



Campus São Mateus
UNIVERSIDADE FEDERAL DO ESPÍRITO SANTO



Pró-Reitoria de Pesquisa e Pós-Graduação
Universidade Federal do Espírito Santo

ARTIGO ORIGINAL

OPEN ACCESS

CLASSIFICAÇÃO DA PERCEÇÃO DE SERVIDORES PÚBLICOS FEDERAIS EM RELAÇÃO A ATOS DE CORRUPÇÃO UTILIZANDO ALGORITMOS DE APRENDIZADO DE MÁQUINA

PERCEPTION CLASSIFICATION OF FEDERAL PUBLIC SERVANTS REGARDING ACTS OF CORRUPTION USING MACHINE LEARNING ALGORITHMS

CLASIFICACIÓN DE LA PERCEPCIÓN DE LOS SERVIDORES PÚBLICOS FEDERALES EN RELACIÓN A ACTOS DE CORRUPCIÓN UTILIZANDO ALGORITMOS DE APRENDIZAJE DE MÁQUINA

Vinicius Matheus Pimentel Ariza ^{1*} & Bruno Samways dos Santos ²

^{1,2} Universidade Tecnológica Federal do Paraná - UTFPR

^{1*} vinciusariza@alunos.utfpr.edu.br ² bruno.santos@utfpr.edu.br

ARTIGO INFO.

Recebido: 12.08.2023

Aprovado: 17.09.2023

Disponibilizado: 08.11.2023

PALAVRAS-CHAVE: Mineração de Dados; Corrupção; Serviço Público Federal.

KEYWORDS: Data Mining; Corruption; Federal Public Service.

PALABRAS CLAVE: Minería de Datos; Corrupción; Servicio Público Federal.

*Autor Correspondente: Ariza, V. M. P.

RESUMO

Técnicas computacionais têm-se mostrado úteis na luta contra a corrupção no setor público, permitindo a detecção precoce de atividades suspeitas. Sob este pressuposto, o objetivo deste trabalho foi comparar algoritmos de aprendizado de máquina no contexto da observação de atos de corrupção no Serviço Público. Nesse sentido, foram analisados dados extraídos de uma pesquisa realizada pelo Banco Mundial em 2021 sobre o tema “Ética e Corrupção no Serviço Público”, com cerca de 22.000 respondentes, sendo proposto o desenvolvimento de modelos que auxiliem na promoção da transparência e da integridade no serviço público brasileiro. Os resultados mostraram a viabilidade do uso de técnicas de aprendizado de máquina, com a Regressão Logística se mostrando a melhor opção para o cenário estudado, com acurácia de 82%. O modelo desenvolvido e as análises geradas podem ser usados para auxiliar na identificação de atividades suspeitas de corrupção no setor público, contribuindo para a detecção precoce e a prevenção de práticas ilegais. Os resultados também destacam a importância do desenvolvimento de políticas públicas para promover a ética e a integridade no serviço público, bem como o papel das tecnologias avançadas na melhoria da governança e da confiança da sociedade nas instituições públicas.

ABSTRACT

Computational techniques have proven useful in the fight against corruption in the public sector, enabling the early detection of suspicious activities. The aim of this study was to compare machine learning algorithms in the context of

observing acts of corruption in the Public Service. In this regard, data extracted from a survey conducted by the World Bank in 2021 on the topic of Ethics and Corruption in the Public Service were analyzed, involving approximately 22,000 respondents. The development of models aimed at promoting transparency and integrity in the Brazilian public service is proposed. The results demonstrated the feasibility of using machine learning techniques, with Logistic Regression proving to be the best option for the studied scenario, with an accuracy of 82%. The developed model and generated analysis can be used to assist in the identification of suspicious corruption activities in the public sector, contributing to early detection and prevention of illegal practices. The results also highlight the importance of developing public policies to promote ethics and integrity in public service, as well as the role of advanced technologies in improving governance and society's trust in public institutions.

RESUMEN

Las técnicas computacionales se han mostrado útiles en la lucha contra la corrupción en el sector público, permitiendo la detección temprana de actividades sospechosas. En este sentido, se analizaron datos extraídos de una encuesta realizada por el Banco Mundial en 2021 sobre el tema de Ética y Corrupción en el Servicio Público, con la participación de aproximadamente 22,000 encuestados. Se propone el desarrollo de modelos destinados a promover la transparencia y la integridad en el servicio público brasileño. Los resultados mostraron la viabilidad del uso de técnicas de aprendizaje de máquina, siendo la Regresión Logística la mejor opción para el escenario estudiado, con una precisión del 82%. El modelo desarrollado y los análisis generados pueden ser utilizados para ayudar en la identificación de actividades sospechosas de corrupción en el sector público, contribuyendo a la detección temprana y prevención de prácticas ilegales. Los resultados también resaltan la importancia del desarrollo de políticas públicas para promover la ética y la integridad en el servicio público, así como el papel de las tecnologías avanzadas en la mejora de la gobernanza y la confianza de la sociedad en las instituciones públicas.



1 INTRODUÇÃO

A corrupção e suas implicações na gestão pública têm ganhado destaque em vários meios de comunicação, e, especificamente no Brasil, o assunto alcançou uma elevada proporção, em virtude de frequentes episódios de denúncias envolvendo integrantes dos três poderes do Estado (Macedo & Valadares, 2021). Tópico recorrente nas redes sociais e no jornalismo político, as investigações sobre atos de corrupção em setores públicos, por governantes e servidores, destacam o país no cenário local e internacional. Ilustrando este cenário, de acordo com Gehrke et al. (2017), analisando-se os artigos publicados sobre o Brasil em revistas internacionais (Der Spiegel, L'Obs, The Economist e Time) no período de 2003 a 2014, ao menos 7,3% abordavam tópicos ligados à corrupção, sendo as publicações crescentes ao longo dos anos. Além disso, segundo a organização não governamental Transparência Internacional, o Brasil ocupa a 94ª posição em um *ranking* de 180 países no Índice de Percepção da Corrupção de 2022, o que indica que o problema ainda é muito presente no país.

No entanto, a corrupção é um problema de difícil solução, pois envolve diversos atores e variáveis (Jancsics, 2019), sendo um fenômeno que afeta indivíduos, indústrias, organizações e governos (Chen & Liao, 2011). Ademais, muitas vezes, é difícil detectar e comprovar casos de corrupção, o que torna o combate a esse problema ainda mais desafiador. Em contrapartida, atualmente, estudos indicam que a Tecnologia da Informação e Comunicação (TIC) podem promover a transparência, responsabilidade e participação dos cidadãos no processo anticorrupção, a partir de ferramentas baseadas em análises de *big data* e inteligência artificial (IA) (Adam & Fazekas, 2021).

A partir disso, o uso de técnicas de IA e Aprendizado de Máquina (AM) em detecção de fraude e corrupção governamental tem ganhado popularidade em anos recentes, com aplicações como a de Lima e Delen (2020), que aplicaram Random Forest, Support Vector Machine (SVM) e Redes Neurais Artificiais (RNA) para encontrar os preditores mais importantes para o Índice de Corrupção Percebida (Corruption Perception Index - CPI) em diversos países. De modo semelhante, este índice foi explorado por Domashova e Politova (2021), que aplicaram técnicas de clusterização e classificação para identificar os sinais e causas principais no CPI. Cita-se, ainda, a pesquisa de Li et al. (2020), que aplicaram métodos de Processamento de Linguagem Natural (PLN) e técnicas de aprendizagem não supervisionada para detectar autodeclaração de experiências com corrupção no Twitter, enquanto o trabalho publicado por De Blasio et al. (2022), em que aplicaram algoritmos para prever crimes de corrupção em municípios italianos com dados de 2011. Percebe-se, portanto, que existem diferentes formas de se explorar dados sobre corrupção, mas poucos trabalhos aplicaram algoritmos de classificação neste contexto, com nenhum deles explorando diretamente agentes públicos.

Preocupado com essa questão e seus impactos, em 2021, em uma parceria da Controladoria-Geral da União (CGU), Ministério da Economia e a Escola Nacional de Administração Pública (ENAP), o Banco Mundial desenvolveu a Pesquisa sobre Ética e Corrupção no Serviço Público (2021), contando com a participação de cerca de 22 mil respondentes. O conjunto de dados obtido na referida pesquisa foi utilizado para a aplicação de técnicas de mineração de dados,



utilizando-se algoritmos de AM para classificar os servidores públicos federais que presenciaram ou não atos de corrupção nos últimos três anos de seu exercício profissional.

Ante o exposto, a fim de investigar o tema e propor soluções, neste artigo, foram aplicados algoritmos de AM na classificação da avaliação de servidores públicos federais em relação aos atos de corrupção, além de identificar os principais fatores que contribuem para a percepção de tais atos. Para tanto, foram utilizadas quatro técnicas: Random Forest, K-Nearest Neighbors (KNN), RNA e Regressão Logística.

Desse modo, além desta introdução, o presente artigo está estruturado com mais quatro seções. A segunda seção apresenta uma revisão teórica das técnicas de classificação utilizadas, destacando sua aplicação no contexto da pesquisa, seguida das métricas de avaliação utilizadas. Na terceira seção, são descritos os materiais e métodos utilizados, abrangendo a descrição do conjunto de dados, das ferramentas e das etapas da pesquisa. Na quarta seção, são discutidos os resultados obtidos com a aplicação das técnicas de classificação nos dados da pesquisa. Por fim, a quinta e última seção traz a conclusão sobre a modelagem computacional e as técnicas aplicadas no contexto de percepção de atividades de corrupção por parte dos servidores públicos federais, ressaltando suas implicações para a promoção da transparência e da integridade no serviço público.

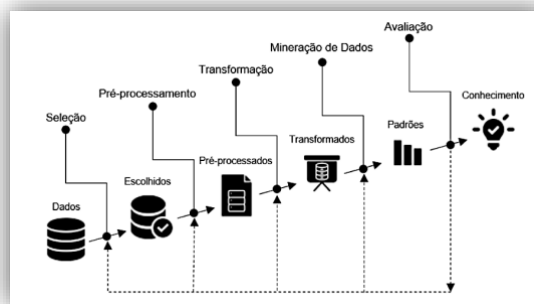
2 METODOLOGIA

Esta seção descreve alguns fundamentos sobre a metodologia de mineração de dados denominada Knowledge Discovery in Databases (KDD), as técnicas implementadas e as métricas para avaliação.

2.1 KNOWLEDGE DISCOVERY IN DATABASES

Com o avanço tecnológico, a quantidade de dados gerados diariamente torna-se cada vez maior, mostrando-se essencial o uso de recursos computacionais para processar essas informações. Dentre esses, o processo KDD, que é utilizado para descoberta de conhecimento em bases de dados, permite que algoritmos compreendam comportamentos e padrões a partir dos dados (Gamarra et al., 2016). Para isso, são necessárias algumas etapas, como a remoção de dados irrelevantes ou incompletos, seleção dos dados relevantes, transformação dos dados no formato adequado para a mineração, aplicação das técnicas de mineração de dados, avaliação dos padrões encontrados, apresentação e assimilação do conhecimento adquirido, definindo vantagens e desvantagens (Fayyad et al., 1996) (Figura 1).

Figura 1. Ilustração das etapas do processo KDD.



Fonte: Fayyad et al. (1996), adaptado pelos autores.



Dentro do processo KDD, para o desenvolvimento de modelos preditivos, pode-se utilizar a aprendizagem de máquina. A aprendizagem de máquina pode ser dividida em dois principais tipos: aprendizagem supervisionada e não supervisionada.

Na aprendizagem supervisionada tem-se o objetivo de treinar modelos com base em conjuntos de dados onde se conhece os pares de entrada-saída. Esses modelos devem ser capazes de aprender uma função que, ao receber um conjunto de dados desconhecidos, mapeie uma saída próxima ou equivalente aos padrões dos pares utilizados no treinamento. Dentre as tarefas de aprendizagem supervisionada, destaca-se a classificação, na qual o modelo treinado atribui determinadas classes a novos dados desconhecidos (Aggarwal, 2015). Já na aprendizagem não supervisionada, os dados de treinamento não possuem um rótulo específico, não necessitando de uma variável de saída como referência para o algoritmo. Tarefas comuns para este caso são o agrupamento de registros e a formação de regras de associação (Géron, 2019). Destaca-se também que é uma área desafiadora, pois não há uma resposta correta para esses algoritmos uma vez que não é possível checar os resultados com rótulos pré-definidos (James et al., 2023).

No contexto deste trabalho, as técnicas de classificação serão empregadas para que a partir dos parâmetros selecionados na modelagem computacional, seja possível prever se o servidor público federal presenciou atos de corrupção durante suas atividades profissionais.

2.2 TÉCNICAS UTILIZADAS

Dentre as técnicas utilizadas no processo de AM, tem-se o KNN, o Random Forest, as RNAs e a Regressão Logística.

O método KNN é uma técnica de AM não paramétrica baseada em distância. Segundo Castro e Ferrari (2016), o KNN avalia os dados que estão à menor distância "k" de dados já classificados e a classificação do dado em questão é determinada pelos dados que tiverem maior representatividade nesta análise. O valor "k" no nome do método refere-se à quantidade de objetos que estão sendo considerados à menor distância do dado que se deseja classificar, ou seja, é o número de vizinhos que serão considerados (Pan & Pan, 2020). O KNN é amplamente utilizado em problemas de classificação de dados e pode ser uma escolha eficaz em problemas com conjuntos de dados pequenos ou médios.

O Random Forest, ou Floresta Aleatória, é uma técnica de AM que combina a predição de vários modelos de árvores de decisão (Speiser et al., 2019). Cada árvore é treinada com uma amostra aleatória do conjunto de dados de treinamento e com um subconjunto aleatório de atributos. O resultado é uma combinação das previsões de todas as árvores individuais (Géron, 2019). Essa técnica é útil para a classificação e regressão em dados complexos, além de ser amplamente utilizada em diversas áreas, como finanças, medicina e ciência ambiental (Breiman, 2021).

A RNA é uma técnica que utiliza ferramentas matemáticas para simular a estrutura neural do ser humano. Mitchell (1997) explica que os métodos de aprendizado baseados em redes neurais são “uma abordagem robusta para aproximar funções de destino com valores reais,



discretos e vetoriais”. A técnica é importante para a interpretação de dados complexos de sensores do mundo real e pode ser utilizada para diversas aplicações, como por exemplo, a classificação de respostas como deste estudo.

Por fim, a Regressão Logística é uma técnica estatística que busca, por meio de modelagens matemáticas, prever a probabilidade de um evento ocorrer com base na análise da relação existente entre variáveis. Conforme definido por Gonzalez (2018), a Regressão Logística é um modelo capaz de prever valores de uma variável categórica, frequentemente binária, a partir de uma ou mais variáveis independentes (Raschka, 2015).

Neste estudo, será utilizado o modelo de RNA Perceptron multicamadas para a classificação binária de dados históricos, ou seja, para dividir os dados em duas categorias distintas. O modelo de Perceptron multicamadas é um tipo de RNA que possui múltiplas camadas de neurônios interconectados e pode ser utilizado para problemas de classificação e regressão (Aggarwal, 2015).

2.3 MÉTRICAS DE AVALIAÇÃO

Para avaliar os modelos gerados, foram adotados três métricas: acurácia, precisão e recall. A acurácia avalia a capacidade do modelo em classificar corretamente as instâncias, enquanto a precisão mede a proporção de classificações corretas para a classe positiva. Já o recall indica a proporção de instâncias da classe de interesse corretamente classificadas pelo modelo. As equações apresentadas (1), (2) e (3) correspondem às definições matemáticas de cada uma das métricas utilizadas neste trabalho.

$$Acurácia = \frac{VP + VN}{VP + VN + FP + FN} \quad (1)$$

$$Precisão = \frac{VP + VN}{VP + FP} \quad (2)$$

$$Recall = \frac{VP}{VP + FN} \quad (3)$$

Nas equações, os termos indicados nas expressões correspondem aos valores verdadeiros positivos (VP), verdadeiros negativos (VN), falsos positivos (FP) e falsos negativos (FN).

Além das métricas mencionadas anteriormente, também foi utilizada a curva Receiver Operating Characteristic (ROC) para avaliar o desempenho dos modelos. A curva ROC é uma representação gráfica da relação entre a taxa de verdadeiros positivos e a taxa de falsos positivos para diferentes valores de limiar de classificação. Essa curva permite avaliar a sensibilidade e especificidade do modelo em diferentes pontos de corte. Em geral, quanto mais próxima à curva ROC estiver do canto superior esquerdo do gráfico, melhor será o desempenho do modelo. A área sob a curva ROC (AUC) também é uma medida comumente utilizada para comparar o desempenho de diferentes modelos (Fawcett, 2006).

2.4 DESCRIÇÃO DO CONJUNTO DE DADOS

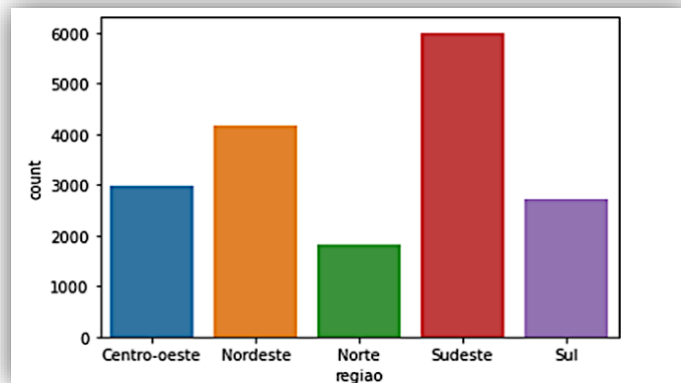
O conjunto de dados original (extraído diretamente do site da pesquisa, no endereço: <https://microdata.worldbank.org/index.php/catalog/4300>) possui originalmente 21.356 instâncias e 55 atributos. Todos os dados utilizados da pesquisa possuem formato de texto,



isto é, sem respostas numéricas. Dos atributos iniciais contidos na base, 15 foram removidos, pois não eram interessantes para o modelo ou estavam sem identificação na base de dados. Além disso, dois atributos foram excluídos, pois ofereciam possíveis vieses aos modelos, visto que perguntavam de forma explícita se os respondentes já denunciaram atos de corrupção e qual a sua repercussão.

Adicionalmente, algumas informações dos respondentes foram obtidas na base de dados, o que permite compreender alguns pontos do perfil dos respondentes, tais como região, gênero e nível de formação. Essas informações estão representadas nas Figuras 2, 3 e 4.

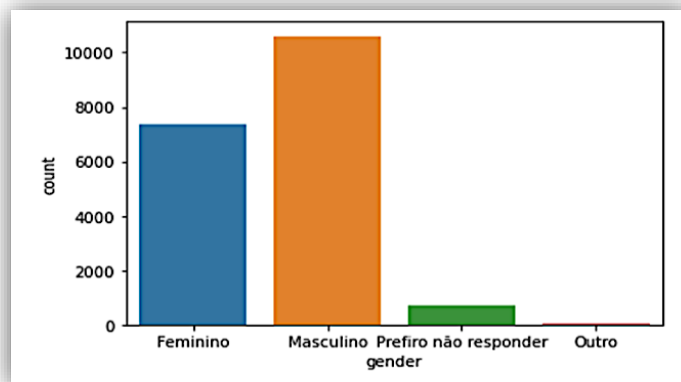
Figura 2. Regiões do país em que atuam os respondentes da pesquisa.



Fonte: Autores (2023).

Observa-se na Figura 2 uma distribuição homogênea de respondentes, sendo que todas as porções do território nacional são representadas na pesquisa de forma considerável.

Figura 3. Quantidade de respondentes da pesquisa por gênero.

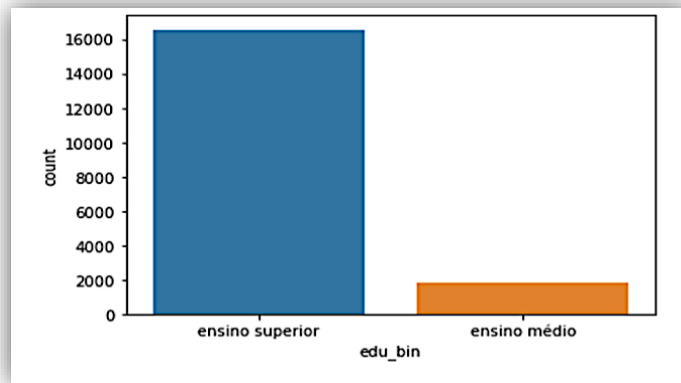


Fonte: Autores (2023).

Com relação ao gênero dos respondentes da pesquisa, ilustrado na Figura 3, observa-se que existe um balanceamento entre os respondentes que se identificam com o gênero masculino e feminino, embora em maior representatividade os respondentes do gênero masculino, nota-se uma amostra considerável de ambos os gêneros.



Figura 4. Quantidade de respondentes da pesquisa por nível de escolaridade.



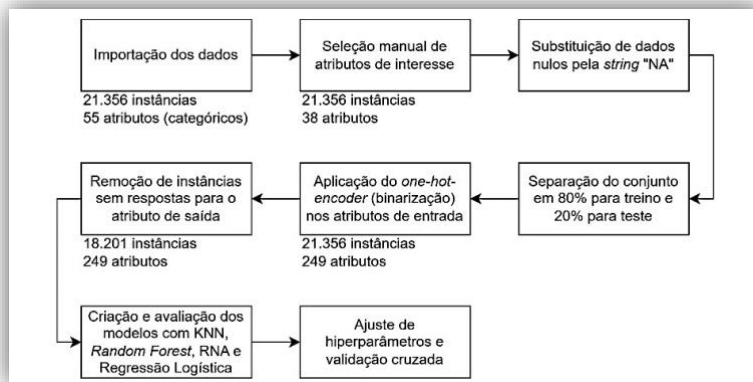
Fonte: Autores (2023).

Com relação à escolaridade dos respondentes, observa-se na Figura 4 que, de forma expressiva, a maior parte dos participantes da pesquisa possui ensino superior completo.

2.5 ETAPAS DA PESQUISA

No processo de modelagem computacional, para que os padrões sejam identificados, e a classificação desejada a respeito dos servidores públicos federais que presenciaram atos de corrupção nos últimos três anos fosse possível, algumas etapas foram seguidas (Figura 5).

Figura 5. Ilustração das etapas aplicadas nesta pesquisa



Fonte: Autores (2023).

Inicialmente, utilizando o Jupyter Notebook – ambiente de desenvolvimento que permite a criação de códigos de programação – foi realizada a importação de bibliotecas relevantes para a programação na linguagem Python. As bibliotecas importadas no processo foram as seguintes: *pandas 1.4.3* (biblioteca utilizada para estruturação, manipulação e refinamento dos dados, além de ser útil para o processamento de dados e construções gráficas), *numpy 1.23.0* (biblioteca que oferece um conjunto de funções e operações que facilitam cálculos numéricos), *matplotlib 3.5.2* (biblioteca útil para a criação de diversos tipos de gráficos) e *seaborn 0.11.2* (biblioteca adequada para análises estatísticas e construção de gráfico). Para a aplicação dos modelos e avaliação dos resultados foi utilizada a biblioteca *scikit-learn*.

Dentre as perguntas do questionário que gerou a base de dados estudada, destaca-se a do atributo nomeado como "ee_2", que questiona os servidores quais atos de corrupção foram

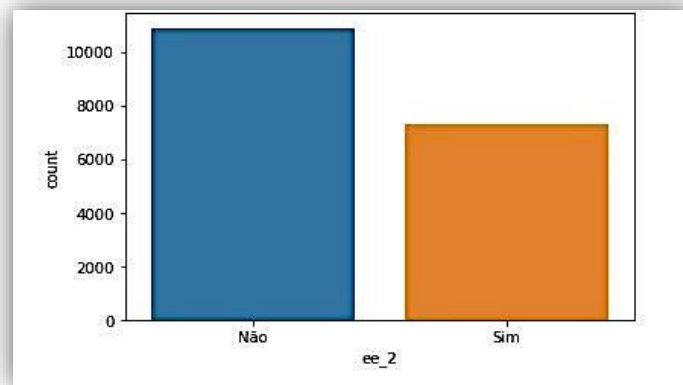


observados por eles nos últimos três anos, o que norteou a modelagem computacional deste estudo. Com a base de dados disponível, iniciou-se o processo de seleção de atributos, isto é, foi estabelecido um filtro de informações com o intuito de manter na coleção de dados, apenas as informações relevantes para o problema, eliminando atributos que poderiam afetar o algoritmo, levando a vieses indesejáveis. Desta forma, após este filtro, a base de dados permaneceu tendo as 21.356 linhas (instâncias) iniciais, porém, agora apenas 38 colunas.

Após a seleção dos atributos, iniciou-se o pré-processamento e transformação dos dados coletados. Nesta etapa, foram verificados dados nulos, substituindo essas informações por um dado de texto "NA". Também foi estabelecido o atributo alvo do modelo, no caso o atributo "ee_2", referente aos colaboradores que presenciaram outros servidores com posturas antiéticas ou em atos de corrupção durante o exercício de suas atividades. Nesse contexto foram levantadas as informações dos colaboradores que presenciaram atos antiéticos, dos que não presenciaram atos antiéticos, e eliminados os dados dos servidores que optaram por não responder à pergunta. Após estabelecimento deste parâmetro, como sendo o parâmetro de predição, a base de dados passou a ter 18.201 instâncias.

Verificou-se ainda que não seria necessário o balanceamento do atributo alvo, visto que existem amostras suficientes de ambas as classes ("SIM" e "NÃO"), para o treinamento do modelo. A Figura 6 ilustra a quantidade de respondentes que presenciaram ou não atos de corrupção.

Figura 6. Quantidade de respondentes que responderam ter presenciado atos de corrupção nos últimos três anos de atuação.



Fonte: Autores (2023).

Por fim, para tornar possível a aplicação de todos os modelos propostos, os atributos que permitiam mais de uma resposta foram separados em diferentes colunas, aplicando o método *one-hot-encoder*, o que caracteriza uma binarização, que também foi utilizado nos demais atributos, gerando uma base de dados final com 18.201 instâncias e 249 atributos.

Tendo-se os dados brutos pré-processados, foram aplicados os algoritmos Random Forest, KNN, RNA e Regressão Logística, avaliando diferentes hiperparâmetros, obtendo-se os melhores resultados de acurácia, precisão, *recall* e AUC para cada algoritmo. Para a aplicação das métricas, comparação dos algoritmos e validação cruzada, a base foi separada em treino (80%) e teste (20%).



3. RESULTADOS E DISCUSSÃO

Foi possível constatar que os algoritmos Random Forest, RNA e Regressão Logística obtiveram desempenhos similares, enquanto o KNN apresentou menor desempenho em termos de métricas de avaliação (Tabela 1).

Tabela 1. Resultado das métricas de avaliação dos modelos.

	Random Forest	KNN	RNA	Regressão Logística
Acurácia	0,81	0,79	0,82	0,82
Precisão	0,78	0,79	0,79	0,77
Recall	0,74	0,62	0,73	0,75
AUC	0,88	0,85	0,89	0,89

Fonte: Autores (2023).

Para alcançar os resultados da Tabela 1, os testes foram realizados com diversas combinações de parâmetros em cada algoritmo, caracterizando a etapa de investigação de hiperparâmetros. Optou-se por manter alguns parâmetros com a opção padrão, como por exemplo, a taxa de aprendizagem constante de 0,001 e 200 épocas para a RNA (Tabela 2).

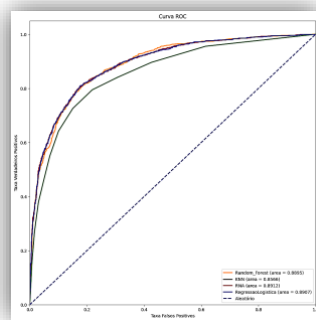
Tabela 2. Parâmetros selecionados após a avaliação de hiperparâmetros dos modelos que apresentaram os melhores resultados iniciais com testes empíricos.

Algoritmo	Parâmetro	Valores testados	Valor selecionado
RNA	activation	relu, logistic	logistic
	solver	adam, sgd	sgd
	batch_size	10, 20	20
Random Forest	criterion	gini, entropy	entropy
	n_estimators	10, 100, 1000	1000
Regressão Logística	solver	lbfgs, sag, saga	sag
	C	1.0, 1.5	1.0

Fonte: Autores (2023).

Com o objetivo de garantir a validade dos resultados obtidos, os três algoritmos com melhor desempenho foram submetidos a uma validação cruzada utilizando o método *k-fold* com 30 divisões cada. A Regressão Logística e RNA obtiveram uma acurácia média de aproximadamente 0,82, enquanto o algoritmo Random Forest apresentou uma acurácia média de 0,81. Adicionalmente, foi realizada uma análise das curvas ROC para os algoritmos utilizados no estudo, e os resultados indicaram que as curvas apresentaram comportamento semelhante, indicando que os algoritmos possuem desempenho equivalente em termos de classificação. Esses resultados reforçam a robustez dos resultados obtidos anteriormente por meio da validação cruzada. Assim, é possível afirmar que os três algoritmos apresentam um desempenho semelhante na tarefa de classificação empregada neste estudo (Figura 7).

Figura 7. Ilustração do gráfico da curva ROC.



Fonte: Autores (2023).



Como mencionado, foi possível constatar que os diferentes algoritmos apresentaram desempenhos similares em relação à tarefa de classificação. No entanto, devido à sua simplicidade e facilidade de interpretação, chegou-se à conclusão de que a utilização do algoritmo de Regressão Logística é a opção mais adequada. A partir da avaliação das métricas de desempenho e das características de cada algoritmo, foi possível identificar que a Regressão Logística apresenta o melhor equilíbrio entre simplicidade, eficiência e capacidade de generalização para a tarefa de classificação em questão. Portanto, a escolha do algoritmo de Regressão Logística se mostrou a melhor opção para o caso em questão.

A Regressão Logística é uma técnica estatística amplamente utilizada em estudos de classificação binária. Isso se deve à sua capacidade de lidar com variáveis dependentes dicotômicas, ou seja, aquelas que apresentam apenas duas categorias possíveis, como “Sim” ou “Não”. Além disso, essa técnica é conhecida por apresentar bons resultados em bases de dados com essas características, o que pode explicar o desempenho satisfatório obtido no estudo em questão. Em um estudo relacionado, Fernandes et al. (2021) afirmam que a Regressão Logística é a melhor ferramenta para lidar com variáveis dependentes dicotômicas, como “eleito” ou “não eleito”, “adotou a política” ou “não adotou”, entre outras. Essa afirmação se aplica diretamente ao estudo em análise, que teve como objetivo classificar servidores públicos em “Sim” ou “Não” para a percepção de atos de corrupção. Os resultados obtidos na análise dos algoritmos indicaram que a Regressão Logística foi capaz de realizar essa classificação de forma satisfatória.

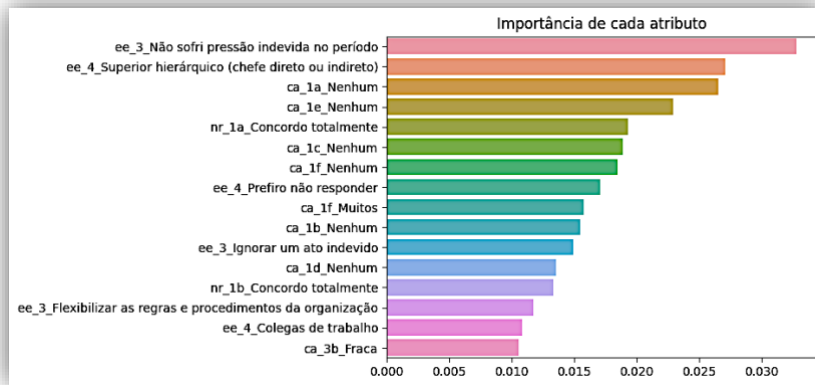
Os resultados encontrados neste estudo mostram que a utilização de modelos de aprendizado de máquina, como a Regressão Logística, pode ser uma ferramenta promissora para fiscalizar diversas situações no serviço público brasileiro. Desde que aplicados corretamente, esses modelos podem se adaptar a diferentes tipos de dados e obter resultados satisfatórios, como demonstrado neste estudo.

Este cenário sobre os modelos se torna ainda mais relevante quando se avalia a transparência de organizações públicas, que pode ser analisada por informações via *websites*. Estas informações são de pouco valor caso os cidadãos não entendam aquilo que foi publicado ou divulgado (Jeong et al., 2023). Assim, estes modelos podem ter uma contribuição prática para auxiliar na transparência no combate à corrupção. Também, a partir do conhecimento da literatura científica e do local onde se analisa, é possível construir mecanismos de prevenção aos riscos de corrupção, tanto na esfera pública quanto na esfera privada (Goutte et al., 2022).

Por fim, utilizou-se uma das funcionalidades do algoritmo Random Forest para identificar os atributos de maior importância para a tomada de decisão. Essa análise foi realizada por meio do gráfico de importâncias, apresentado na Figura 8, que permitiu uma avaliação mais detalhada e precisa dos atributos que influenciam de maneira mais significativa a classificação dos dados.



Figura 8. Ranking de importância dos atributos para a decisão do modelo.



Fonte: Autores (2023).

O atributo mais relevante para a classificação do modelo foi relacionado à pergunta que permitia ao respondente selecionar se sofreu algum tipo de pressão indevida nos últimos três anos. Em particular, a resposta “Não sofreu pressão indevida no período” foi identificada como de grande importância para o algoritmo, pois essa resposta indica que o servidor não foi exposto a situações que poderiam influenciar sua percepção sobre atos de corrupção. Essa pergunta se mostrou relevante para o modelo, uma vez que o comportamento ético dos servidores pode ser influenciado por pressões externas, e a ausência dessas pressões pode indicar uma menor probabilidade de observação de atos de corrupção.

O segundo atributo mais importante identificado pela análise do gráfico de importâncias do modelo estava relacionado à pergunta sobre qual agente exerceu pressão indevida sobre o servidor. A resposta mais relevante para o modelo foi “Superior hierárquico (chefe direto ou indireto)”. Essa informação é importante para o modelo porque pode indicar a presença de uma cultura organizacional que tolera ou até mesmo incentiva comportamentos antiéticos, além de indicar a necessidade de medidas de controle e fiscalização das ações de chefes e gestores dentro da instituição pública. Essa informação pode ser utilizada para desenvolver estratégias de prevenção e combate à corrupção no ambiente de trabalho.

Os dois primeiros atributos podem estar diretamente relacionados com o que Jackson e Köbis (2018) comentam sobre “pressões normativas”, fazendo parte da pressão vertical e pressão horizontal. A primeira reside no fato de existir uma pressão vinda “de cima”, ou seja, de pessoas de níveis organizacionais mais altos. Já a segunda diz respeito ao que os colegas de trabalho fazem, pois a presença de atos de corrupção pode também existir entre grupos de trabalho e a pressão pode advir de outros membros, não sendo necessariamente restrita aos superiores.

Os dois atributos mais relevantes após os dois primeiros foram identificados na pergunta que indagava a opinião dos servidores sobre o número de agentes públicos em sua organização que praticavam certas condutas antiéticas. O primeiro deles estava relacionado à opção “Aceitar dinheiro ou presentes de particulares para cumprir suas funções”, onde a resposta mais relevante para o modelo foi “Nenhum”. O segundo atributo estava relacionado à opção “Contratar uma empresa porque mantém vínculo remunerado ou de amizade nessa



empresa”, tendo também a resposta “Nenhum” como a mais relevante para a classificação do modelo.

Esses resultados sugerem que a percepção dos servidores sobre a prática de condutas antiéticas por outros agentes públicos na organização pode ter um grande impacto na classificação de um possível ato de corrupção. A resposta “Nenhum” indica que os servidores percebem que essas condutas não são comuns ou aceitas dentro da organização, o que pode indicar um ambiente de trabalho ético e saudável. Por outro lado, se houver uma resposta diferente de “Nenhum”, pode ser um indicativo de que a organização possui problemas éticos e, portanto, pode ser mais propensa à corrupção, o que ressalta a importância de se promover uma cultura ética nas organizações públicas e de se implementar medidas para prevenir e combater a corrupção.

Seguindo nesta mesma linha de pensamento, a literatura mostra que há várias raízes da corrupção, de diferentes formas, que podem recair exatamente no caso da aceitação de dinheiro ou beneficiamento de contratos de empresas. Por exemplo, o nepotismo, suborno e desvio de ativos são métodos de benefícios privados, incluídos nas práticas corruptas em níveis institucional, organizacional e individual (Ashforth et al., 2008), e, no Brasil, não são práticas isoladas ou raras de acontecer (Lino et al., 2022).

4 CONSIDERAÇÕES FINAIS

O estudo propôs o desenvolvimento de modelos de AM para classificar servidores públicos que presenciaram ou não atos de corrupção nos últimos três anos (entre abril-maio 2018 e abril-maio 2021). Foram utilizadas técnicas como KNN, Random Forest, RNA e Regressão Logística, e métricas como acurácia, precisão, *recall* e a curva ROC para avaliar os modelos. A base de dados foi pré-processada e transformada, sendo eliminados dados nulos e selecionados atributos relevantes.

A Regressão Logística mostrou-se a melhor opção, apresentando um bom equilíbrio entre simplicidade, eficiência e capacidade de generalização. O uso dessa abordagem pode auxiliar na identificação de atividades suspeitas de corrupção no setor público, contribuindo para a detecção precoce e a prevenção de práticas ilegais.

Em suma, os resultados obtidos neste estudo demonstram a viabilidade do uso de técnicas de aprendizado de máquina na classificação de servidores públicos que presenciaram ou não atos de corrupção. Esses resultados destacam a importância do uso de tecnologias avançadas para promover a transparência e a integridade no serviço público, o que pode levar a melhorias significativas na governança e na confiança da sociedade nas instituições.

No entanto, vale ressaltar que os resultados obtidos refletem o contexto e as percepções do período em que foi aplicada a pesquisa e que o uso desses modelos deve ser complementado por outras estratégias e abordagens no combate à corrupção no serviço público. Pesquisas futuras podem concentrar-se em combinar outras técnicas, como PLN e modelos de classificação a fim de analisar outras variáveis contextuais.

REFERÊNCIAS

Adam, I. & Fazekas, M. (2021). Are emerging technologies helping win the fight against corruption? A review of the state of evidence.

Information Economics and Policy, 57.

<https://doi.org/10.1016/j.infoecopol.2021.100950>

Aggarwal, C. C. (2015). *Data mining: the textbook* (Vol. 1). New York: springer.



- Breiman, L. (2001). Random forests. *Machine learning*, 45. <https://doi.org/10.1023/A:1010933404324>
- Ashforth, B. E., Gioia, D. A., Robinson, S. L., & Trevino, L. K. (2008). Re-viewing organizational corruption. *Academy of Management review*, 33(3). <https://doi.org/10.5465/amr.2008.32465714>
- Castro, L. N. D., & Ferrari, D. G. (2016). *Introdução à mineração de dados: conceitos básicos, algoritmos e aplicações*. São Paulo: Saraiva, 5.
- Chen, S. H. & Liao, C. C. (2011). Are foreign banks more profitable than domestic banks? Home-and host-country effects of banking market structure, governance, and supervision. *Journal of Banking & Finance*, 35(4). <https://doi.org/10.1016/j.jbusres.2022.03.032>
- Fawcett, T. (2006). An introduction to ROC analysis. *Pattern recognition letters*, 27(8). <https://doi.org/10.1016/j.patrec.2005.10.010>
- Fernandes, A. A. T., Figueiredo Filho, D. B., Rocha, E. C. D., & Nascimento, W. D. S. (2021). Leia este artigo se você quiser aprender regressão logística. *Revista de Sociologia e Política*, 28. <https://doi.org/10.1590/1678-987320287406en>
- Gonzalez, L. D. A. (2018). *Regressão logística e suas aplicações*. Recuperado de <https://monografias.ufma.br/jsui/bitstream/123456789/3572/1/LEANDRO-GONZALEZ.pdf>
- de Blasio, G., D'Ignazio, A., & Letta, M. (2022). Gotham city. Predicting 'corrupted' municipalities with machine learning. *Technological Forecasting and Social Change*, 184. <https://doi.org/10.1016/j.techfore.2022.122016>
- Domashova, J. & Politova, A. (2021). The Corruption Perception Index: analysis of dependence on socio-economic indicators. *Procedia Computer Science*, 190. <https://doi.org/10.1016/j.procs.2021.06.024>
- Fayyad, U., Piatetsky-Shapiro, G., & Smyth, P. (1996). From data mining to knowledge discovery in databases. *AI magazine*, 17(3). <https://doi.org/10.1609/aimag.v17i3.1230>
- Gamarra, C., Guerrero, J. M., & Montero, E. (2016). A knowledge discovery in databases approach for industrial microgrid planning. *Renewable and Sustainable Energy Reviews*, 60. <https://doi.org/10.1016/j.rser.2016.01.091>
- Gehrke, G., Borba, J. A., & Ferreira, D. D. M. (2017). A repercussão da corrupção brasileira na mídia: uma análise comparada das revistas Der Spiegel, L'Obs, The Economist, Time e Veja. *Revista de Administração Pública*, 5. <http://dx.doi.org/10.1590/0034-7612158681>
- Géron, A. (2019). *Mãos à Obra: Aprendizado de Máquina com Scikit-Learn & TensorFlow*. Alta Books.
- Goutte, S., Péran, T., & Porcher, T. (2022). Corruption, economy and governance in Central Africa: An analysis of public and regional drivers of corruption. *Finance Research Letters*, 44. <https://dx.doi.org/10.2139/ssrn.3808716>
- Jackson, D. & Köbis, N. (2018). *Anti-corruption through a social norms lens*. U4 Issue, 7. Recuperado de <https://www.u4.no/publications/anti-corruption-through-a-social-norms-lens#conclusion-a-social-norms-approach-to-anti-corruption>
- James, G., Witten, D., Hastie, T., Tibshirani, R., & Taylor, J. (2023). *An Introduction to Statistical Learning with Applications in Python*. Springer.
- Jancsics, D. (2019). Corruption as resource transfer: An interdisciplinary synthesis. *Public Administration Review*, 79(4). <https://doi.org/10.1111/puar.13024>
- Jeong, D., Shenoy, A., & Zimmermann, L. V. (2023). De Jure versus De Facto transparency: Corruption in local public office in India. *Journal of Public Economics*, 221. <https://doi.org/10.1016/j.jpubeco.2023.104855>
- Li, J., Chen, W. H., Xu, Q., Shah, N., Kohler, J. C., & Mackey, T. K. (2020). Detection of self-reported experiences with corruption on twitter using unsupervised machine learning. *Social Sciences & Humanities Open*, 2(1). <https://doi.org/10.1016/j.ssaho.2020.100060>
- Lima, M. S. M. & Delen, D. (2020). Predicting and explaining corruption across countries: A machine learning approach. *Government Information Quarterly*, 37(1). <https://doi.org/10.1016/j.giq.2019.101407>
- Lino, A. F., Azevedo, R. R., de., Aquino, A. C. B., de., & Steccolini, I. (2022). Fighting or supporting corruption? The role of public sector audit organizations in Brazil. *Critical Perspectives on Accounting*, 83. <https://doi.org/10.1016/j.cpa.2021.102384>
- Macedo, S. V. & Valadares, J. L. (2021). Corrupção: reflexões epistemológicas e contribuições para o campo de públicas. *Organizações & Sociedade*, 28. <https://doi.org/10.1590/1984-92302021v28n9607PT>
- Mitchell, T. M. (1997). *Machine learning*. (Vol. 1). New York: McGraw-hill.
- Pan, Z., Wang, Y., & Pan, Y. (2020). A new locally adaptive k-nearest neighbor algorithm based on discrimination class. *Knowledge-Based Systems*, 204. <https://doi.org/10.1016/j.knosys.2020.106185>
- Raschka, S. (2015). *Python machine learning*. Packt publishing ltd.
- Speiser, J. L., Miller, M. E., Tooze, J., & Ip, E. (2019). A comparison of random forest variable selection methods for classification prediction modeling. *Expert systems with applications*, 134. <https://doi.org/10.1016/j.eswa.2019.05.028>
- Transparência Internacional. (2022). *Índice de Percepção da Corrupção 2022*. Recuperado de <https://www.transparency.org/en/cpi/2022/index/brazil>

