# Prediction of Skin Tumor Invasiveness: A National Analysis Through Explainable Artificial Intelligence (XAI)

**Marcus Augusto Padilha da Mata**[1], **Plínio de Sá Leitão Júnior**[1]

[1]Institute of Informatics (INF) – Federal University of Goias (UFG)
PO Box 74.690-900 – Goiânia – GO – Brazil

*Abstract.* *In Brazil, skin tumors represents the type of neoplasm with the highest incidence rate among the population. Because of this, this study explores the invasiveness of this disease using computational techniques to understand how specific patient characteristics influence its progression. Through the analysis of data provided by the Cancer Hospital Registry (RHC) of the National Cancer Institute José Alencar Gomes da Silva (INCA), and with the aid of Artificial Intelligence (AI) algorithms explained by the SHapley Additive exPlanations (SHAP) approach, the study reveals that the invasiveness of skin cancer is affected in a significantly different way by the individual characteristics of patients compared to analyses based on more general attributes. These findings underline the importance of personalization in medicine, suggesting that a deeper understanding of individual characteristics can lead to more accurate diagnoses and more effective treatments. Furthermore, the research highlights the role of XAI in clarifying these relationships, pointing to the need for more refined approaches in prevention, treatment, and the formulation of public health policies aimed at combating skin tumors, despite limitations such as data imbalance encountered during the study.*

## 1. Introduction

Skin tumor, defined as a disease resulting from the unrestrained proliferation of skin cells, was cited as the most prevalent type of tumor in Brazil and the world in 2021, according to [BRAZIL'S MINISTRY OF HEALTH 2021]. This ailment may manifest in diverse forms, influencing various strata of the epidermis. Additionally, the proliferation of these tumors can amplify the pain in individuals owing to the compression exerted on surrounding tissues and the resultant inflammatory reaction, as indicated in [Fidler 2003]. Hence, distinguishing between invasive and non-invasive tumors is deemed essential for the precise diagnosis and therapy, profoundly impacting in the patient's quality of life, longevity and the treatment's success. In Brazil, skin cancer constitutes 27% of all malign tumors, with projections indicating approximately 220,490 new instances annually from 2023 to 2025, as suggested by [Santos et al. 2023]. In the year 2022, the Brazilian Unified Health System (SUS) dedicated around BRL 47 million to chemotherapy treatments for skin cancer, as documented in [BRAZIL'S MINISTRY OF HEALTH 2022].

In this scenario, the INCA, which operates under the aegis of the Brazil's Ministry of Health, assumes a pivotal role in the ongoing surveillance and support of individuals diagnosed with cancer through the utilization of the RHC. As detailed by [Lopes et al. 2021], this comprehensive system, which has been integrated into hospital infrastructure, is designed to collect, store, efficiently process, and rigorously analyze patient-specific data.

Furthermore, in the current technological environment, there is a notable advancement in the adoption of AI across various research fields, notably in the diagnosis and analysis of cancer, as mentioned in [Subasi et al. 2022]. Statistical models, developed from AI and built through the training of machine learning algorithms with data, facilitate efficient processing and selection of information, resulting in highly accurate predictions. However, beyond the precision of these forecasts, the capability to interpret these models holds crucial importance, as indicated in [Molnar 2020]. In this setting, as discussed in [Silva et al. 2022], the SHAP approach stands out for its ability to assess the impact of each feature on the predictions of a model, showing how specific characteristics influence the outcomes. This contributes to a better understanding of how different attributes affect the predictive dynamics of the model.

Within this context, the study delves into the prediction of skin tumors invasiveness in Brazil, focusing on how specific tumor characteristics influence this invasiveness. Rather than merely examining the overall impact of attributes, this work scrutinizes how each feature within an attribute affects the progression of tumors across the country. This approach facilitates a deeper understanding of skin tumors in Brazil and may lead to more personalized prevention and treatment methods. By identifying the factors that cause some tumors to invade, healthcare professionals can create more effective treatment plans. Furthermore, the study contributes to the advancement of oncology research by improving the understanding of skin tumors in Brazil, potentially leading to more personalized prevention and treatment methodologies.

## 2. Related Works

The studies linked to this research adopt machine learning models and focus on improving the interpretability of medical predictions, taking advantage of the SHAP approach. Each investigation explores a different application of these tools in the context of cancer, demonstrating the versatility and effectiveness of SHAP in revealing the decisive factors in different contexts of the disease.

The study of [Sorayaie Azar et al. 2022] conducted an in-depth analysis on ovarian cancer survival prediction using the SEER database, implementing Random Forest and XGBoost machine learning models. The study utilized the SHAP approach to identify and quantify the contributions of significant attributes such as histologic type of the tumor and year of diagnosis. This analysis highlighted the robust performance of Random Forest and XGBoost for classification and regression tasks, respectively, showcasing SHAP's role in enhancing model interpretability and decision-making transparency.
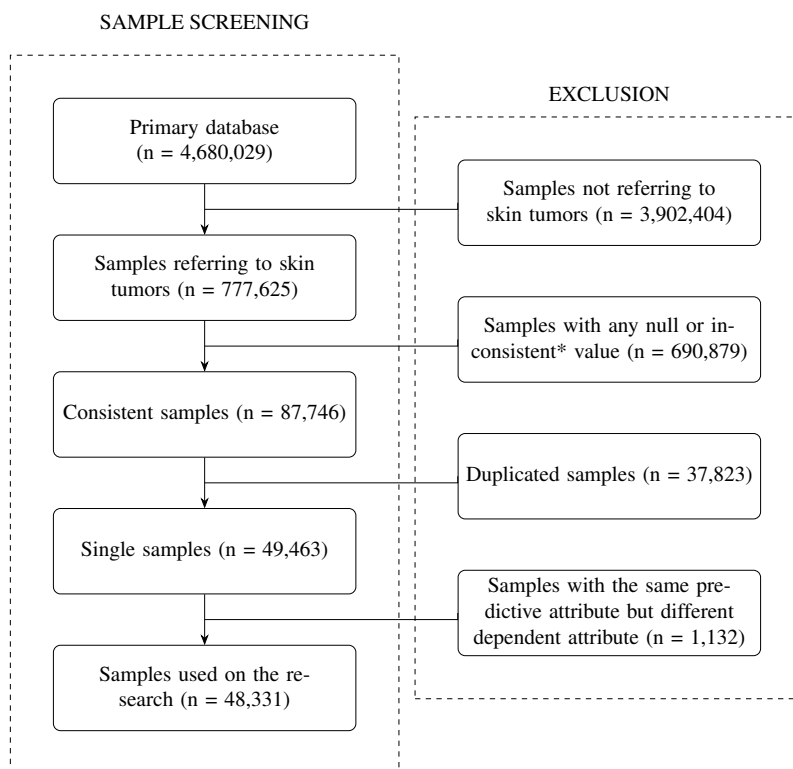
[Lee et al. 2022] focused on assessing the risk of developing a second primary cutaneous tumor in skin cancer survivors, analyzing data from 1248 patients across five cancer registries. The researchers applied various machine learning algorithms, notably achieving optimal results with Random Forest. To address data imbalance, the Synthetic Minority Over-sampling Technique (SMOTE) was employed. The SHAP approach was pivotal in this study, revealing critical predictive attributes such as age, cancer stage, gender, and regional lymph node involvement, thus providing deeper insights into the model's predictive dynamics.

[Alsinglawi et al. 2022] explored the prediction of ICU stay durations for lung cancer patients using the comprehensive MIMIC-III database. The study faced challen-

ges related to data imbalance, which were addressed by employing SMOTE and Adaptive Synthetic (ADASYN) techniques. Predominantly utilizing Random Forest, the application of the SHAP approach enabled the identification of key clinical attributes that significantly influence the prediction of ICU stay durations.

## 3. Materials and Methods

This study employs a dataset comprising 4,680,029 anonymized patient samples diagnosed with cancer, covering 47 attributes. The data originate from RHC, managed by INCA. An initial exploratory analysis of the dataset was conducted to deepen the understanding of the intrinsic characteristics of the collected data and to identify the phenomena represented by the attributes. Subsequently, a screening of the samples was carried out, as illustrated in the flowchart of Figure 1, reducing the initial number of 4,680,029 samples to 48,331. This final dataset includes 46,983 cases of invasive skin tumors and 1,348 cases of non-invasive tumors.

SAMPLE SCREENING

Primary database
(n = 4,680,029)

EXCLUSION

Samples not referring to
skin tumors (n = 3,902,404)

Samples referring to skin
tumors (n = 777,625)

Samples with any null or in-
consistent* value (n = 690,879)

Consistent samples (n = 87,746)

Duplicated samples (n = 37,823)

Single samples (n = 49,463)

Samples with the same pre-
dictive attribute but different
dependent attribute (n = 1,132)

Samples used on the re-
search (n = 48,331)

\* Inconsistent values, for example, entries of '3' in the sex attribute - where '1' represents male and '2' female, as defined in the data dictionary - were considered inconsistent and led to the exclusion of the corresponding samples.

**Figure 1. Screening flowchart of the samples**
Source: Compiled by the author with data extracted from the research.

In addition to sample selection, new attributes were derived from the original data to align with the specific objectives of the study. For instance, the dependent attribute of this study, the invasive behavior of the tumor, was derived from the histological type of the tumor attribute. The original attribute consists of a five-digit code, where the last digit

indicates the tumor's biological behavior. After deriving the new attribute, the original attribute was excluded from further analysis in the next step.

A screening process was also undertaken to ensure the relevance and utility of the attributes in predictive modeling, whereby attributes with negligible variation, deemed irrelevant for prediction, or found to be redundant due to dependency on others were systematically excluded. Following this culling process, the domains of the remaining attributes were encoded, making them suitable for detailed statistical analysis. The *Cramér's V method*[1] was employed to examine the associations between attributes, leading to the identification of attributes with high associations. Attributes with high association were identified, and to avoid redundancies in the model, only one attribute from each pair with strong association was retained.

Table 1 presents all the attributes used in this study. For each listed attribute, the table displays its domain, elucidating the possible characteristics it can assume. Moreover, to facilitate analysis in machine learning models, each characteristic is associated with a specific numeric code in parentheses.

**Table 1. Coding of attribute domains**

| ATTRIBUTE | DOMAIN AND CODE |
|---|---|
| **Primary Tumor Location** | Eyelid (0); External ear (1); Face (2); Scalp and neck (3); Trunk (4); Upper limbs (5); Lower limbs (6); Skin overlay lesion (7) |
| **Geographic Region** | Midwest (0); North (1); Northeast (2); Southeast (3); South (4) |
| **Age Group** | 000-009 (0); 010-019 (1); 020-029 (2); 030-039 (3); 040-049 (4); 050-059 (5); 060-069 (6); 070-079 (7); 080-089 (8); 090-103 (9) |
| **Marital Status** | Single (0); Married (1); Widower (2); Separate (3); Stable union (4) |
| **Education** | None (0); Incomplete fundamental (1); Complete fundamental (2); High School (3); Incomplete higher education (4); Complete higher (5) |
| **Race** | White (0); Black (1); Yellow (2); Brown (3); Indigenous (4) |
| **History of Alcohol Consumption** | Never (0); Former consumer (1); Yes (2) |
| **History of Tobacco Consumption** | Never (0); Former consumer (1); Yes (2) |
| **Origin of Forwarding** | SUS (0); No SUS (1); Came on its own (2) |
| **Case Type** | Analytical (0); Non-analytical (1) |
| **Family History of Cancer** | Yes (0); No (1) |
| **More than One Tumor** | Yes (0); No (1) |
| **Gender** | Male (0); Female (1) |
| **Invasiveness of the Tumor** | Invasive (0); Non-Invasive (1) |

Source: Compiled by the author with data extracted from the research.

---

[1]Cramér's V is a statistical measure that quantifies the association between two nominal variables, ranging from 0 (indicating no association) to 1 (indicating perfect association), based on chi-square statistics. It is utilized to identify potential redundancies among variables within statistical models.

### 3.1. Implementation of the machine learning algorithm

The algorithms selected for analysis include XGBoost, Support Vector Machine, Random Forest, k-Nearest Neighbors, Neural Networks, Naive Bayes, and Decision Tree. This selection was based on previous works by [Ghazal et al. 2022] and [Liu et al. 2021], who utilized these algorithms in similar contexts. Following these works, the approach of 10-fold cross-validation was adopted for the selection of the optimal hyperparameters for each algorithm. The hyperparameter tuning process took into account the imbalance present in the data. The results of this tuning are showed in Table 2.

#### Table 2. Best Hyperparameters for Each Model

| Model | Best Hyperparameters |
| --- | --- |
| Random Forest | criterion = entropy, max_depth = 10, min_samples_leaf = 4, min_samples_split = 10, n_estimators = 200, class_weight = balanced |
| Decision Tree | criterion = gini, max_depth = 10, min_samples_leaf = 4, min_samples_split = 5 |
| XGBoost | colsample_bytree=1.0, gamma = 0.1, learning_rate = 0.2, max_depth = 7, n_estimators = 200, subsample = 0.5, scale_pos_weight = 46331/1348 |
| Naïve Bayes | alpha = 0.0 |
| k - Nearest Neighbors | n_neighbors = 4 |
| Neural Networks | activation = relu, learning_rate = 0.0001, hidden_layer_sizes = 200, solver = adam, batch_size = 64, class_weight = balanced |
| Support-Vector Machine | C = 0.1, gamma = 0.01, kernel = rbf, class_weight = balanced |

Source: Compiled by the author with data extracted from the research.

For evaluating the performance of the models, the following metrics were selected: F1 Score, highlighting the harmony between precision and recall; Area under the ROC Curve (AUC), evaluating the model's ability to effectively distinguish between classes; and Precision, focusing on the accuracy in identifying positive cases. The choice of these metrics was influenced by studies by [Yu et al. 2021] and [Taghizadeh et al. 2022]. Given the imbalanced data in this research, Balanced Accuracy was also used because, as discussed by [Luo et al. 2019], this metric is effective in measuring the model's capability to predict outcomes in imbalanced data.

Among the evaluated models, as seen in Table 3, the Support Vector Machine stood out as the most effective. This result is particularly noteworthy regarding its Balanced Accuracy, highlighting its ability to handle the existing imbalance in the data. Despite the low F1 Score, the significant AUC and Precision values confirm the suitability of this algorithm for accurate and reliable predictions in the context of this research.

### 3.2. Creation of Explainability Models through the SHAP Approach

To understand which attributes most influence the predictions of machine learning models, the SHAP approach was chosen, based on the work of [Lundberg and Lee 2017]. This approach uses Shapley Values to measure the impact of each attribute on predictions. According to [Shapley et al. 1953], these values show the contribution of an attribute considering all possible combinations of attributes, ensuring a fair assessment of its importance.

**Table 3. Ranking of models on a national scope**

| Algorithm | Balanced Accuracy | F1 Score | AUC | Precision |
|---|---|---|---|---|
| Support-Vector Machine | 0.65758 | 0.77442 | 0.72635 | 0.98576 |
| Random Forest | 0.63428 | 0.91048 | 0.71612 | 0.98090 |
| Decision Tree | 0.63141 | 0.80450 | 0.64472 | 0.98255 |
| k - Nearest Neighbors | 0.50528 | 0.98429 | 0.53275 | 0.97238 |
| XGBoost | 0.50000 | 0.98586 | 0.50000 | 0.97211 |
| Naïve Bayes | 0.50000 | 0.98586 | 0.65652 | 0.97211 |
| Neural Networks | 0.50000 | 0.98586 | 0.64297 | 0.97211 |

The SHAP approach distinguishes between *local SHAP values*, which show the contribution of an attribute on a specific prediction, and *global SHAP values*, which indicate the overall impact of an attribute on the dataset. Local SHAP values can be positive or negative, indicating whether the presence of the attribute increases or decreases the chance of the prediction. Global SHAP values are calculated by the absolute average of local SHAP values, showing the overall importance of the attribute.

After selecting the Support Vector Machine algorithm, a model was trained with all the data to improve accuracy and explainability, as performed in the official SHAP documentation by [Lundberg et al. 2020]. Implementing this approach creates an explainability model that identifies and explains the impact of the most significant attributes on the model's predictions.

## 4. Results and Discussion

Table 4 organizes the attributes in order of impact on prediction, based on global SHAP values. It begins with the primary tumor location as the most influential attribute, down to the type of case. This ordering highlights the sequence or priority given to each characteristic by the model, emphasizing the relative importance and specific role that each attribute plays in the outcome of the predictions.

**Table 4. Classification of attributes using global SHAP values.**

| | Utilized Attributes | | Utilized Attributes |
|---|---|---|---|
| 1st | Primary Tumor Location | 8th | Family History of Cancer |
| 2nd | Geographic Region | 9th | History of Tobacco Consumption |
| 3rd | More Than One Tumor | 10th | Race |
| 4th | Gender | 11th | Referral Origin |
| 5th | Age Group | 12th | History of Alcohol Consumption |
| 6th | Education Level | 13th | Type of Case |
| 7th | Marital Status | | |

Figure 2 displays strip plots of the five most impactful attributes on the model, with the features in the charts aligned vertically and the Local SHAP Values on the horizontal axis. Positive values suggest a greater probability of predicting tumor invasiveness, while negative values indicate the opposite. The proximity of these values to the zero axis highlights the intensity of each feature's impact on specific cases. This figure details the

manner in which the characteristics of the most significant attributes contribute to the model's predictions, with each point on the charts representing the unique impact of each sample in the dataset.
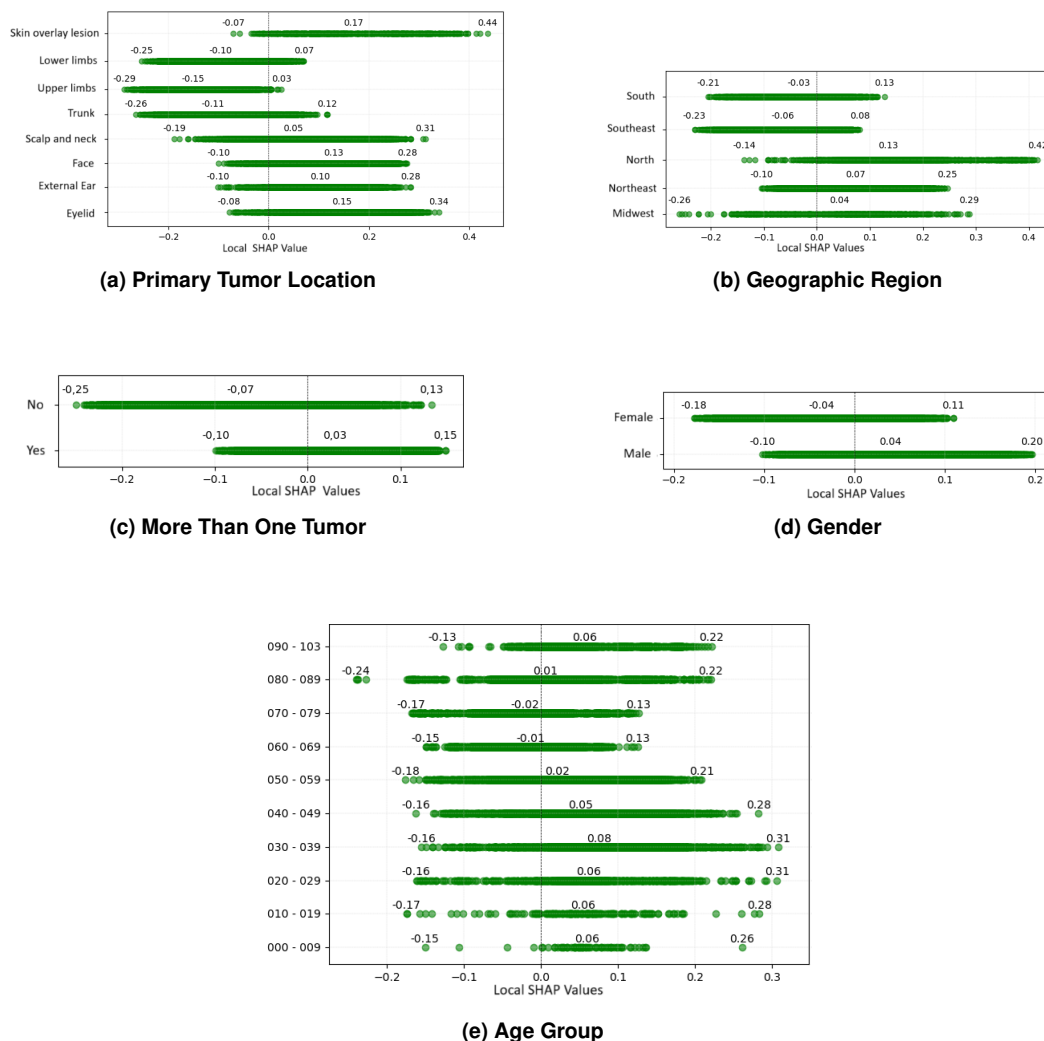


(a) Primary Tumor Location

(b) Geographic Region

(c) More Than One Tumor

(d) Gender

(e) Age Group

**Figure 2. Variation of characteristics contributions**
Source: Compiled by the author with data extracted from the research.

Figure 2a demonstrates the significance of anatomical locations in contributing to the prediction of skin tumor invasiveness. Anatomical sites above the neck, including the eyelid and face, contribute to a higher likelihood of a tumor being classified as invasive, whereas locations such as the lower limbs and trunk contribute to it being predicted as non-invasive. The variability of values indicates the complexity of predicting the nature of the tumor based on its anatomical location. Meanwhile, Figure 2b displays the impact of geographic regions on the model's prediction. Regions such as the Southeast and South are linked to a lower probability of the tumor being predicted as invasive, while the North and Northeast regions increase this chance, especially the North. The Midwest region presents significant variability in these values, contributing to both predictions of invasiveness. Figure 2c compares patients with a single tumor to those with multiple tumors. Patients with one tumor tend to have values indicating a prediction of non-invasiveness,

whereas those with multiple tumors show a slightly higher tendency towards the prevision of invasiveness.

Figure 2d illustrates how gender contributes to the predictive analysis of skin tumor invasiveness. The female gender is associated with a tendency towards contributing to the prediction of non-invasive tumors, albeit with variations that may influence this trend. Conversely, the male gender is more strongly linked to contributions predicting tumor invasiveness. Similarly, Figure 2e highlights the significant role of age groups in determining the classification of skin tumors. Individuals aged between 70 and 103 years exhibit a pronounced contribution towards the prediction of invasiveness, contrasting with those in the 60 to 79 years age bracket, who have a mean contribution leaning towards non-invasiveness. Younger individuals, particularly those up to 49 years, display an opposite tendency, contributing to a higher likelihood of their tumors being classified as invasive, underscoring the nuanced impact of age on the predictive model.

Table 5 highlights the top ten features out of a set of fifty-eight, based on their degree of impact on the model's predictions. These features are ranked according to the average of their individual contributions across all analyzed samples, allowing for a visualization of how each uniquely impacts the model's ability to predict outcomes. The table not only provides an ordered list of these principal features but also serves as a reflection of the interrelationship of the different weights exerted by the SHAP.

### Table 5. Ranking of attribute features on a national scope

|  | Invasive Tumor | Non-Invasive Tumor |
|---|---|---|
| 1st | Skin overlay lesion (*Primary Tumor Location*) | Upper limbs (*Primary Tumor Location*) |
| 2nd | Eyelid (*Primary Tumor Location*) | Trunk (*Primary Tumor Location*) |
| 3rd | Face (*Primary Tumor Location*) | Lower limbs (*Primary Tumor Location*) |
| 4th | External ear (*Primary Tumor Location*) | No (*Occurrence of More Than One Tumor*) |
| 5th | 030 - 039 (*Age Group*) | Female (*Gender*) |
| 6th | None (*Education Level*) | Incomplete Bachelor's Degree (*Education Level*) |
| 7th | 020 - 029 (*Age Group*) | Incomplete Bachelor's Degree (*Education Level*) |
| 8th | 000 - 009 (*Age Group*) | Married (*Marital Status*) |
| 9th | 090 - 103 (*Age Group*) | High School (*Education Level*) |
| 10th | 010 - 019 (*Age Group*) | 070 - 079 (*Age Group*) |

Source: Compiled by the author with data extracted from the research.

## 5. Findings

In Figure 2a, it was found that one of the most relevant factors for predicting a tumor as invasive is the primary tumor's location above the neck. This pattern is corroborated by

the article from [AMERICAN ACADEMY OF DERMATOLOGY 2023], which emphasizes the anatomical complexity of the head and neck region. These areas have a dense network of blood vessels and lymph nodes, increasing the possibility of tumor invasion when it originates in these areas.

Furthermore, the information in Figure 2d indicates that the male gender is more associated with the prediction of invasive skin tumors, and the female gender with the prediction of non-invasive tumors. This aligns with studies, such as the one by [Schwartz et al. 2019], which found a higher incidence and mortality from skin cancer in men than in women. The reasons for these differences may include hormonal factors, behavioral habits, and immune responses, such as men's lower tendency to perform skin self-exams, which can lead to late diagnoses and worse prognoses.

The analysis of the data from Table 4 and Table 5 reveals nuances in the contribution of attributes to the model's prediction that are not immediately evident through traditional attribute explainability analysis by their impacts on the model. This observation suggests an alternative perspective in interpreting the predictive model. For example, while education ranks sixth in the hierarchy of twelve attributes, a specific characteristic from this domain "none", indicating the absence of education, stands out in the same position but within a broader spectrum of fifty-eight characteristics demonstrated in Table 1. This characteristic stands out, for example, in relation to all the characteristics of the attribute referring to the patient's geographic region, which, despite being the second most impactful in the list of attributes, does not have the same degree of impact when considering the breakdown by individual characteristics. This contrast provides a more detailed analysis of the algorithm's behavior in its predictive task.

## 6. Conclusion

This study assessed the impact of attributes at the domain level on the prediction of tumors invasiveness using XAI, employing the SHAP approach. A difference was found between the impact of the attributes on the model and the specific impact of the characteristics within their domains. The findings also highlight important contributions in the field of oncology, underscoring the characteristics with the greatest influence on disease prediction.

Regarding threats to validity, the higher incidence of invasive tumor cases compared to non-invasive ones in the database constitutes an internal threat due to the possibility of generating bias in the predictive model. To mitigate this threat, the selection of the machine learning algorithm was primarily based on its ability to achieve effective Balanced Accuracy, a measure that takes into account the precision of predictions for different states of the dependent attribute. Strategies were also adopted for the tuning of hyperparameters of the algorithms used, in order to mitigate potential biases.

As for future research, it is suggested to expand the application of this study methodology to other types of tumors with high incidence in Brazil, such as breast and prostate tumors, exploring another particularities. The data used in this study were collected on 03/28/2023 and are available at the URL "*https://irhc.inca.gov.br/RHCNet/visualizaTabNetExterno.action*".

# References

Alsinglawi, B., Alshari, O., Alorjani, M., Mubin, O., Alnajjar, F., Novoa, M., and Darwish, O. (2022). An explainable machine learning framework for lung cancer hospital length of stay prediction. *Scientific reports*, 12(1):607.

AMERICAN ACADEMY OF DERMATOLOGY (2023). Melanoma on the head or neck. Acessado em: 17/07/2023.

BRAZIL'S MINISTRY OF HEALTH (2021). Câncer de pele. Accessed on: July 07, 2023.

BRAZIL'S MINISTRY OF HEALTH (2022). Enfrentamento do câncer. Acessado em: 07/07/2023.

Fidler, I. J. (2003). The pathogenesis of cancer metastasis: the'seed and soil'hypothesis revisited. *Nature reviews cancer*, 3(6):453–458.

Ghazal, T. M., Al Hamadi, H., Umar Nasir, M., Gollapalli, M., Zubair, M., Adnan Khan, M., Yeob Yeun, C., et al. (2022). Supervised machine learning empowered multifactorial genetic inheritance disorder prediction. *Computational Intelligence and Neuroscience*, 2022.

Lee, H.-C., Lin, T.-C., Chang, C.-C., Lu, Y.-C. A., Lee, C.-M., and Purevdorj, B. (2022). Clinical risk factor prediction for second primary skin cancer: a hospital-based cancer registry study. *Applied Sciences*, 12(24):12520.

Liu, W.-C., Li, M.-X., Qian, W.-X., Luo, Z.-W., Liao, W.-J., Liu, Z.-L., and Liu, J.-M. (2021). Application of machine learning techniques to predict bone metastasis in patients with prostate cancer. *Cancer Management and Research*, pages 8723–8736.

Lopes, M. C., de Matos Amorim, M., Freitas, V. S., and Calumby, R. T. (2021). Survival prediction for oral cancer patients: A machine learning approach. In *Anais do IX Symposium on Knowledge Discovery, Mining and Learning*, pages 97–104. SBC.

Lundberg, S. and Lee, S. (2017). A unified approach to interpreting model predictions. arxiv: 170507874. *Ar. Xiv*.

Lundberg, S. M. et al. (2020). Shap: Shapley additive explanations. Accessed on: February 12, 2024.

Luo, Y., Tseng, H.-H., Cui, S., Wei, L., Ten Haken, R. K., and El Naqa, I. (2019). Balancing accuracy and interpretability of machine learning approaches for radiation treatment outcomes modeling. *BJR— Open*, 1(1):20190021.

Molnar, C. (2020). *Interpretable machine learning*. Lulu. com.

Santos, M., de Lima, F. C. d. S., Martins, L. F. L., Oliveira, J. F. P., de Almeida, L. M., and de Camargo Cancela, M. (2023). Estimativa de incidência de câncer no brasil, 2023-2025. *Revista Brasileira de Cancerologia*, 69(1).

Schwartz, M. R., Luo, L., and Berwick, M. (2019). Sex differences in melanoma. *Current epidemiology reports*, 6:112–118.

Shapley, L. S. et al. (1953). A value for n-person games.

Silva, W. S., Oliveira, V. T., Araújo, S. S., Vieira, D., and Castro, M. F. (2022). Explainability e auditability: interpretando e validando modelos de machine learning. *Sociedade Brasileira de Computação*.

Sorayaie Azar, A., Babaei Rikan, S., Naemi, A., Bagherzadeh Mohasefi, J., Pirnejad, H., Bagherzadeh Mohasefi, M., and Wiil, U. K. (2022). Application of machine learning techniques for predicting survival in ovarian cancer. *BMC Medical Informatics and Decision Making*, 22(1):345.

Subasi, A., Panigrahi, S. S., Patil, B. S., Canbaz, M. A., and Klén, R. (2022). Advanced pattern recognition tools for disease diagnosis. In *5G IoT and Edge Computing for Smart Healthcare*, pages 195–229. Elsevier.

Taghizadeh, E., Heydarheydari, S., Saberi, A., JafarpoorNesheli, S., and Rezaeijo, S. M. (2022). Breast cancer prediction with transcriptome profiling using feature selection and machine learning methods. *BMC bioinformatics*, 23(1):1–9.

Yu, J., Zhou, Y., Yang, Q., Liu, X., Huang, L., Yu, P., and Chu, S. (2021). Machine learning models for screening carotid atherosclerosis in asymptomatic adults. *Scientific reports*, 11(1):22236.