente, a um titular específico é reduzida.

Compreender a anonimização de dados como um processo a partir de uma abordagem baseada em risco, além de fundamentada na Lei Geral de Proteção de Dados Pessoais (LGPD), é benéfico para o agente responsável pelo tratamento de dados, pois oferece um conjunto mínimo de etapas que podem ser seguidas. Além disso, essa abordagem possui a vantagem de permitir a adaptação às características específicas de cada organização.

Uma das etapas do processo de anonimização está relacionada à escolha e à aplicação do conjunto de técnicas de anonimização. Embora seja comum na literatura encontrar exemplos de técnicas voltadas para dados textuais estruturados, o tratamento de dados atualmente abrange uma ampla gama de formatos, incluindo imagens, áudio, dados não estruturados e dados em fluxo. Portanto, é essencial conhecer, desenvolver e aprimorar técnicas de anonimização que sejam aplicáveis a diversas formas de dados, garantindo a privacidade e a proteção dos dados de forma mais ampla.

O §3º do art. 12 da Lei Geral de Proteção de Dados Pessoais (BRASIL, 2018) estabelece que a Autoridade Nacional de Proteção de Dados (ANPD) poderá dispor sobre padrões e técnicas utilizadas em processos de anonimização. Diante disso, esse documento colabora para o posicionamento desta Autoridade no tema.

Por esses motivos, a Seção 2 apresenta uma proposta de processo genérico para a anonimização, juntamente com uma breve discussão sobre as limitações de algumas técnicas empregadas nesse contexto; essas limitações destacam a importância crucial da escolha correta do conjunto de técnicas de anonimização. Dentro do âmbito das técnicas de anonimização, este documento apresenta no Apêndice I um caderno de técnicas, fornecendo explicações sobre cada uma, juntamente com suas representações em pseudocódigo. Além disso, é destacado para qual formato de dado cada técnica é mais

adequada. Adicionalmente, são apresentadas suas aplicações e limitações, permitindo aos responsáveis pela anonimização uma análise criteriosa para selecionar a melhor abordagem de acordo com os requisitos específicos e considerações de segurança e privacidade aplicáveis.

O presente estudo é uma introdução computacional ao processo de anonimização e compõe a série de estudos técnicos a respeito da anonimização de dados na LGPD. O primeiro documento, “*Estudo técnico sobre a anonimização de dados na LGPD – análise jurídica*”, abordou a anonimização de dados de acordo com sua dimensão jurídico-regulatória e analisou seus fundamentos normativos à luz da LGPD. As principais conclusões apontadas nessa análise jurídica orientam o presente estudo, entre as quais destacam-se: (i) o ato inicial do processo de anonimização de dados configura operação de tratamento de dado pessoal, atraindo, assim, o regime da LGPD; (ii) os conceitos de dado pessoal e de dado anonimizado possuem caráter dinâmico e contextual; (iii) a avaliação do processo de anonimização deve se dar de acordo com um modelo baseado em risco.

## **2. O PROCESSO DE ANONIMIZAÇÃO DE DADOS**

O processo de anonimização, orientado por uma abordagem baseada em risco, tem como objetivo fornecer um conjunto mínimo de etapas que podem servir de guia de boas práticas aos agentes de tratamento de dados. Essas etapas sugerem que o agente identifique e compreenda os riscos envolvidos em sua atividade, bem como adote medidas para mitigá-los.

A seção inicia apresentando o conflito entre a utilidade e o grau de anonimização do dado pessoal. Em seguida, é debatida a importância de se documentar o processo de anonimização para garantir a qualidade dos dados. A subseção 2.3, por sua vez, aborda os riscos e a gestão dos riscos de reidentificação. Por fim, os limites das técnicas clássicas<sup>1</sup> de anonimização para lidarem com dados em outros formatos além do textual estruturado.

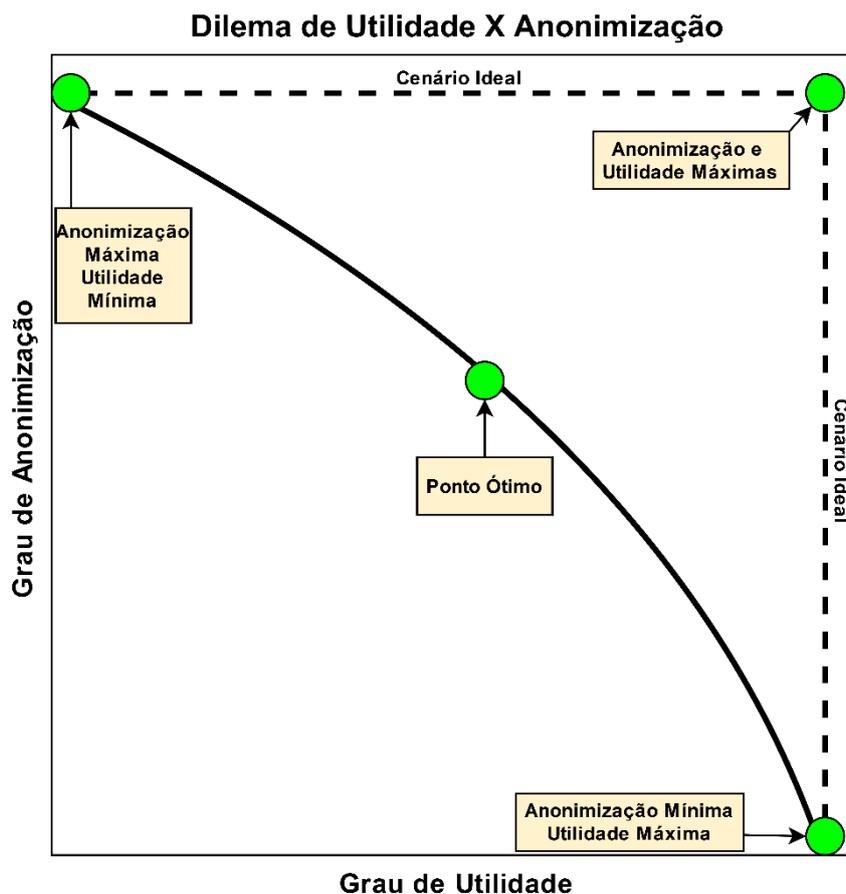
### **2.1. Utilidade do dado pessoal derivada da finalidade da operação de tratamento**

---

<sup>1</sup> Técnicas clássicas são as técnicas que lidam com dados textuais estruturados não em fluxo.

A LGPD (BRASIL, 2018) determina que os dados pessoais devem ser tratados para propósitos legítimos, específicos, explícitos e informados ao titular. Partindo desse enunciado, é possível observar que a atividade de tratamento de dados pessoais precisa estar atrelada a uma finalidade específica, de tal forma compete ao agente de tratamento identificar o grau de utilidade do dado pessoal para alcançar a finalidade especificada, em consequência estabelecer o grau necessário de anonimização dos dados (Figura 1).

Figura 1: Dilema Utilidade x Anonimização.



Em termos teóricos, existe um ponto ótimo em que o grau de utilidade do dado pessoal e o grau de anonimização são simultaneamente máximos. Entretanto, em termos práticos esse ponto ótimo não é fácil de ser alcançado, pois depende de um ajuste fino entre duas variáveis conflitantes. De tal forma, a abordagem da anonimização como um processo contínuo baseado em risco possibilita que o controlador defina, de acordo com seu contexto, o compromisso entre ao grau de utilidade e o grau de anonimização que contemple a finalidade definida no tratamento e minimize o risco de reidentificação do titular.

## 2.2. Documentação do processo de anonimização

Dentro do contexto do processo de anonimização, a documentação das atividades realizadas permite aos agentes de tratamento, em especial ao controlador, terem uma trilha dos estados que o conjunto de dados assumiu durante o processo. Em especial atenção, observa-se que o conjunto de dados resultante deve manter as propriedades estatísticas da base em sua forma original. Caso contrário a qualidade dos dados poderá ser degradada, diminuindo ou até mesmo impossibilitando alcançar a finalidade pretendida.

Considerando que, majoritariamente, as técnicas de anonimização adicionam ruído aos dados, essa adição pode acarretar mudanças nas propriedades estatísticas do conjunto de dados. A guarda dos registros das propriedades estatísticas do conjunto de dados original é um elemento útil na avaliação do dilema exposto na Figura 1 e consequentemente é importante para maximizar a qualidade do conjunto de dados.

## 2.3. Gestão do risco de reidentificação

A anonimização não pode ser entendida como um processo discreto<sup>2</sup> com dois estados distintos, em que o primeiro representaria o cenário em que todos os dados são tratados em sua forma original e o segundo representaria o cenário da anonimização plena, e que após os dados terem transitado do primeiro estado para o segundo, não seria mais possível haver uma transição de retorno para o cenário inicial.

Narayanan e Shmatikov (2010) discorrem sobre os mitos e as falácias associados à anonimização de dados ao afirmarem que não há técnica com eficácia plena, estando todas elas sujeitas a ataques de reidentificação, isto é, riscos de reidentificação. Por esse motivo, a anonimização deve ser entendida como um processo contínuo baseado em riscos, composto por muitas fases e com transições multidirecionais entre essas, havendo inclusive a possibilidade de o processo de anonimização ser revertido por ação de agentes terceiros ou de agentes internos, não sendo factível considerar que o processo de anonimização pode resultar em um cenário de risco zero.

---

<sup>2</sup> A noção de processo discreto está relacionada aos estados serem bem definidos e sem intersecção entre eles, ou seja, significa que entre dois estados do processo não há entre estados.

Há dois cenários conhecidos na literatura para a reidentificação de um indivíduo em um conjunto de dados. O primeiro é conhecido como reidentificação do promotor<sup>3</sup>, que parte da hipótese de que o atacante conhece um indivíduo em particular do conjunto de dados e deseja encontrar o registro relacionado a ele. O segundo é conhecido por reidentificação do jornalista<sup>4</sup>, que parte da hipótese de que o atacante não conhece um indivíduo em particular do conjunto de dados e deseja somente conseguir reidentificar qualquer indivíduo (EL-EMAM, 2013; *National Institute of Standards and Technology*, 2015).

O risco associado à reversibilidade do processo de anonimização é um elemento chave que deve ser considerado pelo agente de tratamento ao realizar a anonimização dos dados pessoais. De tal forma, o processo de anonimização tem como objetivo minimizar os riscos de reidentificação dos dados pessoais mantendo a utilidade deles (AEPD, 2016).

Para isso, a gestão do risco de reidentificação deve ser realizada de forma contínua durante todo o tratamento de dados realizado, para que o controlador tenha evidências suficientes para a tomada de decisão relacionada à proteção de dados e à privacidade dos titulares de dados. Por mais que o cenário de anonimização apresente características heterogêneas, quando se analisa sob o ponto de vista de processo, algumas etapas são essenciais e estão apresentadas na Figura 2.

A primeira etapa consiste na determinação do Risco de Reidentificação Aceitável (RRA) para um determinado conjunto de dados, e tem como objetivo estipular um limite superior para o risco. Um risco de reidentificação superior ao limite estabelecido descaracterizará o conjunto de dados como anonimizado. Essa primeira etapa é de extrema importância e possui uma gama de variáveis dependentes do contexto<sup>5</sup> que devem ser observadas pelo agente de tratamento. Desse modo, não é possível estabelecer uma metodologia padronizada a todos os casos. Pode-se citar como exemplos de variáveis de contexto a existência de dados pessoais sensíveis ou dados financeiros que podem diminuir o limite do risco aceitável.

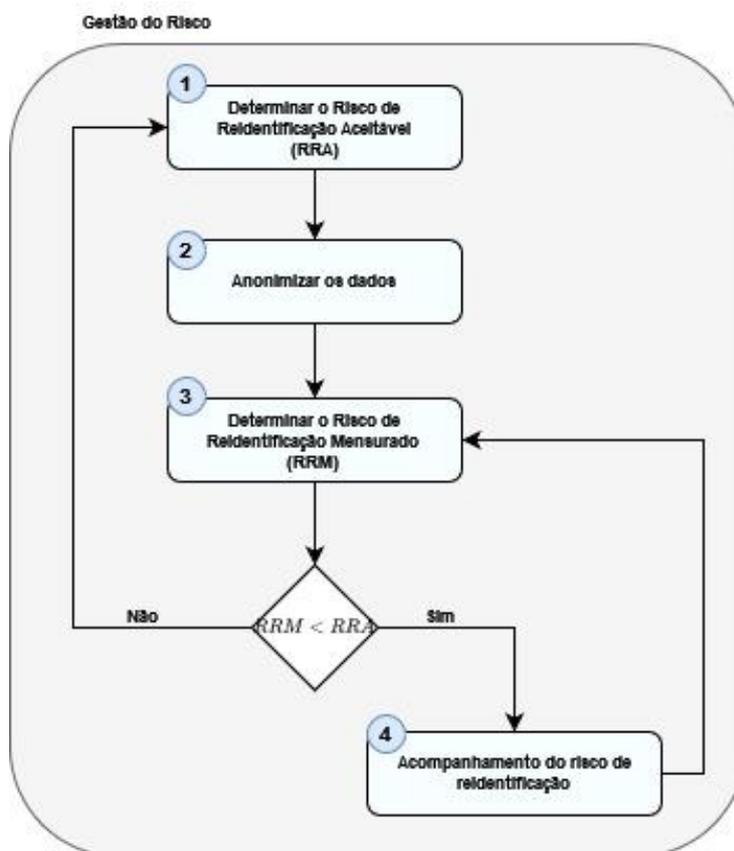
---

<sup>3</sup> O cenário do promotor, do inglês *prosecutor scenario*, refere-se à situação em que o risco de reidentificação de um titular em específico é aumentado pois o atacante sabe que os dados desse titular estão no conjunto de dados.

<sup>4</sup> O cenário do jornalista, do inglês *journalist scenario*, refere-se à situação em que o risco de reidentificação de todos os titulares são iguais, o atacante não possui nenhuma informação *a priori* de quais titulares estão presentes no conjunto de dados.

<sup>5</sup> Variável dependente do contexto é definida como sendo uma característica interna do agente de tratamento.

Figura 2: Etapas essenciais para a gestão do risco de reidentificação.



Fonte: Própria.

A segunda etapa consiste na aplicação do conjunto de técnicas de anonimização escolhido. O objetivo dessa etapa é produzir um conjunto de dados anonimizado que tenha um risco de reidentificação não superior ao limite do risco aceitável definido na etapa anterior.

A terceira etapa consiste em determinar o Risco de Reidentificação Mensurado (RRM) de um ataque de reidentificação ter sucesso no conjunto de dados - o risco de reidentificação deve assumir a forma de probabilidade<sup>6</sup>. De modo semelhante à primeira etapa, diversas variáveis dependentes do contexto podem ser observadas pelo controlador; como exemplo, tem-se a condição do conjunto de dados ser público, compartilhado ou privado.

<sup>6</sup> Uma probabilidade é um valor real no intervalo fechado de 0 a 1.

Adicionalmente, na terceira etapa, métricas contextuais<sup>7</sup> também devem ser observadas na mensuração do risco. El-Eman (2013) e IPCO (2016) afirmam que há diversas métricas contextuais, porém elas majoritariamente derivam de uma métrica base<sup>8</sup> que utiliza o conceito de equivalência de classe para determinar o risco de reidentificação.

A equivalência de classe é um conceito da teoria dos conjuntos que pode ser definido como sendo um subconjunto  $S$  que inclui todos os elementos equivalentes a um dado elemento. A relação da equivalência de classe é reflexiva, simétrica e transitiva (AVELSGAARD, 1989). Dentro do processo de anonimização, a equivalência de classe pode ser traduzida como sendo o subconjunto formado pelos titulares que possuem dados, no conjunto original, que compartilham os mesmos valores.

Na Tabela 1, um exemplo de equivalência de classe em um conjunto de dados textuais tabular ( $A$ ) é apresentado. Observa-se que, ao considerar a coluna Estado de Nascimento, é possível formar 3 subconjuntos do conjunto inicial de acordo com a repetição dos valores da coluna.

Tabela 1: Exemplo de equivalência de classe em um conjunto de dados tabular.

Nome Completo	Estado de Nascimento	Data de Nascimento
FM	Acre	16/12/1944
RJ	Rio de Janeiro	27/03/1960
LB	Rio Grande do Norte	30/12/1932
EEC	Acre	05/11/1938

Fonte: Própria.

Tabela 2: Exemplos de subconjuntos gerados por equivalência de classe.

$A_{EN=ACRE}=2$

Nome Completo	Estado de Nascimento	Data de Nascimento
FM	Acre	16/12/1944
EEC	Acre	05/11/1938

$A_{EN=Rio\ de\ Janeiro}=1$

Nome Completo	Estado de Nascimento	Data de Nascimento
RJ	Rio de Janeiro	27/03/1960

$A_{EN=Rio\ Grande\ do\ Norte}=1$

Nome Completo	Estado de Nascimento	Data de Nascimento
LB	Rio Grande do Norte	30/12/1932

<sup>7</sup> Métrica contextual é definida como sendo uma métrica derivada de uma métrica base, com a incorporação de elementos particulares. Por exemplo: K-Anonimização (SAMARATI e SWEENEY, 1998), T-Proximidade (LI, LI e VENKATASUBRAMANIAN, 2007) e L-Diversidade (AGGARWAL e YU, 2008), que são derivadas da métrica base conhecida por Equivalência de Classe.

<sup>8</sup> Métrica base é um valor definido para mensurar o risco de reidentificação calculado unicamente com base no próprio conjunto de dados. Por exemplo, a Equivalência de Classe.

Fonte: Própria.

Conforme exposto na Tabela 2, o primeiro subconjunto gerado é composto por todos os elementos em que a coluna Estado de Nascimento tem o valor Acre; ele possui a equivalência de classe igual a 2. De forma semelhante, o segundo e o terceiro subconjunto são compostos pelos elementos que possuem, respectivamente, os valores Rio de Janeiro e Rio Grande do Norte para a coluna Estado de Nascimento. Os dois últimos subconjuntos do exemplo possuem a equivalência de classe igual a 1.

É importante observar que, de forma geral, quanto menor o valor da equivalência de classe, maior é o grau de unicidade do dado, por conseguinte, maior o risco de reidentificação. No exemplo, os titulares de dados nascidos no Rio de Janeiro ou no Rio Grande do Norte têm risco de reidentificação igual a 100%, por sua vez os nascidos no Acre têm o risco de reidentificação em 50%. E o risco médio de reidentificação do conjunto de exemplo é de aproximadamente 83,33%, ao considerar a Equivalência de Classe como métrica.

Conforme mencionado, as métricas contextuais de mensuração do risco de reidentificação, majoritariamente, derivam da métrica de equivalência de classe, de tal forma, a métrica contextual pode ser computada para cada um dos elementos do conjunto de dados, e os valores resultantes podem ser ponderados para determinar o valor geral da métrica contextual para o conjunto de dados. O valor resultante pode ser então ponderado pelas variáveis contextuais.

$$\text{Risco de Reidentificação Mensurado} = V_C * \theta, \quad (2)$$

$\theta$  representa o valor geral da métrica contextual e  $V_C$  representa um fator de ponderação das variáveis contextuais, quando existentes, caso não existam  $V_C$  pode assumir o valor de 1.

No contexto de base de dados públicas ou compartilhadas, o risco deve ser majorado e, conseqüentemente, o valor de  $V_C$  deve ser definido de forma adequada a representar a majoração do risco.

O risco mensurado (RRM) deve ser comparado ao limite de risco aceitável (RRA). Caso o RRM seja maior do que o RRA, o conjunto de dados perde a condição de anonimizado, sendo necessário o reinício do processo de anonimização. Caso contrário, é necessário acompanhar o uso do conjunto de dados, especialmente quando operações

realizadas sobre ele possam modificar o risco mensurado, tais como operações de inclusão, alteração ou deleção de dados; havendo essas operações é necessário atualizar o nível de risco mensurado.

#### 2.4. Limitações das técnicas de anonimização

É vasta a literatura que apresenta lista de técnicas de anonimização e casos de aplicação dessas. Geralmente, essa lista contém técnicas de anonimização que são limitadas a dados textuais estruturados e não em fluxo, por exemplo, as técnicas de supressão, generalização, perturbação e permutação. Essas limitações são muitas vezes indicadas nos próprios documentos, a exemplo das publicações da Autoridade de Ontário (IPCO; 2016) e da Autoridade de Singapura (PDPC, 2023).

Entretanto, dados textuais estruturados representam uma parcela não majoritária do formato em que dados podem ser armazenados. Ainda dentro do formato textual estruturado, o avanço da computação ubíqua trouxe o cenário de dados em fluxos. Nesse cenário, por mais que seja possível a aplicação dessas técnicas de anonimização, o custo computacional envolvido torna sua aplicação impraticável.

O cenário de dados em fluxo pode ser definido como aquele em que novos dados são gerados continuamente a uma alta taxa de velocidade, com tamanho potencialmente infinito e necessidade de processamento imediato (YANXIA *et al.*, 2016).

Tabela 3: Exemplo de conjunto de dados tabular textual.

Nome Completo	Cidade de Nascimento	Data de Nascimento
FM	Rio Branco	16/12/1944
AFB	Macaíba	27/03/1960
LB	Natal	30/12/1932
MTL	Xapuri	05/11/1938
CGG	Macaé	23/03/1989
RJ	Rio de Janeiro	28/02/1957

Fonte: Própria.

Para exemplificar a limitação existente na aplicação das técnicas de anonimização em dados não textuais, inicialmente considere o conjunto de dados textual tabular exposto na Tabela 3. Analisando as colunas, é perceptível a adequação da técnica de anonimização de generalização na coluna Cidade de Nascimento, substituindo as cidades pelas Unidades da Federação correspondentes. O conjunto de dados resultante está exposto na Tabela 4.

O conjunto de dados resultante tem como média da métrica base de equivalência de classe igual a 2. Com um risco médio de reidentificação de 50%, já que cada elemento compartilha o mesmo valor da coluna Unidade da Federação de Nascimento com outro elemento.

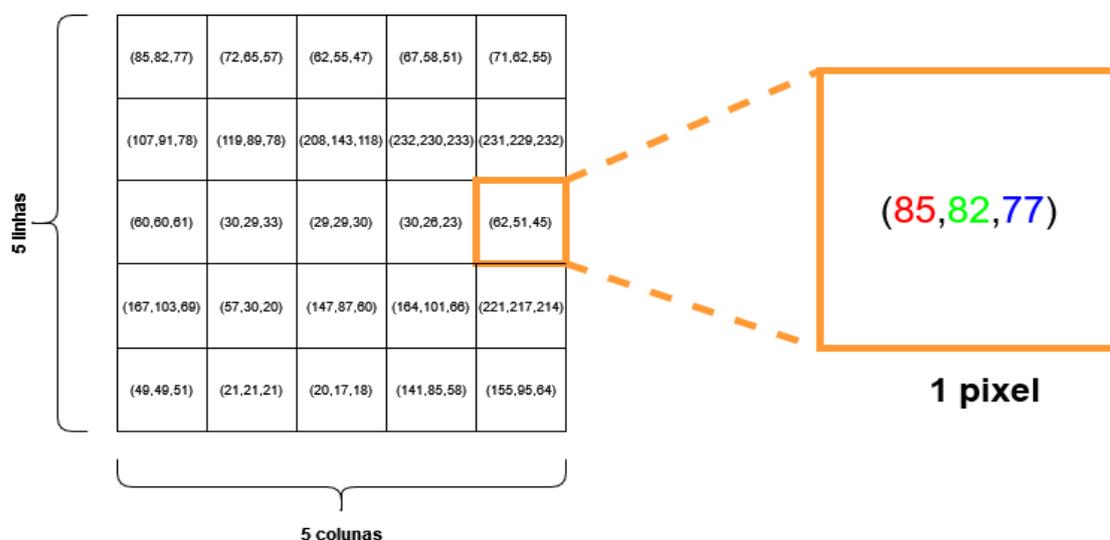
Tabela 4: Conjunto de dados da Tabela 3 após aplicação da técnica de generalização no atributo Cidade de Nascimento.

Nome Completo	Unidade da Federação de Nascimento	Data de Nascimento
FM	Acre	16/12/1944
AFB	Rio Grande do Norte	27/03/1960
LB	Rio Grande do Norte	30/12/1932
MTL	Acre	05/11/1938
CGG	Rio de Janeiro	23/03/1989
RJ	Rio de Janeiro	28/02/1957

Fonte: Própria.

Para avançar sobre o exemplo, é necessária uma breve explicação de como uma imagem colorida é composta. Em uma imagem colorida os dados são estruturados em um tensor, sendo cada posição do tensor representando um *pixel* da imagem. Particularmente para as imagens coloridas, cada *pixel* é composto por outros três valores, sendo um valor para cada canal de cor (vermelha, verde e azul). A Figura 3 apresenta uma representação simplificada de uma imagem de tamanho 5x5, totalizando 25 *pixels*. Em destaque na figura há a representação dos três valores para cada canal de cor contido em um *pixel*.

Figura 3 Representação tensorial de uma imagem.



Fonte: Própria.

Figura 4: Conjunto de dados de imagens (geradas artificialmente).

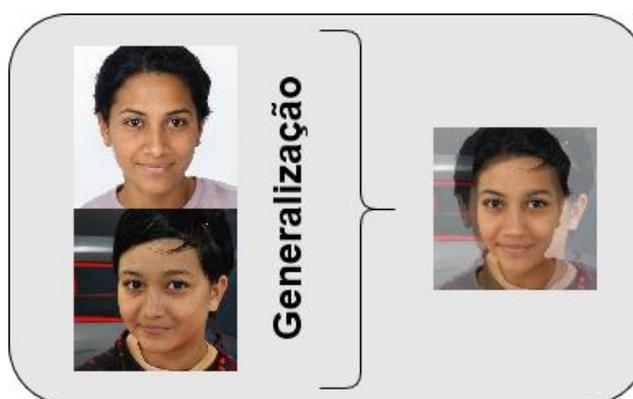


Fonte: Unreal Person (2023).

De forma análoga à aplicação da generalização nos dados textuais na Tabela 4, ao considerar um conjunto de dados de imagens (Figura 4), observa-se que a aplicação da generalização é inviável, não havendo como estabelecer critérios para o agrupamento dos elementos do conjunto. Na hipótese em que algum critério de agrupamento seja de alguma forma estabelecido, por exemplo, a cor do cabelo, a inviabilidade recai na técnica em si, pois o agrupamento de duas imagens de modo que não seja possível identificar as pessoas nelas contidas apresenta pouco sentido.

Uma tentativa de forçar a aplicação da técnica de generalização nesse cenário é considerar a estrutura de uma imagem (Figura 3) e para cada canal de cor de cada um dos *pixels* das duas imagens aplicar alguma medida de centro sobre o par de valores, por exemplo, a média. O valor da medida de centro então é utilizado para formar a imagem resultante da generalização. Entretanto é possível observar no exemplo da Figura 5 que a imagem resultante tem utilidade nula e não garante a privacidade do titular.

Figura 5: Imagem resultante da técnica de generalização (média simples).



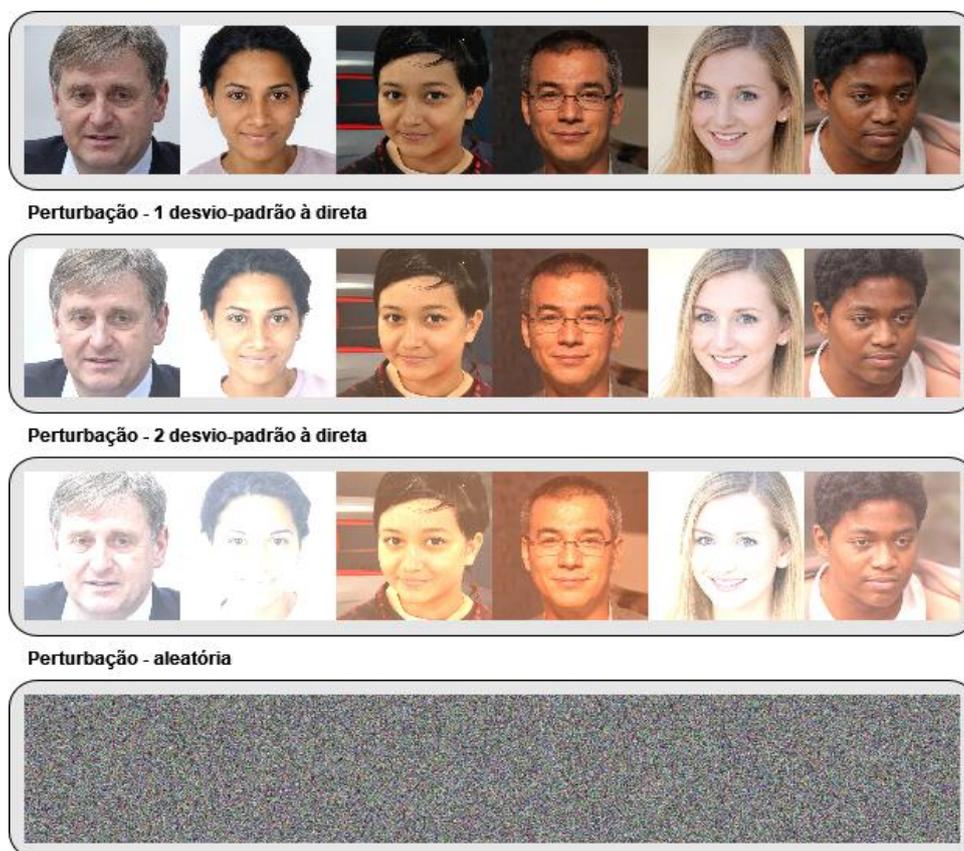
Fonte: Própria.

Expandindo o exemplo, a aplicação da técnica anonimização denominada de perturbação em imagens tem também importantes limitações. Na Figura 6 expõe-se o resultado da aplicação da perturbação em três níveis distintos. Observa-se que nos dois primeiros há o deslocamento dos valores dos *pixels* à direita da média, o que resulta em uma imagem mais clara, porém ainda sendo possível perfeitamente identificar a pessoa

presente na imagem. Por sua vez, o último nível apresenta o resultado da aplicação da perturbação completamente aleatória, o que resulta em uma máxima anonimização, porém anula a utilidade do dado.

De tal forma, o exemplo exposto solidifica o argumento da necessidade de primeiramente abordar a anonimização na forma de processo e não somente uma abordagem restrita à listagem de técnicas de anonimização com indicação de casos de usos, como é comumente encontrado. Adicionalmente, é de extrema importância considerar a ampla gama de contextos distintos de tratamento de dados pessoais que não permite que a discussão ignore os aspectos contextuais que ultrapassam o escopo deste documento.

Figura 6: Aplicação da técnica de perturbação em imagens.



Fonte: Própria.

Todt, Hanisch e Strufe (2022) apresentam uma lista exemplificativa de técnicas de anonimização específicas para imagens, sendo essas denominadas de máscara de olhos, permutação por blocos, ruído gaussiano, desfoque gaussiano, pixelação, DP Pix, Fawkes, CIAGAN e K-RTIO. Outro importante ponto que é apresentado pelos mesmos autores é que as técnicas de anonimização de imagens estão sujeitas a outros tipos de

ataques de reidentificação, a exemplo do ataque adversarial baseados em aprendizado de máquina.

Outro exemplo de anonimização de imagens é discutido em Nature (2022). Nesse trabalho é desenvolvida uma técnica para minimizar os riscos de divulgação inadequada de imagens faciais de pacientes, por meio da aplicação de uma máscara digital para apagar características identificáveis mantendo características relevantes para a doença necessária para o diagnóstico.

Certamente, outros tipos de dados, além de dados textuais e de imagens, possuem características particulares que necessitam serem consideradas quando da aplicação de técnicas de anonimização.

### **3. CONSIDERAÇÕES FINAIS**

O presente estudo técnico foi elaborado com o objetivo de apresentar a anonimização na Lei Geral de Proteção de Dados Pessoais do ponto de vista da ciência da computação. O documento traz como importante orientação aos agentes de tratamento abordar por padrão a anonimização como um processo contínuo baseado em riscos, e não somente limitar-se à aplicação de técnicas, pois essas sempre estão sujeitas a ataques de reidentificação bem-sucedidos.

Do ponto de vista restrito às técnicas, observa-se que muitas das técnicas abordadas na literatura têm suas aplicações extremamente limitadas ao cenário de dados textuais estruturados não em fluxo. De tal forma, limitar a presente análise às técnicas para anonimizar dados textuais estruturados seria de pouca valia aos agentes de tratamento, pois no cenário atual os dados assumem diversos formatos, tais como imagem e áudio.

Para contemplar essa análise mais ampla das técnicas de anonimização, no Apêndice I é apresentado um caderno contendo a análise de algumas técnicas de anonimização selecionadas. As técnicas apresentadas são meramente exemplificativas e a relação está longe de exaurir as técnicas existentes. Uma ressalva importante é que a análise não tem como objetivo indicar técnicas aos agentes de tratamento, limitando-se unicamente a apresentar os detalhes técnicos, estando o agente de tratamento responsável por decidir quais técnicas devem ser utilizadas no caso concreto.

De forma semelhante à parte da análise jurídica desse estudo técnico, o presente documento não objetiva esgotar o tema da anonimização no contexto da ciência da

computação. Ao contrário, lança as bases para a expansão das orientações da ANPD a todo o ecossistema brasileiro de proteção de dados.

## REFERÊNCIAS

AEPD. **Orientaciones y garantías em los procedimientos de anonimización de datos personales**. Agencia Española de Protección de Datos. 2016.

AGGARWAL, Charu C.; YU, Philip S. **A general Survey of Privacy-Preserving Data Mining Models and Algorithms**. Privacy-Preserving Data Mining. Advances in Database System. vol 34. Springer. 2008.

ARBUCKLE, L.; EL-EMAM, K. **Building an Anonymization Pipeline**. O`Reilly Media, Inc., 2020.

AVELSGAARD, Carol. **Foundations for Advanced Mathematics**. Scott, Foresman, 1989.

BRASIL. **Lei nº 13.709, de 14 de agosto de 2018**. Lei Geral de Proteção de Dados Pessoais (LGPD). Disponível em: [https://www.planalto.gov.br/ccivil\\_03/\\_ato2015-2018/2018/lei/113709.htm](https://www.planalto.gov.br/ccivil_03/_ato2015-2018/2018/lei/113709.htm). Acesso em: 09 mai. 2023.

EL-EMAM, Khaled. **Guide to the de-identification of personal health information**. Nova York: CRC Press, 2013.

EL-EMAM, K.; ARBUCKLE, L. **Anonymizing Health Data**, O`Reilly Media, Inc., 2013.

IPCO. **De-identification Guidelines for structured data**. Information and Privacy Commissioner of Ontário, 2016.

PDPC. **Guide to basic anonymisation**. Singapore: SG Digital-PDPC, 2022.

LI Ninghui; LI, Tiancheng; VENKATASUBRAMANIAN Suresh. **T-Closeness: Privacy Beyond k-Anonymity and L-Diversity**. IEEE 23rd Internation Conference on Data Engineering, 2007.

NATURE. **Anonymizing Facial Images to Improve Patient Privacy**. *In* Nature Medicine, 2022.

NARAYANAN, A.; SHMATIKOV, V. **Myths and fallacies of “personally identifiable information**. *In* Communications of the ACM, v. 53, n. 6, p. 24, 1 jun. 2010.

NATIONAL INSTITUTE OF STANDARDS AND TECHNOLOGY. **NISTIR 8053: De-Identification of Personal Information**. 2015.

SAMARATI, P.; SWEENEY, L. **Protecting privacy when disclosing information: k-anonymity and tis enforcement through generalization and suppression**. Technical Report, 1998.

TODT, J., HANISCH, S., STRUFE, T. **Fantômas: Evaluating Reversibility of Face Anonymizations Using a General Deep Learning Attacker**, 2022.

UNREAL PERSON, **This Person does not exist**. Disponível em: <https://www.unrealperson.com/>. Acesso em: 16 de maio de 2023.

YANXIA, L., CUIRONG, W., CONG. W., BINGYUI, L. **Uncertain Data Stream Classification with Concept Drift**. *In*: Conference International on Advanced Cloud and Big Data (CDB), 2016.

## APÊNDICES

### I. CADERNO DE TÉCNICAS DE ANONIMIZAÇÃO DE DADOS

Esse caderno tem como objetivo apresentar uma análise de técnicas de anonimização selecionadas. A análise compreende uma descrição da técnica, exemplo, a representação em pseudocódigo, vantagens e limites. A relação de técnicas apresentadas nesse caderno é meramente exemplificativa e está longe de exaurir todas as técnicas existentes.

É importante destacar que esse documento não tem como objetivo indicar quais técnicas são melhores ou piores do que outras. O caderno objetiva unicamente apresentar os detalhes das técnicas, cabendo ao agente de tratamento identificar as mais adequadas para incorporar ao seu processo de anonimização.

#### I.A. TÉCNICAS PARA ANONIMIZAR DADOS TEXTUAIS ESTRUTURADOS

Técnica de Generalização		
<b>Descrição</b>		
A técnica agrupa os dados com características em comum em um nível de granularidade maior. Os valores dos atributos são substituídos pelos valores do grupo.		
<b>Exemplo</b>		
Dado original		
Nome Completo	Cidade de Nascimento	Idade
FM	Rio Branco	79
AFB	Macaíba	63
LB	Natal	91
MTL	Xapuri	85
CGG	Macaé	34
RJ	Rio de Janeiro	66
Dado anonimizado por meio da generalização		
Nome Completo	Estado de Nascimento	Faixa Etária
FM	Acre	70-79
AFB	Rio Grande do Norte	60-69
LB	Rio Grande do Norte	90-99
MTL	Acre	80-89
CGG	Rio de Janeiro	30-39
RJ	Rio de Janeiro	60-69
<b>Aplicação</b>		
Apropriada quando os dados possuem características em comum que permitem sua representação de forma generalizada, sem perda da utilidade.		
<b>Limites</b>		
<ul style="list-style-type: none"> <li>• Aplicável somente a dados textuais estruturados.</li> <li>• Aplicável somente em dados que compartilham características em comum.</li> </ul>		

- Não aplicável quando houver dados com valores únicos.
- Perda da precisão dos dados.
- Custo computacional elevado para aplicação em dados textos estruturados em fluxo

### Pseudocódigo

**Algoritmo:** Pseudocódigo da técnica de generalização

**Entrada:** *dt*: Dados tabular textual; *vGen*: Lista de valores generalizados < *valorColuna* : *valorAntigo* : *valorNovo* >;

**Saída:** *dtta*: Dados tabular textual anonimizados

```

1  $N \leftarrow tamanho(dt)$ 
2  $i \leftarrow 1$ 
   // Para cada instância nos dados de entrada, substitui o valor original
   // pelo valor generalizado.
3 while  $i \neq N$  do
4   forall coluna c de dt[i] do
5     forall elemento e de vGen do
6       if  $c = e.valorColuna$  then
7         if  $dt[i][c].valor = e.valorAntigo$  then
8            $dt[i][c].valor \leftarrow e.valorNovo$ 
9        $i \leftarrow i + 1$ 
10  $dtta \leftarrow dt$ 
11 return dtta

```

## Técnica de Mascaramento

### Descrição

A técnica consiste em substituir uma parte dos caracteres dos dados por um caractere símbolo (por exemplo \* ou x).

### Exemplo

Dado original

Nome Completo	CPF	Idade
FM	111.111.111-11	79
AFB	222.222.222-22	63
LB	333.333.333-33	91
MTL	444.444.444-44	85
CGG	555.555.555-66	34
RJ	666.666.666-66	66

Dado anonimizado por meio do mascaramento

Nome Completo	CPF	Idade
FM	111.***.***-11	79
AFB	222.***.***-22	63
LB	333.***.***-33	91
MTL	444.***.***-44	85
CGG	555.***.***-66	34

RJ	666.***.***-66	66
<b>Aplicação</b>		
Apropriada quando a substituição de parte dos dados por um caractere simbólico fornece o nível desejado de anonimização, sem perda da utilidade.		
<b>Limites</b>		
<ul style="list-style-type: none"> <li>• Aplicável somente a dados textuais.</li> <li>• Perda da precisão dos dados.</li> <li>• Custo computacional elevado para aplicação em dados textos estruturados em fluxo.</li> <li>• Escolha do padrão adequado que permita desvincular o titular do dado anonimizado.</li> <li>• Em casos de dados públicos ou compartilhados, como não há um padrão para o mascaramento, é possível que partes distintas dos dados estejam visíveis e por consequência os dados originais sejam reconstruídos.</li> </ul>		
<b>Pseudocódigo</b>		
<b>Algoritmo:</b> Pseudocódigo da técnica de mascaramento		
<b>Entrada:</b> <i>dt</i> : Dados tabular textual; <i>colMascara</i> : Lista de colunas com a respectiva máscara a ser aplicada. $\langle \text{valorColuna} : \text{valorMascara} \rangle$ ; <b>Saída:</b> <i>dt</i> : Dados tabular textual anonimizados		
<pre> 1 <math>N \leftarrow \text{tamanho}(dt)</math> 2 <math>i \leftarrow 1</math>   // Para cada instância nos dados de entrada, aplica a máscara no valor original. 3 while <math>i \neq N</math> do 4   forall <i>coluna c</i> de <i>dt</i>[<i>i</i>] do 5     forall <i>elemento e</i> de <i>colMascara</i> do 6       if <math>c = e.\text{valorColuna}</math> then 7         // Função para aplicar a máscara no valor 8         <math>dt[i][c].\text{valor} \leftarrow \text{funcMascaramento}(dt[i][c].\text{valor}, e.\text{valorMascara})</math> 9       <math>i \leftarrow i + 1</math> 10 <math>dt \leftarrow dt</math> 11 return <i>dt</i> </pre>		

## Técnica de Permutação

### Descrição

A técnica consiste em reorganizar os valores dos dados dentro do conjunto de dados, de tal forma que os valores originais ainda são representados, mas geralmente não mais associado ao seu titular.

### Exemplo

Dado original

Nome Completo	Profissão	Tempo de Profissão
FM	Professor	20
AFB	Vendedor	12
LB	Advogado	15

MTL	Engenheiro de Software	8
CGG	Enfermeiro	25
RJ	Médico Veterinário	12

Dado anonimizado por meio da permutação

Nome Completo	Profissão	Tempo de Profissão
FM	Advogado	15
AFB	Professor	20
LB	Enfermeiro	25
MTL	Vendedor	12
CGG	Médico Veterinário	12
RJ	Engenheiro de Software	8

### Aplicação

Apropriada somente quando a análise dos dados precisa ser feita de forma agregada, pois a técnica elimina a possibilidade de analisar os dados ao nível do titular.

### Limites

- Aplicável somente a dados textuais estruturados.
- Aplicável somente para análise agregada.
- Perda da precisão dos dados
- Custo computacional elevado para aplicação em dados textos estruturados em fluxo.

### Pseudocódigo

---

#### Algoritmo: Pseudocódigo da técnica de permutação

---

**Entrada:** *dt*: Dados tabular textual; *colunas*: Lista de colunas para permutação

**Saída:** *dt*: Dados tabular textual anonimizados

*// Para as colunas selecionadas do dados de entrada, aplicar a permutação.*

```

1 forall coluna c de dt do
2   forall elemento e de colunas do
3     if c = e then
4       // Função para aplicar a permutação nos valores da coluna c.
4       dt[[c] ← funcPermutacao(dt[[c])
5 dt ← dt
6 return dt

```

---

## Técnica de Adição de Ruído

### Descrição

A técnica consiste em realizar pequenas modificações nos dados originais adicionando ruído nos dados. Normalmente utilizada em dados numéricos.

### Exemplo

Dado original

Nome Completo	Idade	Altura	Peso	Qtd. Filhos
FM	30	1,55	53	2
AFB	36	1,60	67	2

LB	20	1,80	92	0
MTL	22	1,68	52	1
CGG	44	1,71	61	3
RJ	27	1,75	72	1

Dado anonimizado por meio da adição de ruído.

Ao valor original é somando 1 desvio-padrão do intervalo de valores.

Para as colunas Idade e Quantidade filhos o valor do ruído foi truncando.

Nome Completo	Idade	Altura	Peso	Qtd. Filhos
FM	37	1,65	71,63	2
AFB	43	1,70	85,63	2
LB	27	1,90	110,63	0
MTL	29	1,78	70,63	1
CGG	37	1,62	71,63	2
RJ	43	1,70	85,63	2

### Aplicação

Apropriada para dados numéricos em cenários em que a precisão dos dados não é essencial para o alcance da finalidade pretendida. A adição de ruído pode fazer com que o conjunto de dados perda suas propriedades estatísticas e o invalide para a finalidade pretendida.

### Limites

- Aplicável preferencialmente em dados numéricos.
- Perda da precisão dos dados, a adição de ruído pode descaracterizar os dados com a perda de sua utilidade.
- Não deve ser utilizada quando a precisão dos dados é essencial.

### Pseudocódigo

#### Algoritmo: Pseudocódigo da técnica de adição de ruído

**Entrada:** dtt: Dados tabular textual; colunas: Lista de colunas para adição de ruído

**Saída:** dtta: Dados tabular textual anonimizados

*// Para as colunas selecionadas do dados de entrada, aplicar a adição de ruído.*

```

1 forall coluna c de dtt do
2   forall elemento e de colunas do
3     if c = e then
4       // Função para aplicar o ruído nos valores da coluna c.
5       dtt[[c] ← funcRuido(dtt[[c])
6 dtta ← dtt
return dtta

```

## I.B. TÉCNICAS PARA ANONIMIZAR IMAGENS

### Técnica de Desfoque Gaussiano (*blur*)

#### Descrição

A técnica consiste em aplicar um filtro de convolução nos *pixels* com o objetivo de desfocar uma área de interesse na imagem.

<p>Dado Original</p> 	<p>Dado Anonimizado com a técnica de desfoque gaussiano.</p> 
<p><b>Aplicação</b> Apropriada para dados de imagem ou vídeo em que se deseja desfocar regiões de interesse, em geral faces, para minimizar o risco de identificação do titular.</p>	
<p><b>Limites</b></p> <ul style="list-style-type: none"> <li>• Aplicável em dados de imagens ou frame de vídeo.</li> <li>• Dificuldade de definir o limite dos parâmetros do desfoque para garantir a utilidade do dado e ao mesmo tempo preservar a privacidade.</li> <li>• Dificuldade em identificar quais regiões da imagem o desfoque deve ser aplicado para garantir a utilidade do dado e ao mesmo tempo preservar a privacidade.</li> </ul>	
<p><b>Pseudocódigo</b></p> <hr/> <p><b>Algoritmo:</b> Pseudocódigo da técnica de desfoque gaussiano (<i>blur</i>)</p> <hr/> <p><b>Entrada:</b> imagem: Imagem original; filtro: Filtro de <i>kernel</i> para aplicar a imagem</p> <p><b>Saída:</b> imagemAnonimizada: Imagem com desfoque gaussiano</p> <ol style="list-style-type: none"> <li>1 <math>imagemAnonimizada \leftarrow funcDesfoque(imagem, kernel)</math></li> <li>2 <b>return</b> <i>imagemAnonimizada</i></li> </ol>	

### Técnica de Pixelização

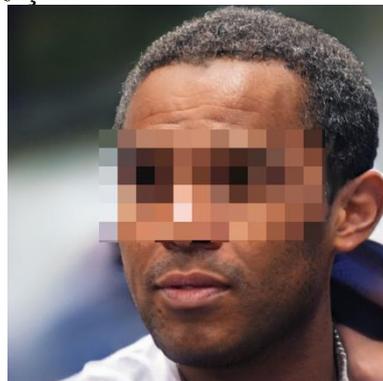
#### Descrição

A técnica consiste diminuir a resolução da imagem, ou em uma área de interesse dessa para reduzir a nitidez da imagem.

Dado Original



Dado Anonimizado com a técnica de pixelização.



#### Aplicação

Apropriada para dados de imagem ou vídeo em que se deseja *pixelização* regiões de interesse, em geral faces, para minimizar o risco de identificação do titular.

#### **Limites**

- Aplicável em dados de imagens ou frame de vídeo.
- Dificuldade de definir o limite dos parâmetros da *pixelização* para garantir a utilidade do dado e ao mesmo tempo preservar a privacidade.
- Dificuldade em identificar quais regiões da imagem a *pixelização* deve ser aplicada para garantir a utilidade do dado e ao mesmo tempo preservar a privacidade.

#### **Pseudocódigo**

**Algoritmo:** Pseudocódigo da técnica de *pixelização*

**Entrada:** imagem: Imagem original; filtro: Filtro de *kernel* para aplicar a imagem

**Saída:** imagemAnonimizada: Imagem com *pixelização*

1 *imagemAnonimizada* ← *funcPixelizacao(imagem, kernel)*

2 **return** *imagemAnonimizada*

## II. GLOSSÁRIO

**Agente de Tratamento:** O controlador e o operador.

**Anonimização:** Utilização de meios técnicos razoáveis e disponíveis no momento do tratamento, por meio dos quais um dado perde a possibilidade de associação, direta ou indireta, a um indivíduo.

**Banco de Dados:** Conjunto estruturado de dados pessoais, estabelecido em um ou em vários locais, em suporte eletrônico ou físico.

**Cenário do Jornalista:** Situação em que o risco de reidentificação de todos os titulares é igual, pois o atacante não possui nenhuma informação a priori de quais titulares estão presentes no conjunto de dados.

**Cenário do Promotor:** Situação em que o risco de reidentificação de um titular em específico é aumentado, pois o atacante sabe que os dados desse titular estão no conjunto de dados.

**Conjunto de Dados:** Vide Banco de Dados.

**Controlador:** Pessoa natural ou jurídica, de direito público ou privado, a quem competem as decisões referentes ao tratamento de dados pessoais.

**Dado Anonimizado:** Dado relativo ao titular que não possa ser identificado, considerando a utilização de meios técnicos razoáveis e disponíveis na ocasião de seu tratamento.

**Dado em Fluxo:** Dado gerado continuamente a uma alta taxa de velocidade, com tamanho potencialmente infinito e necessidade de processamento imediato.

**Dado Pessoal:** Informação relacionada a pessoa natural identificada ou identificável.

**Dado Pessoal Sensível:** Dado pessoal sobre origem racial ou étnica, convicção religiosa, opinião política, filiação a sindicato ou a organização de caráter religioso, filosófico ou político, dado referente à saúde ou à vida sexual, dado genético ou biométrico, quando vinculado a uma pessoa natural.

**Equivalência de Classe:** Subconjunto de um conjunto que contém todos os elementos com algum valor de atributo igual a todos os elementos.

**Métrica Base:** Uma métrica base é um valor definido para mensurar o risco de reidentificação calculado unicamente com base no próprio conjunto de dados. Por exemplo, a Equivalência de Classe.

**Métrica Contextual:** Uma métrica contextual é definida como sendo uma métrica derivada de uma métrica base, com a incorporação de elementos particulares.

**Operador:** Pessoa natural ou jurídica, de direito público ou privado, que realiza o tratamento de dados pessoais em nome do controlador.

**Tensor:** Estrutura matemática que estende uma matriz.

**Titular:** Pessoa natural a quem se referem os dados pessoais que são objeto de tratamento.

**Tratamento:** Toda a operação realizada com dados pessoais, como as que se referem a coleta, produção, recepção, classificação, utilização, acesso, reprodução, transmissão, distribuição, processamento, arquivamento, armazenamento, eliminação, avaliação ou controle da informação, modificação, comunicação, transferência, difusão ou extração.

**Variável Dependente do Contexto:** Uma variável dependente do contexto é definida como sendo uma característica interna do agente de tratamento que pode afetar o cálculo do risco de reidentificação.