

Research Article

Machine Learning for Predicting Distant Metastasis of Medullary Thyroid Carcinoma Using the SEER Database

Zhen-Tian Guo,¹ Kun Tian,¹ Xi-Yuan Xie,² Yu-Hang Zhang,³ and De-Bao Fang ⁶

¹Department of General Surgery, Beijing Electric Power Hospital, State Grid Corporation China, Capital Medical University, Beijing 100073, China

²Fujian Provincial Hospital, Fuzhou, Fujian 350001, China

³Mudanjiang Medical University, Mudanjiang, Heilongjiang 157000, China

⁴Hefei National Laboratory for Physical Sciences at Microscale, School of Basic Medical Sciences,

Division of Life Sciences and Medicine, University of Science and Technology of China, Hefei 230027, Anhui, China

Correspondence should be addressed to De-Bao Fang; debao@ustc.edu.cn

Received 24 May 2023; Revised 19 December 2023; Accepted 21 December 2023; Published 30 December 2023

Academic Editor: Alexander Schreiber

Copyright © 2023 Zhen-Tian Guo et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Objectives. We aimed to establish an effective machine learning (ML) model for predicting the risk of distant metastasis (DM) in medullary thyroid carcinoma (MTC). *Methods.* Demographic data of MTC patients were extracted from the Surveillance, Epidemiology, and End Results (SEER) database of the National Institutes of Health between 2004 and 2015 to develop six ML algorithm models. Models were evaluated based on accuracy, precision, recall rate, F1-score, and area under the receiver operating characteristic curve (AUC). The association between clinicopathological characteristics and target variables was interpreted. Analyses were performed using traditional logistic regression (LR). *Results.* In total, 2049 patients were included and 138 developed DM. Multivariable LR showed that age, sex, tumor size, extrathyroidal extension, and lymph node metastasis were predictive features for DM in MTC. Among the six ML models, the random forest (RF) had the best predictability in assessing the risk of DM in MTC, with an accuracy, precision, recall rate, F1-score, and AUC higher than those of the traditional binary LR model. *Conclusion.* RF was superior to traditional LR in predicting the risk of DM in MTC and can provide a valuable reference for clinicians in decision-making.

1. Introduction

As a result of changes in living environments, heightened health awareness, and advances in detection technology, the incidence of thyroid cancer has experienced a considerable increase in most parts of the world [1]. Medullary thyroid carcinoma (MTC) is a relatively rare malignancy, constituting approximately 5% of all thyroid malignancies. Patients with MTC generally exhibit a poorer prognosis than those with differentiated thyroid cancer (DTC), with MTC accounting for approximately 13% of all thyroid cancerrelated fatalities [2, 3]. Roughly 75% of MTC cases are sporadic, while around 25% are autosomal dominant [4]. Research has demonstrated that mutations in *RET*, a proto-oncogene, are present in approximately 6% of sporadic MTC patients and up to 98% of familial-inherited MTC patients [5]. Studies have indicated that extrathyroidal extension and distant metastasis (DM) are significant predictors of poor prognosis in patients [6, 7]. At the time of initial diagnosis, 10%–15% of MTC patients present with DM [8]. DM of MTC may involve the bones, lungs, and liver [9]. The American Thyroid Association's guidelines for the management of medullary thyroid cancer recommend various imaging examinations for MTC, potentially involving DM, including enhanced CT, MRI, abdominal ultrasound, and bone scans [10]. These diagnostic methods have a sensitivity of approximately 50%–80% for metastatic diseases. In recent years, the clinical application of drugs targeting *RET* proto-oncogene mutations has been proven to be effective in treating MTC patients with *RET* mutations [11]. Consequently, early diagnosis of MTC with DM and early intervention for high-risk patients may significantly improve patient survival.

Machine learning (ML) is a subfield of artificial intelligence technology. Compared to traditional predictive models, ML can enhance the accuracy of models by uncovering nonlinear relationships in large datasets [12, 13]. During medical treatment, vast amounts of data from patients are generated. Therefore, processing and analyzing these data using ML can offer a reliable reference for clinicians to diagnose diseases and prognosticate outcomes. Thus, our study aimed to develop a model based on the Surveillance, Epidemiology, and End Results (SEER) database to predict the occurrence of DM in patients with MTC.

2. Materials and Methods

2.1. Data Sources and Study Population. Data for this study were acquired from the SEER public databases, utilizing SEER*Stat 8.4.0.1 software for data extraction. Our study focused on patients diagnosed with MTC in the United States between 2004 and 2015. We excluded patients with missing data, unclear clinical and pathological conditions, uncertain histological classifications, or other types of thyroid cancer (TC). The histological types were restricted to medullary carcinomas. According to the International Classification of Diseases (ICD) for Oncology-3, patients' histological codes are 8345/3 and 8510/3, adopting AJCC 7th edition TNM stage. Variables included age, sex (male or female), race (White, Black, and others), year of diagnosis, Spanish-Hispanic origin, laterality (unilateral and bilateral), multifocality (solitary and multifocal), tumor size, extrathyroidal extension, lymph node metastasis, MTC subtypes, and DM. Distant metastasis means that the tumor invades at least one or more target organs such as brain, bone, liver, lung, and so on. As the SEER database contains public data, informed consent from relevant patients for the use of the SEER database for research purposes was not required, nor was the ethical approval. Our request for access to the SEER data was approved by the National Cancer Institute, USA (reference number 19238-Nov2021).

2.2. Screening for Risk Factors and Model Construction. Statistical analysis was conducted using SPSS software (version 26.0; IBM Corporation). In the univariable analysis, we employed Pearson's correlation analysis to examine the association between predictor variables, with results being presented in the form of heat maps. The predictive factors related to DM were initially screened through univariable analysis (p < 0.05), and the variables that met the criteria were incorporated into a multivariable logistic regression (LR) analysis. The receiver operating characteristic (ROC) curve was plotted and analyzed based on the results. An area under the ROC curve (AUC) greater than 0.5 was considered meaningful. All computed p values were two-sided, and statistical significance was accepted at <0.05.

The rate of DM of patients with MTC in the SEER database was low, resulting in an unbalanced original dataset. To establish a more accurate prediction model, it is essential to address this imbalance. In this study, we employed two techniques for processing the original dataset: oversampling and undersampling. We then used a correlation matrix to analyze the original and processed data. The synthetic minority oversampling technique (SMOTE) and undersampling are standard approaches for balancing class distribution in imbalanced datasets, widely used to improve prediction models [14]. The distribution of the target variables after the sampling process is illustrated in Figure 1. After data processing, the correlation between variables became more apparent, as demonstrated in Figure 2.

We used Python software (version 3.9.12, Python Software Foundation) to incorporate the selected variables include all variables in the ML model and construct a prediction model. The technically processed data (oversampled and undersampled data) were randomly divided into a training set (80%) and a test set (20%). The training set employed six commonly used ML algorithms: decision tree (DT), support vector machine (SVM), random forest (RF), k-nearest neighbors (KNN), extreme gradient boosting (XGBoost), and gradient boosting machine (GBM). Model evaluation was primarily based on accuracy, precision, recall, F1-score, and AUC value. The model with the highest AUC value was selected as the optimal model.

3. Results

3.1. Analysis of Patient Information. This study included a total of 2049 MTC patients, of which 138 (6.7%) developed DM and the remaining 1911 (93.3%) did not. The baseline characteristics of all patients are presented in Table 1.

In the univariable LR analysis, DM was significantly associated with age, sex, multifocality, tumor size, extrathyroidal extension, and lymph node metastasis (p < 0.05) (Table 2). These characteristic variables were incorporated into the multivariable LR analysis.

In the multivariable LR analysis, age [15] sex, extrathyroidal extension, lymph node metastasis, and tumor size were identified as independent predictors of DM in MTC. However, multifocality was not an independent predictive factor for the occurrence of DM in MTC. Further details can be found in Table 2. The ROC curve was plotted based on traditional multivariable LR results (AUC = 0.838, 95% confidence interval (CI): 0.808–0.868, p < 0.001). Detailed information is summarized in Figure 3.

For the analysis of the ML algorithm, six ML models were constructed and evaluated based on accuracy, precision, recall rate, F1-score, and AUC value. It was observed that ML models constructed after data oversampling outperformed those constructed after undersampling. Tables 3 and 4 provide details on the six ML models constructed from the over- and undersampled data. The ROC curves of the six ML models, constructed by oversampling and undersampling in the training and test sets, are depicted in Figure 4. In the models established using oversampled data, the AUC of all models was greater than



FIGURE 1: The distribution of the target variables after the sampling process. (a) Oversampling data, (b) undersampling data, and (c) target variable distribution of original data.

											1.0
Age	- 1	0.062	-0.047	-0.2	0.089	0.021	-0.019	-0.2	0.0051	0.16	- 1.0
Gender	0.062	1	-0.05	-0.029	0.19	0.15	0.26	0.037	0.024	0.25	- 0.8
Race	-0.047	-0.05	1	-0.13	0.051	0.039	0.014	0.015	0.025	0.0014	
Spanish-Hispanic	-0.2	-0.029	-0.13	1	0.053	0.0068	0.037	0.04	-0.049	-0.042	- 0.6
Tumor size	0.089	0.19	0.051	0.053	1	0.34	0.32	0.042	-0.01	0.42	0.4
Extrathyroidal extension	0.021	0.15	0.039	0.0068	0.34	1	0.47	0.22	-0.026	0.44	- 0.4
Lymph node metastasis	-0.019	0.26	0.014	0.037	0.32	0.47	1	0.22	-0.024	0.43	- 0.2
Multifocality	-0.2	0.037	0.015	0.04	0.042	0.22	0.22	1	0.017	0.1	
MTC subtypes	0.0051	0.024	0.025	-0.049	-0.01	-0.026	-0.024	0.017	1	0.03	- 0.0
Distant metastasis	0.16	0.25	0.0014	-0.042	0.42	0.44	0.43	0.1	0.03	1	0.2
	Age -	Gender -	Race -	Spanish-Hispanic	Tumor size -	Extrathyroidal extension	Lymph node metastasis -	Multifocality -	MTC subtypes -	Distant metastasis	-0.2
					()						

FIGURE 2: Continued.

Age	• 1	0.093	-0.021	-0.18	0.084	0.061	0.031	-0.15	0.0076	0.17
Gender	0.093	1	-0.053	-0.058	0.2	0.19	0.23	-0.027	-0.049	0.21
Race	-0.021	-0.053	1	-0.11	-0.03	-0.021	-0.034	-0.018	0.036	1.4e-15
Spanish-Hispanic	-0.18	-0.058	-0.11	1	0.037	0.027	0.016	0.023	0.048	0.011
Tumor size	0.084	0.2	-0.03	0.037	1	0.4	0.36	0.065	-0.054	0.36
Extrathyroidal extension	0.061	0.19	-0.021	0.027	0.4	1	0.49	0.28	-0.033	0.48
Lymph node metastasis	0.031	0.23	0.034	0.016	0.36	0.49	1	0.18	-0.081	0.45
Multifocality	-0.15	-0.027	-0.018	0.023	0.065	0.28	0.18	1	-0.056	0.16
MTC subtypes	0.0076	-0.049	0.036	0.0048	-0.054	-0.033	-0.081	-0.056	1	0.02
Distant metastasis	0.17	0.21	1.4e-15	0.011	0.36	0.48	0.45	0.16	0.02	1
	Age	Gender	Race	Spanish-Hispanic	Tumor size	Extrathyroidal extension	Lymph node metastasis	Multifocality	MTC subtypes	Distant metastasis
					(b)					
Distant metastasis	1	0.081	0.12	-0.0023	-0.02	0.22	0.27	0.23	0.058	0.011
Age	0.081	1	0.038	-0.064	-0.1	0.045	0.0037	-0.044	-0.2	-0.025
Gender	0.12	0.038	1	-0.023	0.0057	0.17	0.14	0.23	0.083	-0.01
Race	-0.0023	-0.064	-0.023	1	-0.11	0.012	-0.028	-0.011	-0.0065	-0.0063
Spanish-Hispanic	-0.02	-0.1	0.0057	-0.11	1	0.064	0.038	0.066	0.046	-0.02
Tumor size	0.22	0.045	0.17	0.012	0.064	1	0.32	0.31	0.0055	-0.0028
Extrathyroidal extension	0.27	0.0037	0.14	-0.028	0.038	0.32	1	0.48	0.17	-0.018
Lymph node metastasis	0.23	-0.044	0.23	-0.011	0.066	0.31	0.48	1	0.21	-0.03
Multitocality	0.058	-0.2	0.083	0.0065	0.046	0.0055	0.17	0.21	1	-0.011
MIC subtypes	0.011	-0.025	-0.01	-0.063	-0.02	-0.0028	-0.018	-0.03	-0.011	1
	Distant metastasis	Age	Gender	Race	Spanish-Hispanic	Tumor size	Extrathyroidal extension	Lymph node metastasis	Multifocality	MTC subtypes
					(c)					

FIGURE 2: Heatmaps of the correlation between characteristic features of the patients in different datasets. (a) Oversampling data, (b) undersampling data, and (c) original data.

0.850, with the RF model performing better than the other models. The RF model demonstrated accuracy, precision, recall rate, *F*1-score, and AUC value of 0.890, 0.847,0.946, 0.894, and 0.946, respectively, as well as a higher AUC value than the LR model. This indicates that the diagnostic efficiency of the ML algorithm surpasses that of the traditional LR model and exhibits excellent prediction performance. Employing RF for

feature selection, as illustrated in Figure 5, revealed that lymph node metastasis was the most critical factor in determining whether MTC patients also have DM.

This study developed an online network calculator for evaluating the risk of distant metastasis in MTC patients, which can be applied to clinical patients (https://121.43.117. 60:8000/).

TABLE 1: The detailed demographic information of the patients with MTC.

Categories	With DM $(n = 138)$	Without DM (<i>n</i> = 1911)	<i>p</i> value
Age, <i>n</i> (%)			< 0.001
<55	47 (34.1%)	959 (50.2%)	
≥55	91 (65.9%)	952 (49.8%)	
Gender (<i>n</i> %)			< 0.001
Female	50 (36.2%)	1149 (60.1%)	
Male	88 (63.8%)	762 (39.9%)	
Race (<i>n</i> %)			0.891
White	116 (84.1%)	1619 (84.7%)	
Black	14 (10.1%)	159 (8.3%)	
Other	8 (5.8%)	133 (7.0%)	
Year of diagnosis			0.367
2004-2009	54 (39.1%)	823 (43.1%)	
2010-2015	84 (60.9%)	1088 (56.9%)	
Spanish-Hispanic-Latino (<i>n</i> %)			0.372
Yes	17 (12.3%)	289 (15.1%)	
No	121 (87.7%)	1622 (84.9%)	
MTC subtypes (<i>n</i> %)			0.622
MTC with amyloid stroma	4 (2.9%)	71 (3.7%)	
MTC NOS ^a	134 (97.1%)	1840 (93.2%)	
Laterality (<i>n</i> %)			0.419
Unilateral	138 (100%)	1902 (99.5%)	
Bilateral	0 (0%)	9 (0.5%)	
Multifocality (<i>n</i> %)			0.009
Solitary tumor	86 (62.3%)	1388 (72.6%)	
Multifocal tumor	52 (37.7%)	523 (37.7%)	
Tumor size (<i>n</i> %)			< 0.001
≤2	27 (19.6%)	1030 (53.9%)	
2-4	46 (33.3%)	594 (31.1%)	
≥4	65 (47.1%)	287 (15.0%)	
Extrathyroidal extension (<i>n</i> %)			< 0.001
Yes	57 (41.3%)	309 (16.2%)	
No	81 (58.7%)	1602 (83.8%)	
Lymph node metastasis (<i>n</i> %)			< 0.001
No	24 (17.4%)	1206 (63.1%)	
Cervical central lymph node	36 (26.1%)	288 (15.1%)	
Cervical lateral lymph node	70 (50.7%)	359 (18.8%)	
Yes NOS	8 (5.8%)	58 (3.0%)	

MTC, medullary thyroid carcinoma; DM, distant metastasis; NOS, not otherwise specified.

TABLE 2: Univariable analysis and multivariable analysis of variables related to distant metastasis.

		Univariable analysis			Multivariable analysis			
	OR	95% CI	p value	OR	95% CI	p value		
Age (year)								
<55	0.513	0.357-0.513	< 0.001	0.480	0.323-0.713	< 0.001		
≥55	Ref			Ref				
Gender								
Female	0.377	0.263-0.540	< 0.001	0.664	0.450-0.980	0.039		
Male	Ref			Ref				
Race								
White	1.191	0.569-2.491	0.642					
Black	1.464	0.596-3.596	0.406					
Other	Ref							
Spanish-Hispanic								
Yes	Ref							
No	1.268	0.752-2.139	0.373					
MTC subtypes								
MTC with amyloid stroma	0.774	0.278-2.150	0.623					
MTC NOS	Ref							
Multifocality								
Solitary tumor	0.623	0.435-0.892	0.010	0.866	0.578-1.298	0.486		

		Univariable analysi	s	Multivariable analysis			
	OR	95% CI	p value	OR	95% CI	p value	
Multifocal tumor	Ref			Ref			
Tumor size (cm)							
≤2	0.116	0.073-0.185	< 0.001	0.287	0.173-0.476	< 0.001	
2-4	0.342	0.229-0.512	< 0.001	0.555	0.360-0.0855	0.008	
≥ 4	Ref			Ref			
Extrathyroidal extension							
Yes	Ref			Ref			
No	0.136	0.095-0.195	< 0.001	0.364	0.240-0.554	< 0.001	
LNM							
No	0.144	0.062-0.335	< 0.001	0.327	0.132-0.806	0.015	
Cervical central lymph node	0.906	0.401 - 2.050	0.813	1.021	0.437-2.385	0.962	
Cervical lateral lymph node	1.414	0.647-3.091	0.386	1.269	0.564-2.858	0.565	
Yes NOS	Ref			Ref			

TABLE 2: Continued.

MTC, medullary thyroid carcinoma; NOS, not otherwise specified; OR, odds ratio; CI, confidence interval.



FIGURE 3: LR models predict the ROC curve of distant metastasis in MTC patients.

TABLE 3: Comparison of prediction performance between different models constructed from oversampling data.

Model	Accuracy	AUC	Precision	Recall rate	F1-score
DT	0.764	0.930	0.782	0.724	0.752
RF	0.890	0.946	0.847	0.946	0.894
SVC	0.781	0.853	0.761	0.811	0.785
KNN	0.830	0.918	0.777	0.917	0.836
GBM	0.813	0.883	0.788	0.848	0.817
XGBoost	0.879	0.934	0.851	0.915	0.882

AUC, area under the receiver operating characteristic curve; DT, decision tree; SVM, support vector machine; RF, random forest; KNN, k-nearest neighbors; XGBoost, extreme gradient boosting; GBM, gradient boosting machine.

4. Discussion

Patients with MTC account for only 5% of the total number of individuals newly diagnosed with TC, while the global incidence rate of MTC is rising rapidly. Deaths from MTC comprise approximately 13% of the total mortality rate of TC, and the 10-year overall survival rate of MTC ranges between 65% and 71%. However, when MTC occurs with DM, the 10-year overall survival rate can decrease to 40–44% [15, 16]. MTC neither concentrates radioactive iodine nor is it inhibited by thyroxine [17]. Total thyroidectomy is the primary treatment method for MTC, with the decision to

Model	Accuracy	AUC	Precision	Recall rate	F1-score
DT	0.803	0.751	0.827	0.800	0.813
RF	0.732	0.769	0.741	0.766	0.754
SVC	0.750	0.784	0.750	0.800	0.774
KNN	0.714	0.789	0.750	0.700	0.760
GBM	0.785	0.814	0.781	0.833	0.806
XGBoost	0.767	0.768	0.774	0.800	0.786

TABLE 4: Comparison of prediction performance between different models constructed from undersampling data.

AUC, are under the receiver operating characteristic curve; DT, decision tree; SVM, support vector machine; RF, random forest; KNN, k-nearest neighbors; XGBoost, extreme gradient boosting; GBM, gradient boosting machine.



FIGURE 4: ROC curves of six ML algorithms in different datasets. (a) The ROC curves of the six ML algorithms model in the test set with oversampling. (b) The ROC curves of the six ML algorithms model in the training set with oversampling. (c) The ROC curves of the six ML algorithms model in the test set with undersampling. (d) The ROC curves of the six ML algorithms model in the training set with undersampling. ROC, receiver operating characteristic; ML, machine learning; AUC, area under the receiver operating characteristic curve.

perform lymph node dissection depending on the specific situation. Adjuvant radiation therapy can be considered for MTC patients with incomplete resection, a high risk of local recurrence, or DM [10]. Radiotherapy can provide continuous control in patients with DM and prevent further progression [18]. However, the impact of radiotherapy on patients' survival rates remains controversial. In patients without DM, radiotherapy may cause more harm than good [19]. Some perspectives suggest that the role of radiation therapy in MTC is limited to patients who are ineligible or have contraindications for surgical treatment or targeted drugs [20]. Targeted drugs are recommended for patients with DM, particularly because studies have demonstrated [11, 21] that *RET*-specific inhibitors (selpercatinib and

pralsetinib) are effective and promising therapies for MTC patients with DM and progression. The prognosis and treatment effectiveness of MTC are largely related to tumor staging; therefore, early diagnosis is a crucial objective in the management of MTC patients [22]. Previous research on MTC has mostly focused on prognosis and analysis of survival [23, 24].

However, there are few studies on the DM of MTC. Utilizing independent predictors to predict DM can help physicians better evaluate patients with MTC and provide them with more effective individualized treatment options.

Univariable analysis showed that age, sex, multifocality, tumor size, extrathyroidal extension, and lymph node metastasis were independent predictors of DM. However,



FIGURE 5: Feature importance derived from the RF model. The plot shows the relative importance of the variables in the RF model. MTC, medullary thyroid carcinoma.

multivariable analysis indicated that multifocality could not serve as an independent predictor of DM in patients with MTC. This finding is consistent with the conclusion of the RF feature selection, and it is generally believed that multilocality has an independent predictive effect on cervical lymph node metastasis in MTC [25]. Nonetheless, multifocality had a relatively small impact on predicting the occurrence of DM in patients with MTC, which aligns with findings of previous research [25, 26]. RF feature selection revealed that extrathyroidal extension was a key factor in predicting DM, while lymph node metastasis was the most important predictor of DM, consistent with a previous study [26]. We also identified tumor size was an important predictor. Compared with tumors larger than 4 cm, the odds ratio (OR) for tumors of 2–4 cm and \leq 2 cm was 0.555 and 0.287, respectively. As tumor size gradually increases, the risk of DM in MTC also increases. Tumor size significantly impacts the recurrence and long-term survival rates of MTC [24]. Extrathyroidal extension and tumor size are also crucial predictive factors for lymph node and DM in MTC [6, 16]. Meanwhile, extrathyroidal extension and tumor size are directly related to T staging in TNM staging, suggesting that tumor stage can also serve as a predictive factor for DM. Contrary to a previous study [27], sex was considered as an independent predictor of DM. We also discovered that female sex was a protective factor for DM. This conclusion is similar to that of a previous study [26]. In our study, 55 years of age was used as the cutoff age [27] and it showed that older patients were more likely to develop DM than younger patients. Therefore, older patients should be actively followed up and regularly examined. In this study, race could not independently predict DM in patients with MTC, which is consistent with results of previous research [26, 27]. In traditional LR, MTC subtypes and Spanish-Hispanic could not be used as independent predictors, and their influence on the feature selection of RF was also small.

We constructed six predictive models based on the SEER database to predict DM in patients with MTC and evaluated six algorithmic models based on accuracy, precision, recall rate, F1-score, and AUC value. We employed the SMOTE technique to address unbalanced datasets and concluded that, for unbalanced datasets used to build ML models, SOMTE is superior to undersampling [14]. By oversampling and undersampling, we enhanced the performance of the model and determined that the prediction model established by oversampling outperformed the one established by undersampling. This may be attributed to fewer patients with DM among MTC patients, resulting in limited ability of the model to identify key predictive factors for patients with combined DM. This study established six ML algorithms, among which RF demonstrated excellent predictive performance (AUC = 0.946), surpassing that of the traditional LR model (AUC = 0.838). Therefore, RF was the best model for predicting MTC patients with DM using the SEER database.

5. Limitations

However, there are some limitations to this study. First, as this study is based on demographics of North American, other populations should be used for validation in future research. Second, the predictive performance of the model warrants further optimization, and additional predictive factors potentially related to DM should be incorporated into the prediction model in future studies. Finally, due to the limitations of the database, tumor markers such as CEA and AFP were not included in MTC patients. We will continue to improve and supplement the model in future studies.

6. Conclusions

In conclusion, this study aimed to identify independent predictors of DM in patients with MTC and to develop a prediction model utilizing ML algorithms. Our analysis, based on the SEER database, demonstrated that age, sex, tumor size, extrathyroidal extension, and lymph node metastasis were significant independent predictors of DM in MTC patients. The RF ML algorithm outperformed the traditional LR model in predicting DM, providing a more accurate and reliable tool for clinical use.

The application of the SMOTE technique for addressing unbalanced datasets was proven to be effective in enhancing the performance of the prediction model. Our findings underscore the importance of early diagnosis and individualized treatment plans for MTC patients, ultimately contributing to improved patient outcomes.

Data Availability

The dataset presented in this study can be found at https:// seer.cancer.gov. Further inquiries can be directed to the corresponding author.

Conflicts of Interest

The authors declare that they have no conflicts of interest.

Acknowledgments

We are very grateful to Professor Xu Zhang, a biostatistician from the First Affiliated Hospital of Anhui Medical University, for evaluating the experimental design and analysis of this article and providing valuable feedback. We would like to thank Editage (https://www.editage.com/) for English language editing.

References

- C. La Vecchia, M. Malvezzi, C. Bosetti et al., "Thyroid cancer mortality and incidence: a global overview," *International Journal of Cancer*, vol. 136, no. 9, pp. 2187–2195, 2015.
- [2] T. Kondo, S. Ezzat, and S. L. Asa, "Pathogenetic mechanisms in thyroid follicular-cell neoplasia," *Nature Reviews Cancer*, vol. 6, no. 4, pp. 292–306, 2006.
- [3] S. Roman, R. Lin, and J. A. Sosa, "Prognosis of medullary thyroid carcinoma: demographic, clinical, and pathologic predictors of survival in 1252 cases," *Cancer*, vol. 107, no. 9, pp. 2134–2142, 2006.
- [4] A. Matrone, C. Gambale, A. Prete, and R. Elisei, "Sporadic medullary thyroid carcinoma: towards a precision medicine," *Frontiers in Endocrinology*, vol. 13, Article ID 864253, 2022.
- [5] R. Elisei, A. Tacito, T. Ramone et al., "Twenty-five years experience on RET genetic screening on hereditary MTC: an update on the prevalence of germline RET mutations," *Genes*, vol. 10, no. 9, p. 698, 2019.
- [6] A. Kotwal, D. Erickson, J. R. Geske, I. D. Hay, and M. R. Castro, "Predicting outcomes in sporadic and hereditary medullary thyroid carcinoma over two decades," *Thyroid*, vol. 31, no. 4, pp. 616–626, 2021.
- [7] O. Twito, S. Grozinsky-Glasberg, S. Levy et al., "Clinicopathologic and dynamic prognostic factors in sporadic and familial medullary thyroid carcinoma: an Israeli multi-center study," *European Journal of Endocrinology*, vol. 181, no. 1, pp. 13–21, 2019.

9

- [8] R. S. Sippel, M. Kunnimalaiyaan, and H. Chen, "Current management of medullary thyroid cancer," *The Oncologist*, vol. 13, no. 5, pp. 539–547, 2008.
- [9] C. Nashed, S. V. Sakpal, S. Cherneykin, and R. S. Chamberlain, "Medullary thyroid carcinoma metastatic to skin," *Journal of Cutaneous Pathology*, vol. 37, no. 12, pp. 1237–1240, 2010.
- [10] S. A. Wells Jr, S. L. Asa, H. Dralle et al., "Revised American Thyroid Association guidelines for the management of medullary thyroid carcinoma," *Thyroid*, vol. 25, no. 6, pp. 567–610, 2015.
- [11] L. J. Wirth, E. Sherman, B. Robinson et al., "Efficacy of selpercatinib in RET-altered thyroid cancers," *New England Journal of Medicine*, vol. 383, no. 9, pp. 825–835, 2020.
- [12] A. M. Darcy, A. K. Louie, and L. W. Roberts, "Machine learning and the profession of medicine," *Journal of the American Medical Association*, vol. 315, no. 6, pp. 551-552, 2016.
- [13] Z. Obermeyer and E. J. Emanuel, "Predicting the future- big data, machine learning, and clinical medicine," *New England Journal of Medicine*, vol. 375, no. 13, pp. 1216–1219, 2016.
- [14] W. Liu, S. Wang, Z. Ye, P. Xu, X. Xia, and M. Guo, "Prediction of lung metastases in thyroid cancer using machine learning based on SEER database," *Cancer Medicine*, vol. 11, no. 12, pp. 2503–2515, 2022.
- [15] Z. Chen, Y. Mao, T. You, and G. Chen, "Establishment and validation of a nomogram model for predicting distant metastasis in medullary thyroid carcinoma: an analysis of the SEER database based on the AJCC 8th TNM staging system," *Frontiers in Endocrinology*, vol. 14, Article ID 1119656, 2023.
- [16] R. W. Randle, C. J. Balentine, G. E. Leverson et al., "Trends in the presentation, treatment, and survival of patients with medullary thyroid cancer over the past 30 years," *Surgery*, vol. 161, no. 1, pp. 137–146, 2017.
- [17] Z. T. Sahli, J. K. Canner, M. A. Zeiger, and A. Mathur, "Association between age and disease specific mortality in medullary thyroid cancer," *The American Journal of Surgery*, vol. 221, no. 2, pp. 478–484, 2021.
- [18] O. Hamdy, S. Awny, and I. H. Metwally, "Medullary thyroid cancer: epidemiological pattern and factors contributing to recurrence and metastasis," *Annals of the Royal College of Surgeons of England*, vol. 102, no. 7, pp. 499–503, 2020.
- [19] J. A. Call, J. S. Caudill, B. Mciver, and R. L. Foote, "A role for radiotherapy in the management of advanced medullary thyroid carcinoma: the mayo clinic experience," *Rare Tumors*, vol. 5, no. 3, pp. 128–131, 2013.
- [20] S. Huang, J. Zhong, Z. Zhang et al., "Prognosis of radiotherapy in medullary thyroid carcinoma patients without distant metastasis," *Translational Cancer Research*, vol. 10, no. 11, pp. 4714–4726, 2021.
- [21] A. Kukulska, J. Krajewska, Z. Kołosza et al., "Stereotactic radiotherapy is a useful treatment option for patients with medullary thyroid cancer," *Bone Marrow Concentrate Endocrine Disorders*, vol. 21, no. 1, p. 160, 2021.
- [22] V. Subbiah, D. Yang, V. Velcheti, A. Drilon, and F. Meric-Bernstam, "State-of-the-art strategies for targeting RETdependent cancers," *Journal of Clinical Oncology*, vol. 38, no. 11, pp. 1209–1221, 2020.
- [23] F. Orlandi, P. Caraci, A. Mussa, E. Saggiorato, G. Pancani, and A. Angeli, "Treatment of medullary thyroid carcinoma: an update," *Endocrine-Related Cancer*, vol. 8, no. 2, pp. 135–147, 2001.
- [24] J. Tang, S. Jiang, L. Gao et al., "Construction and validation of a nomogram based on the log odds of positive lymph nodes to predict the prognosis of medullary thyroid carcinoma after

surgery," Annals of Surgical Oncology, vol. 28, no. 8, pp. 4360-4370, 2021.

- [25] L. Chen, Y. Wang, K. Zhao, Y. Wang, and X. He, "Postoperative nomogram for predicting cancer-specific and overall survival among patients with medullary thyroid cancer," *International Journal of Endocrinology*, vol. 2020, Article ID 8888677, 13 pages, 2020.
- [26] W. Fan, C. Xiao, and F. Wu, "Analysis of risk factors for cervical lymph node metastases in patients with sporadic medullary thyroid carcinoma," *Journal of International Medical Research*, vol. 46, no. 5, pp. 1982–1989, 2018.
- [27] M. K. Le, M. Kawai, T. Odate, H. G. Vuong, N. Oishi, and T. Kondo, "Metastatic risk stratification of 2526 medullary thyroid carcinoma patients: a study based on surveillance, epidemiology, and end results database," *Endocrine Pathology*, vol. 33, no. 3, pp. 348–358, 2022.