

## **O EMPREGO DE ANÁLISE POR PRINCIPAIS COMPONENTES NA INSTRUMENTALIZAÇÃO DE CLASSIFICADORES DE DESASTRES**

*Samir Batista Fernandes<sup>1</sup>*

<https://orcid.org/0000-0001-9976-5318>

*Wagner Dos Anjos Carvalho<sup>2</sup>*

<https://orcid.org/0000-0003-2754-4414>

*Marcos dos Santos<sup>3</sup>*

<https://orcid.org/0000-0002-2074-6323>

### **RESUMO**

Este estudo abordou a crescente complexidade dos desastres, visando aprimorar as estratégias de classificação da intensidade de desastres nos municípios do Rio de Janeiro. A pesquisa visou preencher a lacuna de métodos eficientes e objetivos de classificação que integram a análise por principais componentes para uma melhor gestão de riscos e resposta a desastres. Foi adotada uma abordagem quantitativa, utilizando dados do Sistema Integrado de Informações sobre Desastres (S2ID), para categorizar desastres com base em sua intensidade e impacto. O estudo foi estruturado em quatro seções, iniciando com a contextualização da problemática de desastres no sistema brasileiro, seguido pelos materiais e métodos estatísticos empregados na análise. Os resultados indicaram a formação de quatro classes distintas de desastres, refletindo as capacidades de investimento e a severidade dos impactos nos municípios. A conclusão ressaltou a relevância da classificação para o planejamento de ações e alocação de recursos pela Defesa Civil. A contribuição central do trabalho reside no desenvolvimento de um modelo classificatório que pode orientar políticas públicas e estratégias de prevenção e mitigação, alinhado ao escopo de estudos em gestão de riscos de desastres.

**Palavras-chave:** Defesa Civil; Análise dos componentes principais; Intensidade de desastres.

---

<sup>1</sup> Universidade de São Paulo (artigo derivado de trabalho de conclusão de curso), Oficial do CBMERJ, Mestre em Segurança e Defesa Civil, Especialista em de Ciência de dados – USP, Especialista em Defesa Civil – Unisul. Pesquisador do NEPED UFF e da Defesa Civil do Estado do Rio de Janeiro. Brasil. samirfernandes@usp.br

<sup>2</sup> Professor Assistente, coordenador da pós-graduação Indústria 4.0 e professor dos cursos de MBA em Gestão e Tecnologias Educacionais e Controladoria e Finanças da Faculdade Presbiteriana Mackenzie Rio. Coordenador da pós-graduação EAD em Engenharia de Custos, Engenharia Calculista e Inteligência Artificial para Contabilidade na Faculdade Unyleya. Professor do programa de pós-graduação em Engenharia de Segurança do Trabalho e de Gestão de Projetos da UNIVERSO. Professor do MBA em Gerenciamento de projetos (ESTÁCIO). Professor do MBA em Gestão Estratégica de Marketing e Vendas (ESTÁCIO). Professor do Instituto Brasileiro da Administração Municipal (IBAM). Brasil. wagner.acarvalho@gmail.com

<sup>3</sup> pós-doutorado em Ciências e Tecnologias Espaciais pelo Instituto Tecnológico de Aeronáutica (ITA) e outro em Engenharia de Produção pela Universidade Federal Fluminense (UFF). É Doutor em Engenharia de Produção (pela UFF) na Linha de Pesquisa: Sistemas, Apoio à Decisão e Logística. É Mestre em Engenharia de Produção (pela COPPE/UFRJ) na Linha de Pesquisa de Pesquisa Operacional. Coordenador do primeiro MBA em Pesquisa Operacional e Tomada de Decisão do Brasil. Professor da graduação e do mestrado em Engenharia da Computação do Instituto Militar de Engenharia (IME). Professor do Programa de Pós-graduação em Engenharia de Produção da UFF. Professor do MBA em Data Science e Analytics da Universidade de São Paulo (USP). Professor de Data Driven Business do IBMEC. Brasil. marcosdossantos\_coppe\_uffrj@yahoo.com.br

## **THE USE OF PRINCIPAL COMPONENT ANALYSIS IN THE INSTRUMENTATION OF DISASTER CLASSIFIERS**

### **ABSTRACT**

This study addressed the growing complexity of disasters, aiming to improve disaster intensity classification strategies in the municipalities of Rio de Janeiro. The research aimed to fill the gap of efficient and objective classification methods that integrate principal component analysis for better risk management and disaster response. A quantitative approach was adopted, using data from the Integrated Disaster Information System (S2ID), to categorize disasters based on their intensity and impact. The study was structured into four sections, starting with the contextualization of the problem of disasters in the Brazilian system, followed by the materials and statistical methods used in the analysis. The results indicated the formation of four distinct classes of disasters, reflecting investment capacities and the severity of impacts in municipalities. The conclusion highlighted the relevance of the classification for action planning and resource allocation by Civil Defense. The central contribution of the work lies in the development of a classificatory model that can guide public policies and prevention and mitigation strategies, aligned with the scope of studies in disaster risk management.

**Keywords:** Civil Defense; Principal component analysis; Disaster intensity.

**Artigo Recebido em 26/01/2024**

**Aceito em 01/03/2024**

**Publicado em 30/03/2024**

## 1- INTRODUÇÃO

Os desastres são eventos que produzem danos e prejuízos, além de comprometerem a capacidade de resposta local (KOBİYAMA, MENDONÇA, *et al.*, 2006). No Estado do Rio de Janeiro os desastres aumentam em frequência e intensidade, além de ceifar inúmeras vidas (DOURADO, ARRAES e SILVA, 2012). Os desastres são registrados no sistema integrado de informações sobre desastres – [S2ID], no qual os dados são inseridos em processos digitais com o objetivo duplo de validação do reconhecimento do desastre e da possibilidade de obtenção de instrumentos para a recuperação da área afetada (FERNANDES, SANTOS e SILVA, 2022).

Portanto, nas localidades em que exista a declaração de um desastre, deve-se considerar que houve variáveis mensuráveis, tais como: quantidade de precipitação pluviométrica, quantidade de óbitos, quantidade de pessoas desalojadas, número de solicitações de vistorias emergenciais, um valor monetário que represente os danos e prejuízos etc. (CHMUTINA e BOSHER, 2017).

A quantificação do desastre é mensurada na plataforma [S2ID] e posteriormente os processos de reconhecimento no Brasil são regulamentados por ritos normativos sendo avaliados segundo a interpretação de um ser humano, por meio de documentações comprobatórias dos dados inseridos. Entretanto, o ser humano tem dificuldade de analisar muitos dados e encontrar informações úteis. Por isso, usa-se *data mining*, que é uma técnica que descobre padrões, modelos e informações em bancos de dados. Cabena *et al.* (1997) foram os primeiros a explicar o que é *data mining*, como fazer e para que serve. Depois, surgiram novas técnicas de *data mining*. Han, Kamber e Pei (2011) atualizaram o livro de *data mining* e mostraram as novas técnicas e as aplicações em várias áreas, como comércio, saúde, educação, segurança, biologia, entre outras. (CABENA *et al.*, 1997; HAN, KAMBER e PEI, 2011).

Pesquisas que combinam análises de dados pretéritos, o emprego de técnicas estatísticas e a busca por padrões não triviais por meio da análise por principais componentes que possam categorizar classes de desastres num banco de dados são vantajosas e constituem uma ferramenta excelente visto que o ser humano teria dificuldade em realizar inferências em grande volume de dados ou análises complexas. Esse é o objetivo do processo de descoberta de conhecimento em banco de dados [DCBD], que foi definido por Fayyad, Piatetsky-Shapiro e Smyth (1996) como o processo não trivial de descobrir padrões novos, válidos, úteis e compreensíveis a partir de bancos de dados. Desde então, muitas aplicações de DCBD foram realizadas em diversas áreas, incluindo a proteção e defesa civil, como mostra o estudo de Fernandes, Santos e Silva (2022), que utiliza a técnica de DCBD para detectar anomalias associadas a ocorrências de desastres em um município do Rio de Janeiro. (FAYYAD, SHAPIRO e SMYTH, 1996; FERNANDES, SANTOS e SILVA, 2022).

Desta forma, torna-se cada vez mais necessário estudos que identifiquem as variáveis que possuam relação com os elementos deflagradores de desastres e permitam cada vez mais a sociedade se proteger, além de serem elaboradas políticas públicas mais efetivas. As pesquisas que envolvem desastres representam um aspecto relevante sobre a sociedade e tem o potencial de auxiliar a salvar vidas, porém quando se aborda sobre o tema desastre, muitos campos do conhecimento lançam sua visão a partir de uma ciência específica. Ou seja, a geotecnia discorre do desastre sobre o olhar das ciências duras, a sociologia e ciências afins aborda o tema dos desastres sob olhar do campo social e assim com todas as ciências. A ciência dos desastres é uma área do conhecimento relativamente nova e que se apropria de outros campos para construção de seus avanços (SOUSA, SANTOS e SOUZA, 2020).

Apesar do tema ser tão importante, por se tratar de um assunto que envolve a proteção humana, a classificação da intensidade dos desastres no

[S2ID] só passou a ser exigida por lei após 2022 conforme resposta da Secretaria Nacional de Proteção e Defesa Civil em manifestação pública no decorrer desta pesquisa. Como os dados de interesse e que são inseridos no [S2ID] são, em sua maioria, métricas quantitativas a hipótese é que existe a possibilidade de elaboração de classes de desastres. Para Azevedo (2018) o problema a ser resolvido pode ser pesquisado, desde que tenha um objetivo claro, factível e resultados atingíveis por meio de um processo científico.

Dessa forma, o objetivo da pesquisa é agrupar os desastres por classes utilizando a análise por principais componentes considerando as variáveis selecionadas na plataforma [S2ID] dos municípios do Rio de Janeiro. A pesquisa contribuirá com uma técnica de classificação de desastres e poderá ajudar no planejamento de melhores políticas públicas visando a redução dos riscos de desastres. Portanto a presente pesquisa destina-se aos profissionais que atuam na gestão de redução dos riscos de desastres e aos demais colaboradores envolvidos no sistema de proteção e defesa civil.

O trabalho está estruturado em quatro seções na qual a primeira é contextualizada uma breve problematização da ausência de classificadores de desastres utilizando análise por principais componentes no sistema brasileiro, seguido pela segunda seção na qual são apresentados os materiais e métodos implementados na pesquisa. Em seguida são apresentado e discutido os resultados na terceira seção e finalizado com as considerações finais na quarta seção.

## 2- DESENVOLVIMENTO

Neste tópico, serão abordados os materiais e métodos utilizados para a implementação de um algoritmo de aprendizado de máquina, começando pelo objeto de estudo, a escolha de uma base de dados, os materiais e os métodos utilizados (testes, técnicas e análises).

## **2.1 - Objeto de estudo**

O Estado do Rio de Janeiro, dividido em 92 municípios e possuindo uma área de aproximadamente 43.780 km<sup>2</sup>, é uma região propensa a desastres naturais, como deslizamentos de terra e inundações. A extensão territorial do estado e a grande quantidade de municípios tornam a região especialmente vulnerável a esses eventos.

Os desastres são registrados na plataforma do governo federal [S2ID] para registro e análise visando reconhecimento junto ao governo estadual ou federal. O processo de análise dos desastres é realizado por meio de análise documental e cumprimento de ritos normativos utilizando o conhecimento técnico do ser humano.

## **2.2 - Material utilizado**

A pesquisa utilizou a base de dados [S2ID] do governo federal brasileiro. O acesso à base de dados exigiu uma solicitação formal ao governo federal, por meio de uma manifestação pública. Foram selecionados os desastres com as seguintes características:

- Reconhecidos pelo governo federal tipificados segundo a Classificação e codificação brasileira de desastres – [COBRADE] como:(1) Tempestade convectiva ou chuvas intensas, (2) frentes frias ou zonas de convergência, (3) Deslizamentos, (4) inundações e (5) enxurradas;
- Danos e prejuízos relacionados à precipitação pluviométrica em seu território;
- Informação sobre a intensidade de precipitação pluviométrica em 24h e
- Registrado na base de dados [S2ID], por meio do formulário de informações de desastres – FIDE e no parecer técnico.

Com base no critério acima foram selecionados noventa e seis (96) processos com registros sobre desastres no [S2ID]. A partir da análise dos processos foram selecionadas variáveis de interesse e registradas em uma nova base de dados, na forma de planilha, conforme a tabela 1 abaixo:

**Tabela 1– Variáveis selecionadas para implementação de algoritmo**

<b>Variável</b>	<b>Significado da variável</b>
Município	Nome do ente federativo que sofreu o desastre
População	Quantidade absoluta da população municipal declarada no processo
Área	Extensão territorial municipal
Orçamento anual	Total de receita projetada para o exercício fiscal anual
Receita Anual	Total de receita recebida para o exercício fiscal anual
Danos e Prejuízos	Valor monetário apurado decorrente dos efeitos do desastre
Afetados	Total de população apurada impactada pelo desastre
Capacidade de investimento na resposta ao desastre	Razão entre a receita anual e os danos e prejuízos apurados por desastre
Capacidade de projeção financeira	Razão entre a receita anual e o orçamento anual
Densidade populacional de afetados	Razão entre os afetados por desastres e a população total

Fonte: Autores (2023)

### 2.3 - Limitações do estudo

Um aspecto importante da nossa abordagem metodológica foi o tratamento de dados ausentes, conhecidos como *missing values*. Conforme Little e Rubin (2002) destacam, a ausência de dados pode introduzir viés significativo em análises estatísticas, afetando a validade dos resultados. Em nossa pesquisa, optamos por excluir da base de dados 16 processos com informações sobre precipitação pluviométrica ausentes. Esta decisão, embora estratégica para manter a integridade analítica do modelo, limitou a classificação dos 16 processos de desastres que poderiam ser devidamente agrupados. Essa limitação refletiu a falta de padrão na inserção de dados para obtenção do reconhecimento dos desastres naturais, pois embora excluídos os processos, poderiam oferecer informações relevantes.

Ademais, a tipologia dos desastres foi um fator determinante em nossa seleção de dados. Seguindo a perspectiva de Alexander (2000), que ressalta a

importância de entender as especificidades de cada tipo de desastre na gestão eficaz de riscos, focamos em desastres de natureza súbita, como chuvas intensas. Desastres de natureza gradual, ou de somatório de efeitos parciais, foram excluídos devido à complexidade de sua modelagem e análise. Esta abordagem focada, embora útil para a análise de eventos com impactos imediatos e mensuráveis, limitou o escopo da pesquisa ao não incluir desastres de evolução mais lenta e que não estivessem relacionado à precipitação pluviométrica.

A obtenção de dados do governo federal brasileiro também representou um desafio significativo. Embora os dados não fossem sensíveis, o processo de solicitação e acesso a estes dados junto ao órgão federal responsável foi dificultado por uma série de entraves burocráticos, refletindo uma limitação em termos de disponibilidade e acessibilidade de dados governamentais. Esta situação evidencia a dependência de pesquisas sobre a cooperação de órgãos governamentais para o acesso a informações importantes.

Além disso, é importante reconhecer que os resultados obtidos, embora valiosos, representam apenas uma parcela do universo dos desastres analisados no banco de dados. Este fato sublinha a necessidade de considerar as limitações do estudo ao interpretar os resultados. Assim, destacamos a importância de conduzir pesquisas futuras mais abrangentes que possam oferecer uma visão mais completa e diversificada do panorama dos desastres. Estudos futuros devem se esforçar para incluir uma gama mais ampla de desastres e considerar a aplicação prática dos métodos para profissionais da área, a fim de fortalecer a gestão de riscos de desastres e dos gerenciamentos dos desastres.

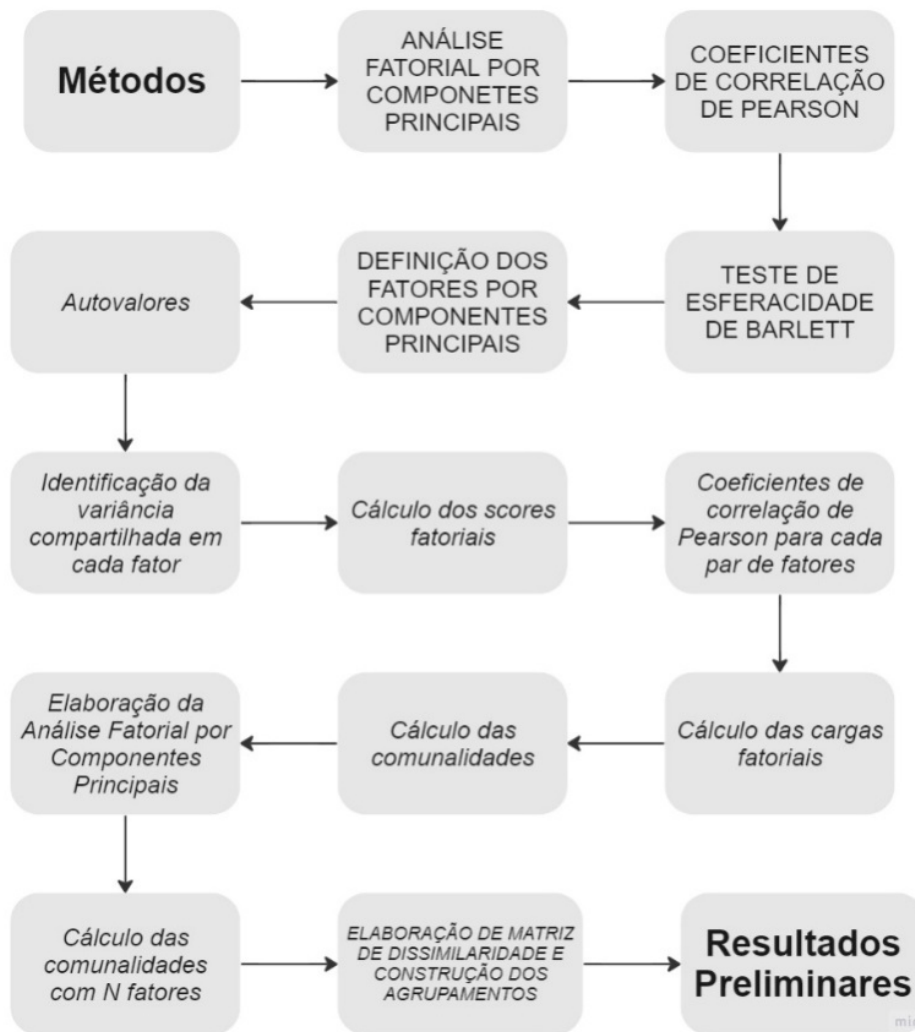
### **3- METODOLOGIA UTILIZADA**

A metodologia aplicada, ilustrada na Figura 1, segue as etapas sugeridas por Fávero e Belfiore (2017) para implementar o algoritmo de



classificação de desastres. Essas etapas envolvem a utilização sequencial de técnicas e testes estatísticos por meio do *software* R, visando à obtenção do resultado.

**Figura 1**– Metodologia utilizada na classificação de desastres



Fonte: Os autores (2023)

### 3.1 - Análise de agrupamentos

A análise de agrupamentos, também conhecida como análise de conglomerados ou análise de clusters, é um conjunto de técnicas exploratórias úteis para verificar a existência de comportamentos semelhantes entre

observações e criar grupos homogêneos em relação a determinadas variáveis(KING, 2014). Essas técnicas não são preditivas para outras observações fora da amostra inicial, permitem comparar a medida de similaridade entre as observações dentro de um mesmo grupo e possui a limitação da inclusão de novas observações, ou seja, novas variáveis exige uma reaplicação da modelagem(FÁVERO e BELFIORE, 2017). Para a realização da pesquisa foi selecionado o método de análise fatorial por componentes principais que teve grande contribuição de Karl Pearson(ANDERSON, 2003; LOVRIC e ARSHAM, 2011).

### **3.2 - Análise fatorial por componentes principais**

A análise fatorial é uma técnica multivariada exploratória que busca estabelecer novas variáveis, chamadas de fatores, a partir do agrupamento de variáveis originais que apresentam coeficientes de correlação elevados. Essa técnica é útil para reduzir a dimensão dos dados e identificar relações entre as variáveis originais, mas não possui caráter preditivo para outras observações não presentes na amostra(FÁVERO e BELFIORE, 2017).

A análise fatorial por componentes principais é o método mais utilizado na redução de dimensionalidade, pois permite determinar outro conjunto de variáveis resultantes da combinação linear do primeiro conjunto de variáveis(RENCHER e CHRISTENSEN, 2012). A análise fatorial por componentes principais tem quatro objetivos principais: (1) redução estrutural; (2) verificação da validade de constructos previamente estabelecidos; (3) elaboração de “rankings”; e (4) extração de fatores ortogonais para posterior uso em técnicas multivariadas confirmatórias que necessitam de ausência de multicolinearidade. Os fatores extraídos podem ser utilizados como variáveis explicativas de outras variáveis, como em modelos multivariados confirmatórios, como regressão múltipla. Em resumo, a análise fatorial foi uma técnica útil para identificar relações entre variáveis originais, reduzir a

dimensão dos dados e criar variáveis a partir do agrupamento de variáveis originais com coeficientes de correlação elevados(FÁVERO e BELFIORE, 2017).

### 3.3 - Coeficientes de correlação de Pearson para cada par de variáveis

A partir de um base de dados que apresente  $n$  observações (na pesquisa foram chamadas de processos inseridos no [S2ID]), e para cada observação  $i$  ( $i=1,n$ ) valores correspondentes a cada uma das  $K$  variáveis métricas  $X$ , conforme mostra a tabela 2 foi possível correlacionar a cada processo analisado como uma observação. Para cada processo foi utilizado as variáveis município, população, área, orçamento anual, receita anual, danos e prejuízos, afetados, capacidade de investimento na resposta ao desastre, capacidade de projeção financeira e densidade populacional de afetados.

**Tabela 2– Modelo geral de um banco de dados para elaboração de análise fatorial**

Observação $i$	$X_{1i}$	$X_{2i}$	...	$X_{ki}$
1	$X_{11}$	$X_{21}$		$X_{k1}$
2	$X_{12}$	$X_{22}$		$X_{k2}$
3	$X_{13}$	$X_{23}$		$X_{k3}$
$\vdots$	$\vdots$	$\vdots$	...	$\vdots$
$n$	$X_{1n}$	$X_{2n}$		$X_{kn}$

Fonte: FÁVERO e BELFIORE(2017)

Dado um processo  $p$  e  $q$  distintos na base de dados, pode-se representar um grau de similaridade entre eles pela seguinte expressão:

$$\rho_{pq} = \frac{\sum_{j=1}^k (X_{jp} - \bar{X}_p) \times (X_{jq} - \bar{X}_q)}{\sqrt{\sum_{j=1}^k (X_{jp} - \bar{X}_p)^2} \times \sqrt{\sum_{j=1}^k (X_{jq} - \bar{X}_q)^2}} \quad (1)$$

onde  $\rho_{pq}$  é o valor do coeficiente de Pearson entre dois processos distintos na base de dados e  $\bar{X}_p$  e  $\bar{X}_q$  representam a média de cada uma das linhas da base de dados. A expressão acima permite analisar a similaridade entre as linhas da base de dados e o comportamento em linha das observações para o conjunto de variáveis.

### **3.4 - Teste de esfericidade de Barlett**

O Teste de Esfericidade de Bartlett é um teste estatístico usado para comparar uma matriz de correlação  $\rho$  com uma matriz identidade de mesma dimensão. O teste é usado para avaliar a adequação dos dados para análise fatorial. Para Fávero e Belfiores (2017, p. 386)

“se as diferenças entre os valores correspondentes fora da diagonal principal de cada matriz não forem estatisticamente diferentes de 0, a determinado nível de significância, poderemos considerar que a extração dos fatores não será adequada”.

Portanto, o teste de Esfericidade de Bartlett testa a hipótese nula de que a matriz de correlação é uma matriz identidade, o que significa que não há correlações entre as variáveis. Um resultado significativo indica que os dados são adequados para a análise fatorial (LOVRIC e ARSHAM, 2011).

### **3.5 - Definição dos fatores por componentes principais**

Uma vez identificado adequação para análise fatorial é necessário empregar os fatores que carregam parcialmente uma característica das variáveis originais em agrupamentos de variáveis. Abaixo veremos os fatores que auxiliaram na determinação dos agrupamentos de variáveis.

- *Autovalores* – Os autovalores são medidas numéricas que indicam a importância relativa de cada componente principal na representação da variabilidade dos dados em uma análise fatorial por componentes principais. Eles representam a variância explicada por cada componente e são usados para determinar quantos componentes devem ser retidos na análise. Quanto maior o autovalor, mais importante é a contribuição daquele componente para

a estrutura dos dados(HAIR JR., BLACK, W.C., *et al.*, 2019),(JOHNSON e WICHERN, 2007).

- *Identificação da variância compartilhada em cada fator*– A identificação da variância compartilhada em cada fator na análise fatorial por componentes principais é essencial para entender a contribuição de cada fator para a variação total dos dados. A variação compartilhada é medida pelos autovalores, que indicam a quantidade de variação total dos dados explicada por cada fator. Quanto maior o autovalor, maior é a quantidade de variação explicada pelo fator correspondente. Assim, a identificação da variância compartilhada em cada fator permite selecionar os fatores mais relevantes para a explicação dos dados(TABACHNICK e FIDELL, 2013),(HAIR JR., BLACK, W.C., *et al.*, 2019).

É comum utilizar fatores maiores do que 1 para representar os fatores principais na análise fatorial por componentes principais porque isso permite reter mais variância total do conjunto de dados original. Por exemplo, se for definido apenas um fator principal, ele pode não ser suficiente para representar toda a variação dos dados, resultando em uma perda de informação importante. Portanto, ao aumentar o número de fatores principais, é possível capturar mais variância e informações relevantes do conjunto de dados original(JOLLIFFE, 2002).

- *Cálculo dos scores fatoriais* – O cálculo dos scores fatoriais é uma etapa importante na análise fatorial por componentes principais, pois permite obter valores numéricos para cada indivíduo (ou objeto) em relação a cada fator. Esses scores representam a pontuação de cada indivíduo em relação à importância ou influência de cada fator no conjunto de dados. Os scores são calculados através da combinação linear dos valores das variáveis originais com os pesos dos fatores correspondentes. Esses pesos são chamados de "*loadings*" e indicam a importância de cada variável para cada fator. O resultado é uma matriz de scores fatoriais que pode ser utilizada para análises posteriores, como a visualização dos dados em gráficos de dispersão ou a

comparação de grupos de indivíduos em relação a cada fator(JOLLIFFE, 2002).

- *Coefficientes de correlação de Pearson para cada par de fatores*– Os coeficientes de correlação de Pearson para cada par de fatores na análise fatorial por componentes principais são importantes para entender a relação entre os fatores e como eles se relacionam com as variáveis originais. Esses coeficientes medem a correlação entre os scores fatoriais de cada par de fatores e variam de -1 a 1. Um coeficiente positivo indica uma correlação positiva entre os fatores, enquanto um coeficiente negativo indica uma correlação negativa. Um coeficiente próximo de zero indica uma correlação fraca ou inexistente. A interpretação dos coeficientes de correlação de Pearson é fundamental para a interpretação dos resultados da análise fatorial por componentes principais(JOHNSON e WICHERN, 2007),(TABACHNICK e FIDELL, 2013).

- *Cálculo das cargas fatoriais*– As cargas fatoriais são os coeficientes que mostram a relação entre cada variável original e cada fator extraído na análise fatorial por componentes principais. Elas indicam o quanto cada variável contribui para a definição de cada fator. O cálculo das cargas fatoriais é realizado por meio da correlação entre cada variável original e cada fator, ponderada pelo inverso da raiz quadrada dos autovalores correspondentes. Quanto maior a carga fatorial de uma variável em um fator, maior é a sua contribuição para a definição desse fator(HAIR JR., BLACK, W.C., *et al.*, 2019; TABACHNICK e FIDELL, 2013).

- *Cálculo das comunalidades*– O cálculo das comunalidades na análise fatorial por componentes principais é importante para avaliar a quantidade de variância compartilhada entre as variáveis originais e os fatores extraídos. As comunalidades representam a proporção da variância total de cada variável original que pode ser explicada pelos fatores. Essa medida é calculada somando-se os quadrados das cargas fatoriais de cada variável e pode variar de 0 a 1. Quando as comunalidades são altas, isso indica que as

variáveis originais são bem representadas pelos fatores extraídos. Caso contrário, pode ser necessário reconsiderar o modelo fatorial ou adicionar mais fatores (JOLLIFFE, 2002), (HAIR JR., BLACK, W.C., *et al.*, 2019).

- *Elaboração da Análise Fatorial por Componentes Principais*– Na Análise Fatorial por Componentes Principais, é comum extrair fatores apenas com autovalores maiores do que 1. Essa abordagem é adotada porque, ao extrair apenas os fatores mais importantes, é possível reter a maior parte da variância total do conjunto de dados original. Em seguida, é possível realizar a rotação dos fatores para obter uma estrutura mais clara e interpretável. A interpretação dos fatores extraídos é realizada por meio das cargas fatoriais, que representam a correlação entre cada variável original e cada fator extraído (HU e BENTLER, 1999), (BORGES, ESTEVES, *et al.*, 2018).

- *Cálculo das comunalidades com N fatores*– O cálculo das comunalidades com N fatores é realizado a partir da soma dos quadrados das cargas fatoriais ao quadrado para cada variável em cada um dos N fatores. As comunalidades representam a quantidade de variância total de cada variável que pode ser explicada pelos fatores extraídos. Quando N é igual ao número total de variáveis, as comunalidades correspondem às variâncias das variáveis. As comunalidades são importantes para avaliar a adequação do modelo fatorial e identificar quais variáveis contribuem mais para a explicação dos fatores (FÁVERO e BELFIORE, 2017), (ENAP, 2019).

### **3.6 - Elaboração de matriz de dissimilaridade e construção dos agrupamentos**

Uma vez realizada as etapas na construção dos fatores por componentes principais, a próxima etapa foi construir a matriz de dissimilaridades que é usada na análise de *cluster* hierárquica para agrupar objetos semelhantes em *clusters*. Na análise fatorial por componentes principais, a matriz de dissimilaridades é calculada com base nos fatores

extraídos e nas cargas fatoriais dos objetos. Essa abordagem permite a incorporação da estrutura de dependência entre as variáveis no processo de agrupamento, resultando em análises mais precisas e interpretações mais confiáveis (BORG e GROENEN, 2005), (INDHUMATHI e SATHIYABAMA, 2010).

Na construção de agrupamentos, empregamos técnicas avançadas de análise para entender melhor as semelhanças entre diferentes tipos de desastres, com um foco particular no uso do agrupamento hierárquico dentro da análise fatorial por componentes principais. Este método nos permitiu identificar grupos de desastres que compartilham características similares, com uma relevância estatística, indicando que essas semelhanças não são meras coincidências, mas sim padrões significativos.

Para aprofundar nossa compreensão e validar a consistência desses grupos, utilizamos a Análise de Variância de um Fator (ANOVA). A ANOVA é uma ferramenta estatística projetada para comparar as médias de diferentes grupos, testando se as diferenças observadas entre essas médias são estatisticamente significativas ou se podem ser explicadas pelo acaso. Essa análise é crucial, pois ajuda a confirmar se os padrões que identificamos têm fundamento estatístico sólido.

O resultado dessa análise foi apresentado através de três componentes principais: a estatística F, o valor p e a tabela ANOVA. A estatística F nos ajudou a avaliar a significância das diferenças entre as médias dos grupos identificados. O valor p, por sua vez, indicou a probabilidade de que essas diferenças possam ter ocorrido por acaso, com valores baixos sugerindo que as diferenças são estatisticamente significativas. A tabela ANOVA resumiu a contribuição de cada fonte de variação nas diferenças observadas entre as médias dos grupos.

#### **4- RESULTADOS E DISCUSSÕES**

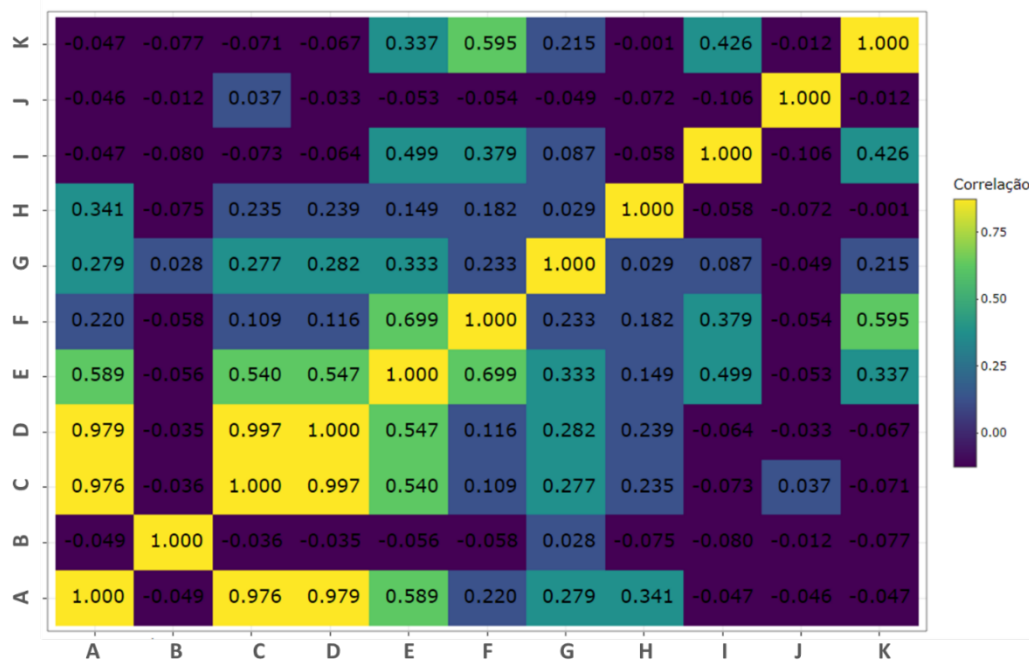


Inicialmente, foi observado baixas correlações entre todas as variáveis, excetuando área, população, receita e orçamento, após aplicamos o teste de esfericidade de Bartlett, que resultou em um valor estatisticamente insignificante, praticamente zero. Esse resultado indicou uma significância elevada, evidenciada pelas correlações fora da diagonal principal de cada matriz analisada. Isso sugeriu que as variáveis não estão correlacionadas ao ponto de invalidar a análise fatorial, um aspecto fundamental para a validade do nosso estudo.

Com base nesse achado, concluímos que a extração de fatores foi apropriada para o nosso conjunto de dados. Isso significou que foi possível, e correto, derivar novos fatores a partir das variáveis originais. Assim, a análise fatorial se mostrou uma abordagem adequada para investigar as dimensões subjacentes dentro dos nossos dados.

Este passo foi crucial, pois permitiu-nos simplificar a complexidade dos dados, reduzindo as variáveis originais a fatores mais gerenciáveis, sem perder informações essenciais. Tal abordagem não apenas facilitou a interpretação dos resultados, mas também reforçou a confiança na robustez e na relevância das conclusões extraídas do estudo, assegurando que a análise fatorial era apropriada e bem fundamentada para explorar as nuances dos dados coletados.

**Figura 2** - Mapa de calor das correlações de Pearson entre as variáveis



Legenda: A: População; B: Área; C: Receita Anual; D: Orçamento anual; E: Danos e Prejuízos; F: Afetados; G: Precipitação pluviométrica 24h; H: Densidade populacional; I: Capacidade de investimento na resposta ao desastre; J: Capacidade de projeção financeira; K: Densidade populacional de afetados;

Fonte: Os autores (2023)

Dada as 11 (onze) variáveis consideradas na base de dados foi realizada a análise dos autovalores com dez valores dos principais componentes – PC para posterior redução obtendo os seguintes resultados: PC1: 3.763; PC2:2.268; PC3: 1.079; PC4; 1.027; PC5; 0.902; PC6: 0.826; PC7: 0.575; PC8: 0.447; PC9: 0.107; PC10:0.011 e PC11: 0.

É possível observar na tabela 3 que PC1, PC2, PC3 e PC4 juntos possuem mais de 70% da variância acumulada e que pelo critério de raiz latente ou de critério de Kaiser foram considerados os PC maiores do que 1 (um), ou seja, PC1, PC2, PC3 e PC4 (FÁVERO e BELFIORE, 2017).

**Tabela 3 - Cálculo das cargas fatoriais**

	PC1	PC2	PC3	PC4	PC5	PC6	PC7	PC8	PC9	PC10	PC11
Autovalores	3.763	2.268	1.079	1.027	0.902	0.826	0.575	0.447	0.107	0.011	0.000
Prop. da Variância	0.342	0.206	0.097	0.093	0.082	0.075	0.052	0.040	0.009	0.001	0.000
Prop. da Variância acumulada	0.342	0.548	0.645	0.738	0.820	0.896	0.948	0.989	0.998	0.999	1.000

**Fonte:** Os autores (2023)

A partir dos autovalores e do cálculo da carga fatorial, foi possível determinar os “scores” fatoriais. Os “scores” permitiram a realização de obtenção do coeficiente de Pearson entre as variáveis originais e cada um dos fatores. Conforme a tabela 4, foi realizado o processo de “Standardized Loadings” (“Pattern Matrix”) e foi verificado que 4 componentes seriam suficientes para a representação da variância das variáveis originais. Dessa forma ficamos com PC1, PC2, PC3 e PC4. Refazendo o teste de esfericidade de Bartlett para a nova matriz criada foi possível verificar que o valor estatisticamente igual a 0 (zero), a elevado nível de significância, correspondentes fora da diagonal principal de cada matriz o que significa que mesmo com as perdas de variância originais houve ainda uma boa representação estatística dessas.

Quando olhamos para os dados de desastres, como a quantidade de pessoas afetadas, os custos dos danos e prejuízos, e o tamanho da população, percebemos que essas informações variam muito. Algumas dessas informações, como o número de pessoas, o dinheiro gasto e recebido, e os custos dos danos são muito importantes para entender o impacto de um desastre. Essas são as informações que mais se destacam na PC1. Outros detalhes, como a quantidade de pessoas afetadas e a densidade da população no local do desastre, se destacam em uma segunda análise (PC2).

Ao analisar todos esses dados juntos, conseguimos ver padrões e entender melhor como cada detalhe contribui para identificar a intensidade dos desastres. Isso nos ajuda a criar um modelo para classificar os desastres de forma mais eficiente, levando em conta todas essas variáveis importantes.

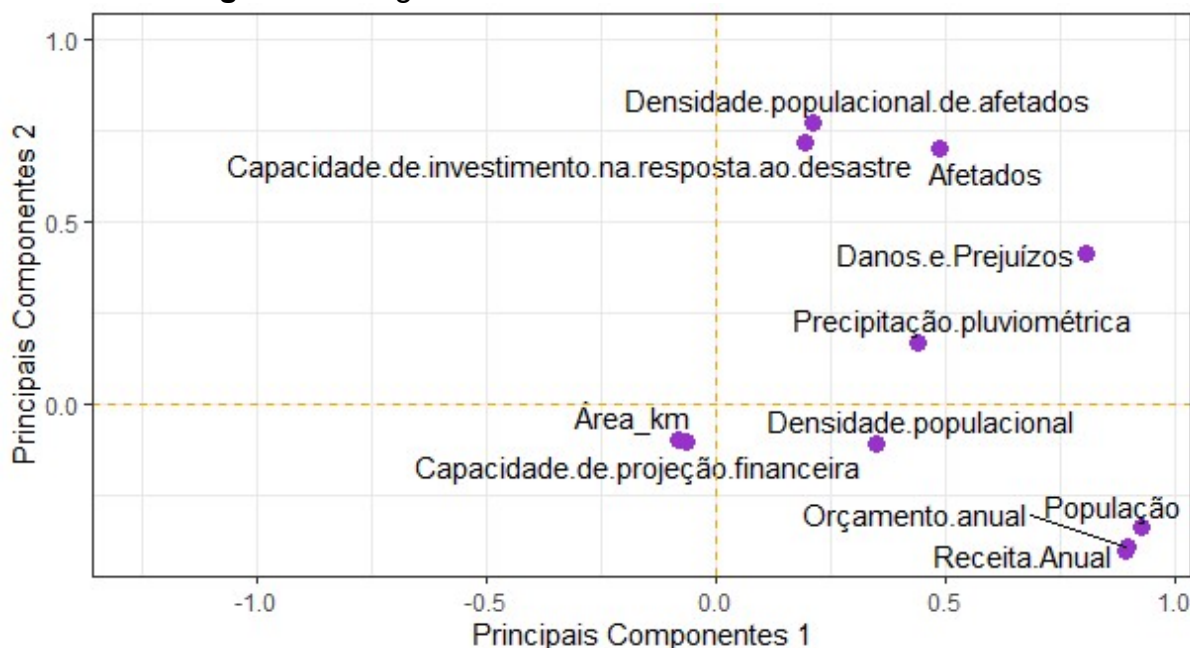
**Tabela 4**–“StandardizedLoadings” (“PatternMatrix”) baseadas na matriz de correlação (cont...)

<b>Variável</b>	<b>PC1</b>	<b>PC2</b>	<b>PC3</b>	<b>PC4</b>	<b>Comunalidades</b>
População	0.93	-0.34	-0.03	0.00	0.98
Área	-0.08	-0.10	0.69	-0.43	0.68
Orçamento anual	0.89	-0.40	0.07	0.08	0.97
Receita Anual	0.90	-0.39	0.05	0.02	0.96
<b>Variável</b>	<b>PC1</b>	<b>PC2</b>	<b>PC3</b>	<b>PC4</b>	<b>Comunalidades</b>
Danos e Prejuízos	0.81	0.41	0.06	0.03	0.83
Afetados	0.49	0.70	-0.06	0.02	0.73
Precipitação pluviométrica	0.44	0.17	0.40	-0.10	0.39
Densidade populacional	0.35	-0.11	-0.59	-0.17	0.52
Capacidade de investimento na resposta ao desastre	0.20	0.72	-0.01	-0.02	0.56
Capacidade de projeção financeira	-0.06	-0.10	0.26	0.88	0.86
Densidade populacional de afetados	0.21	0.77	0.03	0.10	0.65

**Fonte:** Os autores (2023)

Na elaboração visual da carga fatorial dos principais fatores, conforme a figura 3, observa-se que população, orçamento e receita anual estão bem agrupadas e possuem uma correlação positiva para o PC1 elevada. Outra correlação é observada ao se agrupar densidade populacional, desastres e afetados, ou seja, uma maior quantidade de pessoas afetadas no desastre aumenta a densidade de pessoas afetadas no território e comprometem a capacidade de investimento na resposta.

**Figura 3 - Carga Fatorial dos dois Primeiros Fatores**



Fonte: Os autores (2023)

Na análise dos grupos, aplicou-se a Análise de Variância dos Fatores (ANOVA) para avaliar como as características dos desastres variam entre grupos previamente definidos. O método ANOVA possibilitou identificar que as variações nas características dos desastres eram mais acentuadas entre os grupos do que dentro deles. Isso foi observado através da comparação das médias quadráticas dos agrupamentos, que se mostraram superiores às médias quadráticas dos resíduos, indicando diferenças significativas entre os grupos.

Além disso, o valor F conforme a tabela 5, obtido na análise destacou uma variável específica como sendo a mais discriminante entre os grupos, de maneira estatisticamente significativa. Isso denota a capacidade dessa variável de distinguir entre diferentes tipos de desastres. Importante ressaltar que, com os valores de probabilidade [Pr] estatisticamente iguais a zero para todos os fatores analisados, a hipótese nula — de que não há diferença significativa entre as médias dos grupos — foi rejeitada.

Esses resultados permitiram concluir que é possível, com significância estatística, agrupar os desastres em categorias com características similares entre si, mas distintas quando comparadas a outros grupos. Portanto, a pesquisa contribuiu para o avanço na classificação dos desastres, oferecendo uma base sólida para a elaboração de estratégias diferenciadas de gestão e resposta a cada tipo de desastre.

**Tabela 5**–“Output” da ANOVA

Fator	Agrupamento	SomaSq	MédiaSq	Fvalor	Pr(>F)
Fator 1	Agrupamentos	88.83	9.870	137.6	<2e <sup>-16</sup>
	Resíduos	6.17	0.072		
Fator 2	Agrupamentos	84.85	9.427	79.84	<2e <sup>-16</sup>
	Resíduos	10.15	0.0118		
Fator 3	Agrupamentos	86.03	9.559	91.66	<2e <sup>-16</sup>
	Resíduos	8.97	0.104		
Fator 4	Agrupamentos	52.32	5.813	11.71	<13e <sup>-12</sup>
	Resíduos	42.68	0.496		

Fonte: Os autores (2023)

A partir da validação dos resultados acima foi possível criar 10(dez) agrupamentos de desastres classificados com rótulos de 1 até 10 de forma qualitativa, entretanto os grupos de desastres 1,2 e 3 foram mantidos e os desastres de 4 até 10 foram agrupados e receberam o nome de raros representados no apêndice. Cabe destacar que as classes de desastres 1, 2 e 3 possuem a maior quantidade de desastres na qual totalizam 81 processos de desastres.

Desta forma, a classe 1, possuiu uma analogia com os desastres de Nível I ou de pequena intensidade da portaria 260 de 12 de fevereiro de 2023 do governo federal brasileiro. Os dados indicaram que a classe 1

comprometem, em média, a capacidade de investimento em 1,99% da receita corrente líquida, impactam com danos médios de 6,8 milhões de reais e em média possuem uma densidade populacional de 416,9 habitantes por km<sup>2</sup>, com um número médio de 7.017 afetados. A precipitação pluviométrica média foi de 111,5 mm em 24 horas. Esses desastres podem ser gerenciados com os recursos locais e medidas administrativas excepcionais, sem necessidade de aporte significativo de recursos de níveis estadual ou federal (BRASIL, 2022).

A classe 2 possuiu relação com os desastres de Nível II ou de média intensidade e comprometem a capacidade de investimento na média de 9,20% da receita corrente líquida, entretanto apresentou danos médios tratados de 18,08 milhões e uma densidade populacional média de 182,36 habitantes por km<sup>2</sup>. O número médio de afetados foram 13.384, e a precipitação média foi de 114,4 mm. Estes desastres requerem não apenas os recursos locais, mas também o aporte de recursos do estado e/ou da União para restabelecer a normalidade.

A classe 3, referente a desastres de Nível III ou de grande intensidade, em média possuiu o maior comprometimento toda capacidade do investimento. Atingindo o valor de 12,71% da receita corrente líquida, danos médios de 23,6 milhões e a mais alta densidade populacional de 1.766,3 habitantes por km<sup>2</sup>, afetando em média 46.893 pessoas, com uma precipitação de 109,2 mm. Estes eventos podem ser considerados como de grande magnitude, resultando em comprometimento significativo do orçamento local devido aos danos e afetados elevados, exigindo uma ação coordenada das três esferas de atuação do Sistema Nacional de Proteção e Defesa Civil e, em alguns casos, de ajuda internacional.

A classe de desastres raros, potencialmente desastres de Nível III, porém podem ser nível I ou II. Necessitam de investigação detalhadamente.

#### **4.1 - Recomendações Práticas e Implicações Políticas**

A utilização de análises por componentes principais, conforme demonstrado, pode ser uma ferramenta eficaz na classificação rápida e precisa da intensidade dos desastres. Esta técnica pode auxiliar na identificação de padrões e correlações entre as variáveis consideradas. A integração dessa metodologia em trabalhos de avaliação de danos e prejuízos visando a classificação de desastres oferece uma análise mais precisa, permitindo uma melhor identificação e alocação de recursos.

Além disso, é fundamental a capacitação contínua dos profissionais da proteção civil em técnicas de análise de dados avançadas considerando o avanço da ciência dos desastres e suas múltiplas áreas do campo do conhecimento. O conhecimento técnico aprofundado nessas áreas permitirá uma aplicação mais efetiva das ferramentas e metodologias propostas, otimizando as estratégias de prevenção e resposta a desastres.

Os resultados deste estudo também têm implicações políticas significativas. Primeiramente, evidenciam a necessidade de políticas públicas que promovam a adoção de tecnologias avançadas de análise de dados no campo da gestão de desastres. Adicionalmente, as contribuições dessa pesquisa reforçam a importância de uma abordagem baseada em dados na formulação de políticas de proteção civil. Políticas que se apoiam em análises precisas e detalhadas de dados podem levar a uma gestão de desastres mais proativa e eficiente. Isso implica na necessidade de colaboração entre cientistas de dados, profissionais de proteção civil e formuladores de políticas, visando integrar conhecimentos técnicos e experiências práticas no processo decisório.

Por fim, os resultados deste estudo podem ser usados para argumentar a favor de uma revisão e fortalecimento das normativas relacionadas à gestão de dados de desastres. Isso inclui a necessidade de garantir a qualidade, a



acessibilidade e a atualização constante dos dados, aspectos fundamentais para a eficácia das análises realizadas.

## **5- CONSIDERAÇÕES FINAIS**

Utilizando o material e a metodologia empregada e apresentado os resultados foi possível agrupar os desastres por classes utilizando a análise por principais componentes considerando as variáveis selecionadas na plataforma [S2ID] dos municípios do Rio de Janeiro. A criação de classes de desastres é de extrema importância para aprimorar as estratégias de prevenção e resposta a emergências. A classificação por intensidade permite uma análise mais precisa e adequada dos riscos associados a cada tipo de desastre, o que pode auxiliar na definição de medidas preventivas mais eficazes e na mobilização de recursos em situações de crise.

Nesse sentido, a utilização de análise por principais componentes para classificar os desastres pode ser extremamente útil. Algoritmos de aprendizado de máquina podem analisar um grande volume de dados na qual o ser humano seria incapaz de analisar para identificar padrões e correlações que permitam classificar o evento com precisão e rapidez.

Em resumo, a criação de classes de desastres e a utilização de análise por principais componentes para classificar os eventos podem contribuir significativamente para aprimorar a gestão de riscos e emergências. Essas medidas podem salvar vidas, preservar o patrimônio público e privado e garantir a segurança da população em situações de crise.

Embora a classificação de desastres por intensidade através da análise por principais componentes tenha provado ser eficaz, a aplicação de métodos preditivos pode oferecer uma compreensão ainda mais profunda e a capacidade de antecipar eventos de desastre com maior precisão.

Para pesquisas futuras os algoritmos como redes neurais artificiais, máquinas de vetores de suporte, ou até mesmo técnicas de aprendizado

profundo, podem ser utilizados para classificar e prever a intensidade de desastres com base em novas entradas de dados. Isso não se limita a um único método, mas abre um campo vasto para a exploração de várias abordagens e comparações de sua eficácia.

Um aspecto particularmente valioso dessas abordagens preditivas é seu potencial capacidade de identificar "cisnes negros" e "rinocerontes cinzas" – eventos raros, imprevisíveis e de alto impacto, bem como ameaças conhecidas, mas frequentemente negligenciadas. A detecção precoce desses fenômenos pode ser fundamental para mitigar riscos que, embora improváveis, possuem consequências devastadoras.

A análise preditiva, aplicada no contexto de desastres, pode revolucionar a maneira como as estratégias de prevenção e resposta são formuladas. A capacidade de prever com precisão a intensidade e o impacto potencial de um desastre antes de sua ocorrência pode permitir uma mobilização de recursos mais eficiente, uma melhor preparação das equipes de emergência e, em última instância, a minimização do impacto sobre a vida humana e o patrimônio.

Portanto, encorajamos a continuação desta linha de pesquisa, aproveitando os avanços na ciência de dados e aprendizado de máquina para fortalecer ainda mais as capacidades de prevenção e resposta a desastres em âmbito municipal no Brasil e além. A implementação de tais sistemas analíticos preditivos não só complementar os esforços existentes, mas também auxiliará outras práticas mais proativas e baseada em dados na gestão de desastres. A identificação de padrões ocultos e a previsão de eventos extraordinários, como cisnes negros e rinocerontes cinzas, podem ser cruciais para desenvolver estratégias de resposta mais robustas e abrangentes.

### REFERÊNCIAS

ALEXANDER, David. **Confronting Catastrophe: New Perspectives on Natural Disasters**. UK: Oxford University Press, 2000.

ANDERSON, T.W. **An Introduction to Multivariate Statistical Analysis**. California: [S.n.], v. III, 2003.

BORG, I; GROENEN, P. J. **Modern multidimensional scaling: Theory and applications**. [S.l.]: Springer, 2005.

BORGES, Vinicius R. P. *et al.* Using Principal Component Analysis to support students' performance prediction and data analysis. **VII Congresso Brasileiro de Informática na Educação (CBIE 2018)**, Brasília, 2018.

BRASIL. **Portaria MDR n. 260, de 2 de fevereiro de 2022**. Diário Oficial da República Federativa do Brasil. Brasília. 2022.

CHMUTINA, Ksenia; BOSHER, Lee. **Disaster Risk Reduction for the Built Environment**. [S.l.]: [S.n.], 2017.

ENAP. Análise fatorial. **Livro da Fundação Escola Nacional de Administração Pública**, Brasília, 2019.

FÁVERO, Luiz Paulo Lopes; BELFIORE, Patrícia Prado. **Manual de análise de dados: estatística e modelagem multivariada com excel, SPSS e stata**. Rio de Janeiro: Elsevier, 2017.

FAYYAD, Usama ; SHAPIRO, Gregory Piatetsky; SMYTH, Padhraic. From Data Mining to Knowledge Discovery in Databases. **AI Magazine. American Association for Artificial Intelligence**, v. 17, 1996.

FERNANDES, Samir B.; SANTOS, Alexandre B. D.; SILVA, Rodrigo W. d. Descoberta de conhecimento em banco de dados aplicados em proteção e defesa civil: uma análise dos registros de socorros contribuindo para a identificação de desastres. **Revista núcleo do conhecimento**, 27 julho 2022.

HAIR JR., J.F. *et al.* **Multivariate Data Analysis**. 7<sup>a</sup>. ed. [S.l.]: Prentice Hall, 2019.

HAN, Jiawei; KAMBER , Micheline; PEI, Jian. **Data Mining: Concepts and Techniques**. 3. ed. [S.l.]: Morgan Kaufmann Publishers, 2011.

HU, Li-tze; BENTLER, Peter M. Cutoff criteria for fit indexes in covariance structure analysis: Conventional criteria versus new alternatives. **Structural Equation Modeling: A Multidisciplinary Journal**, 1999.

INDHUMATHI, R; SATHIYABAMA, S. Reducing and Clustering high Dimensional Data through Principal Component Analysis. **International Journal of Computer Applications** , World, Dezembro 2010.

JOHNSON, Richard A. ; WICHERN, Dean W. **Applied Multivariate Statistical Analysis**. [S.I.]: Pearson, 2007.

JOLLIFFE, I. T. **Principal Component Analysis**. World: Springer, 2002.

KING, Ronald S. **Cluster Analysis and Data Mining: An Introduction**. [S.I.]: [S.n.], 2014.

LITTLE, Roderick J. A. ; RUBIN, Donald. **Statistical Analysis with Missing Data**. 2ª. ed. [S.I.]: John Wiley & Sons, 2002.

LOVRIC, Miodrag ; ARSHAM, Houssein. **Bartlett's Test**. [S.I.]: [S.n.], 2011.

RENCHE, A. C.; CHRISTENSEN, W. F.. **Methods of Multivariate Analysis**. [S.I.]: ohn Wiley & Sons., 2012.

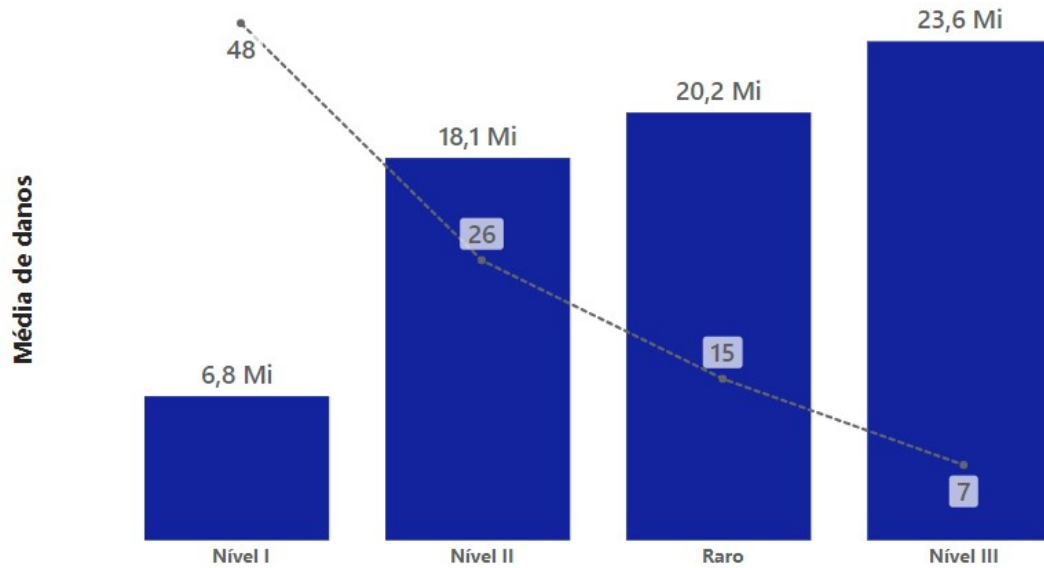
SOUSA, Hélio Marcus Damasceno; SANTOS, Daiana Ferreira de Deus; SOUZA, Gabriel Vieira Damasceno de. Disaster research: challenges and contributions to society. **Journal of Environmental Analysis and Progress**, 2020. 57-66.

TABACHNICK, B. G.; FIDELL, L. S. **Using Multivariate Statistics**. Boston: Pearson, 2013.

**APÊNDICE A**

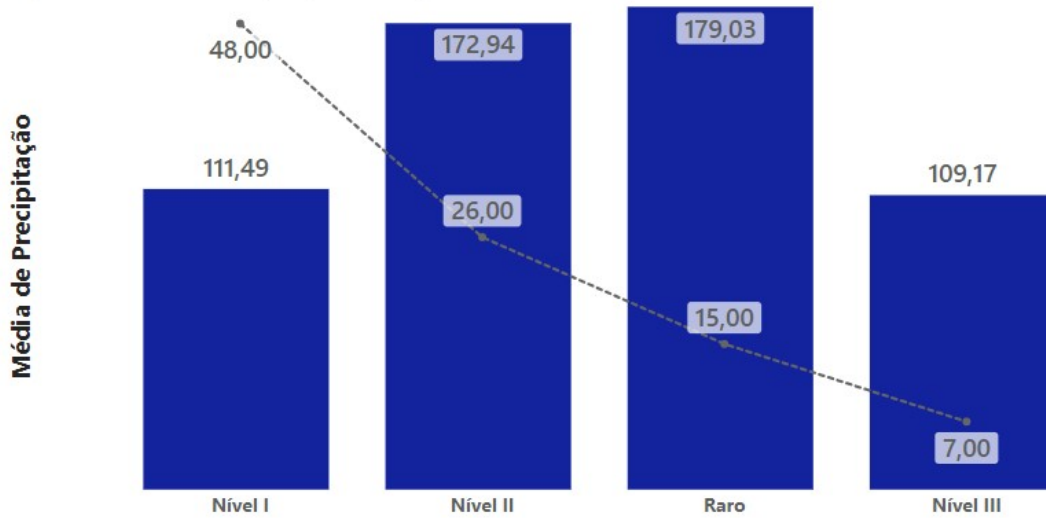
**Média de Danos e Prejuízos (em Milhões de R\$) por agrupamento**

Legenda ● Média de danos ● Contagem de desastres



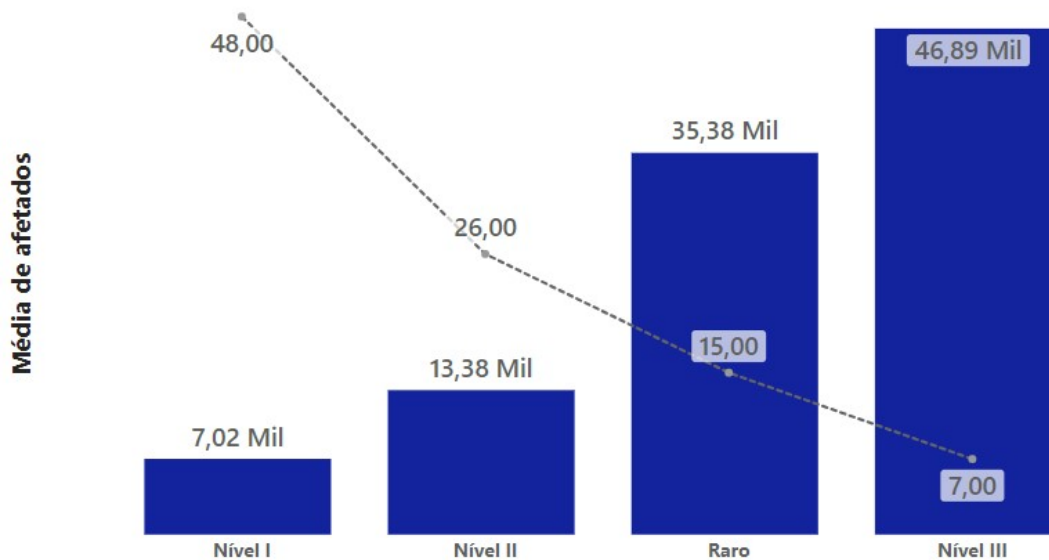
### Média de Precipitação pluviométrica (em mm/24) horas por agrupamentos

Legenda: ● Média de Precipitação ● Contagem de desastres



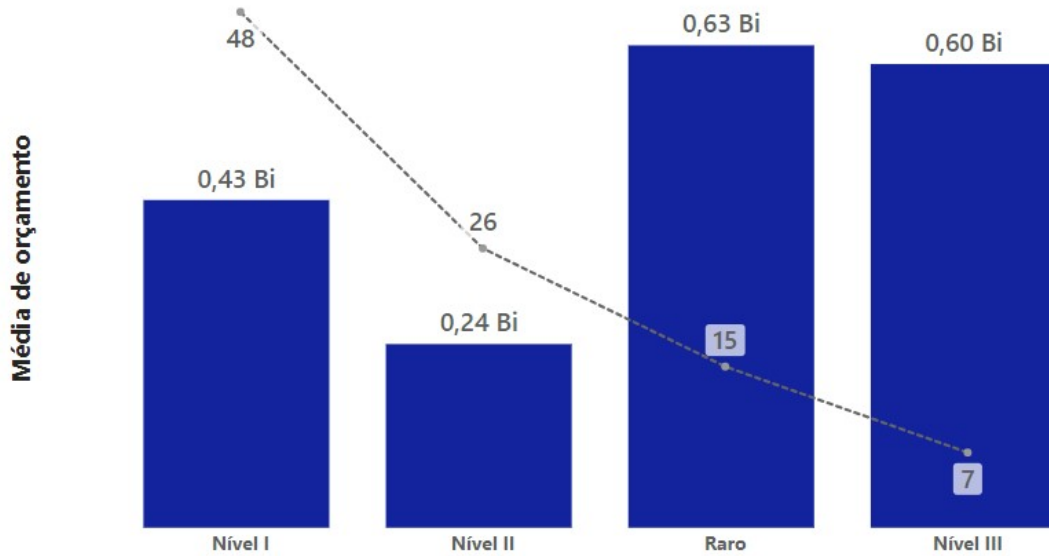
### Média de afetados (em milhares) por agrupamentos

Legenda: ● Média de afetados ● Contagem de desastres



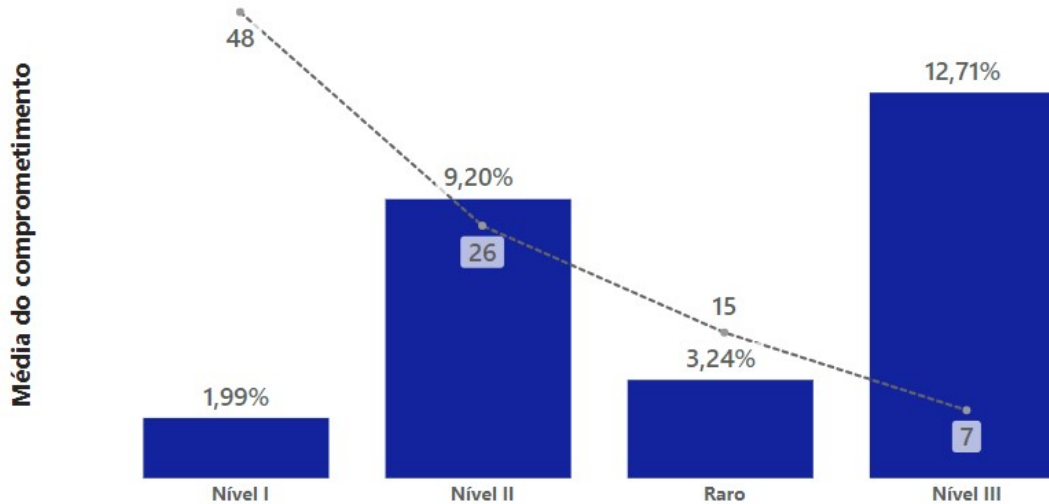
### Média de Receita municipal (em milhões de R\$) por agrupamento

Legenda: ● Média de orçamento ● Contagem de desastres



### Média do comprometimento da capacidade municipal de Investimento na resposta ao desastre por agrupamentos

Legenda: ● Média do comprometimento ● Contagem de desastres



### Média da densidade populacional afetada no desastre por agrupamentos

Legenda: ● Média de densidade populacional ● Contagem de desastres

