#### ACCEPTED MANUSCRIPT • OPEN ACCESS

# Interpretable Machine learning model to predict survival days of malignant brain tumor patients

To cite this article before publication: Snehal Rajput et al 2023 Mach. Learn.: Sci. Technol. in press https://doi.org/10.1088/2632-2153/acd5a9

#### Manuscript version: Accepted Manuscript

Accepted Manuscript is "the version of the article accepted for publication including all changes made as a result of the peer review process, and which may also include the addition to the article by IOP Publishing of a header, an article ID, a cover sheet and/or an 'Accepted Manuscript' watermark, but excluding any other editing, typesetting or other changes made by IOP Publishing and/or its licensors"

This Accepted Manuscript is © 2023 The Author(s). Published by IOP Publishing Ltd.



As the Version of Record of this article is going to be / has been published on a gold open access basis under a CC BY 4.0 licence, this Accepted Manuscript is available for reuse under a CC BY 4.0 licence immediately.

Everyone is permitted to use all or part of the original content in this article, provided that they adhere to all the terms of the licence <a href="https://creativecommons.org/licences/by/4.0">https://creativecommons.org/licences/by/4.0</a>

Although reasonable endeavours have been taken to obtain all necessary permissions from third parties to include their copyrighted content within this article, their full citation and copyright line may not be present in this Accepted Manuscript version. Before using any content from this article, please refer to the Version of Record on IOPscience once published for full citation and copyright details, as permissions may be required. All third party content is fully copyright protected and is not published on a gold open access basis under a CC BY licence, unless that is specifically stated in the figure caption in the Version of Record.

View the article online for updates and enhancements.

#### Interpretable Machine Learning Model to Predict Survival Days of Malignant Brain **Tumor** Patients Snehal Rajput<sup>1</sup>, Rupal A. Kapdi<sup>2</sup>, Mehul S. Raval<sup>3\*</sup> and Mohendra Roy<sup>1\*</sup> <sup>1</sup>SOT, Pandit Deendaval Energy University, PDEU Road, Gandhinagar, 382007, Gujarat, India. <sup>2</sup>Institute of Technology, Nirma University, SG Highway, Ahmedabad, 382481, Gujarat, India. <sup>3</sup>School of Engineering and Applied Science, Ahmedabad University, Commerce Six Roads, Ahmedabad, 380009, Gujarat, India. \*Corresponding author(s). E-mail(s): mehul.raval@ahduni.edu.in; mohendra.roy@ieee.org; Contributing authors: snehal.rphd19@sot.pdpu.ac.in; rupal.kapdi@nirmauni.ac.in; Abstract An artificial intelligence (AI) model's performance is strongly influ-enced by the input features. Therefore, it is vital to find the opti-mal feature set. It is more crucial for the survival prediction of the glioblastoma multiforme (GBM) type of brain tumor. In this study, we identify the best feature set for predicting the survival days (SD) of GBM patients that outrank the current state-of-the-art methodologies. The proposed approach is an end-to-end AI model. This model first segments tumors from healthy brain parts in patients' MRI images, ex-tracts features from the segmented results, performs feature selection, and makes predictions about patients' survival days based on selected features. The extracted features are primarily shape-based, location-based, and radiomics-based features. Additionally, patient metadata is also included as a feature. The selection methods include recursive feature elimination (RFE), permutation importance (PI), and finding

 $\mathbf{2}$ 

the correlation between the features. Finally, we examined features' behavior at local (single sample) and global (all the samples) levels. In this study, we find that out of 1265 extracted features, only 29 dominant features play a crucial role in predicting patients' survival days (SD). Among these 29 features, one is metadata (Age of patient), three are location-based, and the rest are radiomics features. Furthermore, we find explanations of these features using post-hoc interpretability methods to validate the model's robust prediction and understand its decision. Finally, we analyzed the behavioral impact of the top six features on survival prediction, and the findings drawn from the explanations were coherent with the medical domain. We find that after the Age of 50 years, the likelihood of survival of a patient deteriorates, and survival after 80 years is scarce. Again, for locationbased features, the SD is less if the tumor location is in the central or back part of the brain. All these trends derived from the developed AI model are in sync with medically proven facts. The results show an overall 33% improvement in the accuracy of SD prediction compared to the top-performing methods of the BraTS-2020 challenge.

**Keywords:** Brain tumor segmentation, feature importance, survival prediction, interpretability

### 1 Introduction

Brain cancer patients' survival rate is lower than other cancer types. The Glioblastoma Multiforme (GBM), or simply, Glioblastoma, is the most invasive and frequently diagnosed type of brain tumor [1, 2]. Due to its infiltrative and diffuse characteristics, the World Health Organization (WHO) has categorized it as a Type-4 tumor [3]. Following the central-brain-tumor registry of the United States (CBTRUS)-2021 report, there were a total of 83,029 deaths in the USA alone between 2014 and 2018 due to malignant brain tumors and other central nervous system disorders (CNS) tumors [2].

#### 1.1 Brain tumor segmentation

Usually, the brain anatomy analysis is done using MRI images, which are non-invasive, and provide high-resolution and detailed information about soft tissues. Recently, deep-learning-based approaches are becoming more popular for segmentation from medical images due to the introduction of powerful GPUs [4]. UNet-based approaches have generated robust segmentation results, as evidenced by their great performance in the medical image segmentation domain [5–7]. Brain Tumor Segmentation (BTS) separates cancerous tissues from healthy tissues, which can further dissect into necrosis, enhancing tumor or edema. In many standard benchmarks, such as in Brain Tumor Segmentation Challenge (BraTS) [8–10] counts the whole tumor (WT), tumor core (TC)

46 47

1 2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18 19

20

21 22 23

24 25

26 27

28

29

30

31

32

33 34 35

36

37

38

39

40

41

42

43

44

- 48 49
- 50 51 52

and enhancing tumor (ET) subregions for the evaluation of the segmentation methods.

The state-of-the-art BTS methods use 2D, 3D, or hybrid UNet [11]. The UNet performance is further improved by assembling attention blocks, [7, 12–15] residual connections between layers [16] and dense connections between layers of the network [6, 7, 17]. In the BraTS–2020, Isensee et al. [18] proposed an improvised version of the 'No-New Network' model. Likewise, in the BraTS-2021 challenge, an optimized version of the same network was proposed at the conference of Medical-Image Computing and Computer-Assisted Intervention (MICCAI) 2021 [19]. The above segmentation techniques suggest that automatic segmentation is a complicated method due to the high variance in structure, shape, location, texture of tumor tissues, lack of ample images in the available standard dataset, and an imbalance between cancerous and healthy tissues. Thus, a robust segmentation method is desirable to develop an accurate and transparent survival prediction system.

### 1.2 Survival days prediction

The Survival Days (SD) prediction is far more complex as it depends on many factors such as accurately segmented brain tumor [20], ample dataset, clinical information such as age, gender, health condition, treatment, biological char-acteristics, and qualitative image properties from radiographic images [21]. Though hugely challenging, it is crucial to improve early diagnosis, treatment planning, and post-treatment analysis of GBM patients [22, 23]. The GBM patients have a dismal survival record, with a median chance of survival of fewer than 12 months [24]. Various studies also show that the survival of pa-tients varies with their age [2, 20, 25]. SD prediction from the BraTS challenge can be further categorized into long-term survival (where survival days are >450 days), mid-term survival (300 to 450 days), and short-term survival (<300 days). Here, accuracy and Spearman ranking coefficient (SpearmanR) are used to evaluate the performance of the models. 

### 1.3 End-to-end methods for BTS and SD

Since both the tumor segmentation and SD prediction are individually complex, therefore, various research groups are trying to develop an end-to-end model by integrating a tumor segmentation with the SD prediction method to make the system smooth and less complicated for SD prediction [26]. In this regard, Mckinley et al. [7] proposed a 3D-2D densely connected encoder-decoder architecture for the segmentation task and thereby extracted the features. An ensemble of linear regressor and random forest classifiers was trained using age and features extracted from BTS to predict survival days. Bommineni et al. [27] proposed four identical networks for segmentation, where networks were trained on each class label and multiple class labels. They used the linear regressor for SD prediction and trained the model on the surface area, volume, spatial location, age, and resection status features.

### **1.4** Interpretability

Usually, BTS and SD models are not tested for their interpretability. In this regard, an end-to-end model that combines automated segmentation, feature extraction, and survival prediction with interpretability is a promising option. For the BTS task, we implemented 3D U-Net [28]. The features were extracted from segmentation results using various wavelet-based, location-based. shape-based, and radiomic-based filters. The radiomic features provide valuable insights into GBM prognosis but will be limited in providing biological insights. The several reasons for these limitations include - tumor heterogeneity, imaging limitation, and, most importantly, the lack of biological context. They can provide insights into the phenotypic structure but cannot explain the underlying molecular processes. Integrating radiomics with genomics, proteomics, or clinical data is necessary for a holistic view. This task is very complex and requires heavy computational resources and expertise. Therefore, the present work examines interpretability from the phenotypic perspective based on publicly available BraTS 2020 challenge data [29].

In addition, we used recursive feature elimination (RFE). Permutation Importance (PI), and correlation matrix to reduce the number of features. Fur-20 ther, we studied the correlation map, Partial Dependency (PD) plots [30, 31], 21 Shapley Additive exPlanations (SHAP) plots [32, 33], and Kaplan-Meier (KM) plots [34] to analyze the predictions. SHAP identifies the most important feature contributing to the prediction. This can aid the clinician in understanding the decision-making process and making treatment-related decisions. PDP will 25 help to visualize how a particular radiomic feature affects prediction across dif-26 ferent patients. This establishes the relationship between radiomic features and prediction and also reveals nonlinear dependencies amongst features. Thus, 28 both can help make informed decisions and offer valuable insights into GBM prognosis.

In summary, our work focuses on the points listed follows:

- Finding an optimal feature set that augurs well for SD prediction.
- Validation of SD prediction on the BraTS-2020 dataset.
- Providing detailed explanations and rationale for the selection of the dominant features set.
- Interpretation of the model behavior and biomedical inference of the top six most important features.

All the obtained results are validated through the BraTS-Challenge-2020 evaluation platform [29].

## 2 Methods

### 2.1 End-to-End approach for SD prediction

The structural diagram of the proposed end-to-end approach is shown in Figure 1. The multiple parametric MRI images are the input to the model such as

16 17 18

13

14

15

19

22 23

24

27

29

30 31

32

33

34 35

36 37

38

39

42

40 41

43 44

45 46

47

48 49

50

T1-weighted (including contrast agent), T2-weighted, and fluid-attenuated inversion recovery (T2-FLAIR) images. The segmentation model is built on 3D U-Net architecture, known as the "No-new Network" [28]. The architecture relies on 3D UNet, which is a well-proven architecture for biomedical segmentation tasks and is robust for tumor segmentation. The network consists of a symmetric five-layered encoder and decoder structure. It is a simple, easy-to-implement architecture with 8.3M parameters. This makes it suitable for the resource-constrained 16GB GPU and 256GB RAM environment while maintaining good segmentation performance on BraTS 2020 dataset. For detailed architecture, please refer to Supplementary Figure A1. For this segmentation model's training dataset. The obtained mean Dice scores for ROIs are 0.819 (Whole Tumor: WT), 0.766 (Tumor Core: TC), 0.702 (Enhancing Tumor: ET) for BraTS2020 training set and 0.880 (WT), 0.858 (TC), 0.759 (ET) for validation set respectively.



Figure 1 The workflow of the proposed end-to-end approach for the SD prediction.

The network and segmentation of the tumorous tissue from the training set are shown in Figure A1(a) of the supplementary section. In addition, Figure A1(b-d) also exhibits a qualitative comparison between the given input (T2-FLAIR) MRI image predicted image and ground truth. The SD predictor model was trained using the dataset's segmented results and ground truth. In contrast, it was tested on the features extracted from the segmented results of the validation set. The feature selection module finds the best group from these extracted features, which are then used to predict survival days. Finally, the SD prediction module is investigated for its decision, understanding its generic (global/overall) and specific (local/sample-wise) behavior on survival days. The details of the feature extraction, a feature selection module, the survival prediction model, and its interpretability are discussed in the subsections below.

### 2.2 Feature extraction module

The feature extraction module obtained the image-based features [25] and radiomics-based features [35] (Table 1 lists the specifics of the features).

Here, the image-based features are extracted by determining the tumor's shape and location. In contrast, radiomics-based features are extracted from necrotic and non-enhanced tumor regions using wavelet and Laplacian of Gaussian (LoG) filters (with  $\sigma$  value 1 to 5). Here, the lower value of  $\sigma$  highlights fine textures, and the higher  $\sigma$  focuses on coarse textures. The wavelet filters denoise the images and capture spatial and global signals [36]. The LoG filter pinpoints the blob centers and approximates its size, shape, and orientation [37]. Thus, we obtained 1264 features (1225 radiomics-based + 39 image-based). We also considered the metadata, e.g., the Age of patients, as a feature. As a result, 1265 features in total are being taken into account for the evaluation. Since some of these features can be redundant or not contribute to the prediction, a feature selection procedure is essential.

	Image-based features	
Shape-based features (27)	Surface area of ROIs, the volume of ROIs, proportion of ROIs, proportion ratio between each ROI, the area-to-volume ratio of ROIs, and amount of tumor.	
Location-based features (12)	Centroid of ROIs, the distance between the center of ROIs and the center of the brain.	
	Radiomics-based features	
Shape features (13)	Elongation, major axis length, least axis length, mesh volume, flatness, maximum diameter row, maximum diameter column, surface area, sphericity, and surface volume ratio.	
First-order statistical features (144)	Energy, maximum intensity value, minimum intensity value, mean, entropy, absolute deviation, inter-quartile range, vari- ance, skewness, percentile, kurtosis, uniformity, and median.	
Gray-level features (1068)	Neighboring gray-tone difference matrix (NGTDM), Gray- level co-occurrence_matrix (GLCM), Gray-level.size-zone (GLSZ), Gray-level run-length matrix (GLRLM), and Gray- level_dependence_matrix (GLDM).	

 Table 1
 Feature-set lists 1264 features (1225 radiomics based + 39 image based).

Note: The values within the parenthesis represent the Number of features extracted.

### 2.3 Feature selection module

The primary goal of feature selection methods is to eliminate unimportant or repetitive features. Here, we employed Recursive Feature Elimination [38] and

1 Permutation Importance [39] as feature selection methods. RFE is a back-2 3 ward feature selection method that re-fits the model after iteratively ranking the features according to their importance and eliminating the least impor-4 tant features. The description of the chosen dominant features identified by 5 RFE are shown in supplementary table A1. On the other hand, Permuta-6 tion importance finds influence in the model score by randomly re-arranging 7 a single feature value. The pseudo-code of PI is shown in algorithm 1. This 8 technique breaks the connection between the desirable feature and the out-9 put feature. The model's score decline demonstrates how largely it depends on 10 that feature. Thus, we weighted the features according to their importance. In 11 general, dominating features are given greater weight than other features. Zero 12 or negative weights indicate no contribution of the feature for the prediction. 13 Therefore, we removed them, bringing the set down to 180 features. These 14 180 features were further examined using the Spearman correlation coefficient 15 (SpearmanR) with an absolute cutoff value of 0.5. It is clear from the corre-16 lation values that the necrotic, active, and whole tumor centroids are firmly 17 connected, given that they have similar characteristics in common. As a result, 18 we narrowed the set of features to 29 by eliminating redundant features. 19 20

The pipeline for feature selection is as follows:

- We eliminated the features based on the Permutation Importance weights (which define their contribution to the outcome). The threshold value of the weights is 100. Any features with PI weight <100 are eliminated. This results in 180 prominent features.
- Further, we eliminate the weaker features from these 180 features by finding the SpearmanR and a sorting process. For this, (a) we take features one by one (from the 180 feature set), starting with the feature having the least PI weight, and find its SpearmanR with the rest of the 179 features, (b) then we select the features which are having correlations less than 0.5. (c) then from this selected features, we identify the feature which is having highest PI weight value and use it to replace the feature that is having least PI weight (that we chose in step (a)). This process is repeated for each feature in the 180 feature set. That means the loop will run 179 times. Lastly, we find 29 dominant features (having less correlation) out of 180 features.

A detailed description of the selected dominant features using PI is shown in supplementary table A2.

### 2.4 Survival prediction module

The Random Forest Regressor (RFR) [40] is based on ensemble learning, where decision trees (DT) are fundamental building blocks. Each DT was created using random samples from the training set; hence it is called a Random Forest. This method is widely used as it has been proven accurate and robust [40] across multiple complex problems, including SD prediction [41, 42]. The RFR model is often more successful than other models because the outcomes obtained by averaging the prediction from each tree result in lower variability.

47 48 49

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37 38

39 40

41

42

43

44

45

- 50
- 51 52

#### AUTHOR SUBMITTED MANUSCRIPT - MLST-100882.R1

Alg	orithm 1 Permutation Importance (PI) Algorithm:
	Input: Trained model $m$ on the Dataset $D$
	<b>Compute:</b> The metric S of the model m on dataset D (for instance $R^2$
	metric for a regressor model)
	for each feature $j$ : do
	for each repetition $k$ in $K$ : do
	Arbitrarily re-arrange column $j$ of Dataset $D$ to produce a noisy
	variant of the dataset say, $\hat{\mathbf{D}}_{k,j}$ .
	Measure the metric $s_{k,j}$ of model $m$ on variant Dataset $\hat{\mathbf{D}}_{k,j}$ .
	end for
	Measure importance $I$ of each feature $j$ defined as:
	$I_j = s - \frac{1}{k} \sum_{k=1}^k s_{k,j}$
	end for

Additionally, randomization during tree growth and splitting helps prevent overfitting [43]. Hence, the RFR model is robust for predicting brain tumor patients' survival [44]. Here, a five-fold cross-validation technique was used to train the RFR model. Also, the hyper-parameters of the model were fine-tuned using grid search. The fine-tuned parameters are the maximum tree depth, maximum number of features at each split, number of trees, and the minimal sample size required to be at a child node at a split point.

#### 2.5 Interpretability methods for the proposed SD module

Understanding the decisions taken by AI or Machine learning (ML) models is essential. Especially in the medical domain, the interpretability of such an AI model is vital to increase its reliability. Generally, the non-linearity in an AI model makes them hard to decipher. That is why we use model-agnostic methods like SHAP [32, 33] and PDP [30, 31] to find the interpretability of the proposed model.

The primary objective of the SHAP method is to determine how much each feature impacts the prediction for a given instance. The SHAP-value of a feature is the average marginal contribution of that feature to the value of the predecessor set among all possible permutations of the feature set. It can be expressed as in equation 1 [45].

$$(\Phi_j) = \frac{1}{|\Pi(N)|} \sum_{\pi \ \epsilon \ \Pi(N)} \underbrace{(v(\hat{P}_j^{\pi} \cup j) - (v(\hat{P}_j)))}_{(m \ \epsilon \ \Pi(N)}$$
(1)

where,  $(\Phi_j)$  is the SHAP-value of feature of interest j,  $\Pi(N)$  is the possible coalitions of all feature set,  $\pi$  is the specific coalition, feature of interest is j, v is contribution of feature(s),  $(\hat{P}_j^{\pi} \cup j)$  is the predecessor set of feature j in

a particular coalition, including the j feature whereas  $\hat{P}_i$  is predecessor set of feature j in a particular coalition, excluding j feature. E.g., if  $\pi = \{A, B, D\}$ , j = B and  $v\{A\} = 8$ ,  $v\{B\} = 10$ ,  $v\{C\} = 9$ ,  $v\{A, B\} = 18$ ,  $v\{A, D\} = 10$ 20,  $v\{B, D\} = 22$  and  $v\{A, B, D\} = 25$ , whereas the possible predecessor sets in this example) in a particular coalition  $\pi = \{A, B, D\}$  :  $\{\phi, A\}$  and marginal contribution (MC) of j(=B) is calculated as:  $v\{A, B\} - v\{A\} = 20 - 8 = 12$ . Further, calculating the MC of feature j across all the possible coalitions and averaging will give us a SHAP value  $(\Phi_B)$  of feature B. In summary, it shows each feature's influence on predicting survival days. It helps to understand the global behavior of the model by combining the explanation of each sample (Please see the Supplementary Table A4 for a more detailed explanation of this example). Algorithm 2 displays the pseudo code to find the SHAP value for a feature. 

The PDP displays the global effect of the feature on the target. The PDP considers all the samples and can show and examine the global association between survival days and input variables. The partial dependence function is represented as :

$$f(x_s) = E_c[f(x_s, x_c)]$$
(2)

where  $x_s$  are the desirable feature(s) for which we want to plot partial dependency and  $x_c$  are the remaining features used to train the model.  $x_c = x'_s$ and  $X = x_s + x_c$  is the whole feature set. In PDP, we assume that feature subset  $x_s$  and  $x_c$  are uncorrelated to each other and hence can be calculated using average interaction effect [31] as:

$$f(x_s) = \frac{1}{n} \sum_{i=1}^{n} f(x_s, x_e)$$
(3)

Algorithm 3 displays the pseudo code to find the samples' Partial dependency (PD) values.

#### Algorithm 2 Calculating SHAP -value for a feature:

**Input:** Number of feature N and their respective real value v signifying their contribution. The contribution vector v of a particular feature is calculated through perturbation feature values of coalition  $\pi$ . More details can be found here [46]. k is the number of sampling permutations

<b>Output:</b> SHAP value $\phi_i$ for the feature $j \in N$ .
for Iteration : 1, 2, $\mathbf{K}$ : do
Randomly select $\pi$ from set of all permutation $\Pi(N)$
$\mathbf{for}\; j \in N: \mathbf{do}$
Calculating predecessor set $P_i^{\pi} = \{j \in N \mid \pi(j) < \pi(i)\}.$
$\phi_j=\phi_j+rac{v(\hat{P}_j^\pi\cup j)-(v(\hat{P}_j))}{K}$
end for
end for

#### AUTHOR SUBMITTED MANUSCRIPT - MLST-100882.R1

**Algorithm 3** The steps of obtaining PD Value of samples are:

	<b>Input:</b> The unique feature's values $x_s = x_1, x_2,, x_n$ , where x is feature
	of interest
	<b>Ouput:</b> PD plot of desirable feature.
	for $i \in (1, 2)$ k): do
	Beplace the original $r_{i}$ values with the constant $r_{i}$ in the training
	samples
	computes the predicted value vector from the altered copy of the
	training samples.
	compute the average of the prediction to find $f'(x_{1i})$ .
	end for
	The PDP for $x_1$ is obtained by plotting the pairs $\{x_{1i}, f'(x_{1i})\}$ for $i =$
	1, 2, n
2.6	Performance Metrics
Usi	ng multiple metrics for the performance evaluation provides the robust-
nes	s information of the employed model. Hence, we quantified our model
pre	the function of widely used metrics for survival prediction, such as accuracy
[28]	42], mean squared_error (MSE) [28, 42], median squared_error (medianSE)
ran	$[42]$ , standard-deviation standard_error (studil) [26, 42, 47], spearman
ran	
2.7	Dataset BraTS-2020
The	
the	e training BraTS 2020 [8–10] dataset includes 369 3D MRI samples for
viv	e training BraTS 2020 [8–10] dataset includes 369 3D MRI samples for segmentation and metadata (resection status information, Age, and sur-
1 1 11	e training BraTS 2020 [8–10] dataset includes 369 3D MRI samples for segmentation and metadata (resection status information, Age, and sur- l days). Out of this, 236 patients' metadata are provided for the SD
pre	e training BraTS 2020 [8–10] dataset includes 369 3D MRI samples for segmentation and metadata (resection status information, Age, and sur- l days). Out of this, 236 patients' metadata are provided for the SD liction task. The validation BraTS 2020 dataset contains 125 MRI sample
pre ima	e training BraTS 2020 [8–10] dataset includes 369 3D MRI samples for segmentation and metadata (resection status information, Age, and sur- d days). Out of this, 236 patients' metadata are provided for the SD liction task. The validation BraTS 2020 dataset contains 125 MRI sample ges and metadata of 29 patients. Each sample instance includes the Fluid-
pre ima atte	e training BraTS 2020 [8–10] dataset includes 369 3D MRI samples for segmentation and metadata (resection status information, Age, and sur- al days). Out of this, 236 patients' metadata are provided for the SD diction task. The validation BraTS 2020 dataset contains 125 MRI sample ges and metadata of 29 patients. Each sample instance includes the Fluid- nuated-inversion recovery (T2-FLAIR), T2-weighted MRI preoperative
pre ima atte ima	e training BraTS 2020 [8–10] dataset includes 369 3D MRI samples for segmentation and metadata (resection status information, Age, and sur- al days). Out of this, 236 patients' metadata are provided for the SD diction task. The validation BraTS 2020 dataset contains 125 MRI sample ges and metadata of 29 patients. Each sample instance includes the Fluid- enuated-inversion recovery (T2-FLAIR), T2-weighted MRI preoperative ges, T1 weighted (T1), post-contrast T1-weighted (T1-ce), and correspond-
pre ima atte ima ing	e training BraTS 2020 [8–10] dataset includes 369 3D MRI samples for segmentation and metadata (resection status information, Age, and sur- al days). Out of this, 236 patients' metadata are provided for the SD diction task. The validation BraTS 2020 dataset contains 125 MRI sample ges and metadata of 29 patients. Each sample instance includes the Fluid- nuated-inversion recovery (T2-FLAIR), T2-weighted MRI preoperative ges, T1 weighted (T1), post-contrast T1-weighted (T1-ce), and correspond- ground truth. In addition, the dataset is skull-stripped, aligned to the
pre ima atte ima ing ider	e training BraTS 2020 [8–10] dataset includes 369 3D MRI samples for segmentation and metadata (resection status information, Age, and sur- al days). Out of this, 236 patients' metadata are provided for the SD diction task. The validation BraTS 2020 dataset contains 125 MRI sample ges and metadata of 29 patients. Each sample instance includes the Fluid- muated-inversion recovery (T2-FLAIR), T2-weighted MRI preoperative ges, T1 weighted (T1), post-contrast T1-weighted (T1-ce), and correspond- ground truth. In addition, the dataset is skull-stripped, aligned to the trical anatomical structure, and re-sampled to an isotropic resolution. The mertation dataset is the BraTS 2020, are labeled for the dataset
pre ima atte ima idei seg	e training BraTS 2020 [8–10] dataset includes 369 3D MRI samples for segmentation and metadata (resection status information, Age, and sur- al days). Out of this, 236 patients' metadata are provided for the SD diction task. The validation BraTS 2020 dataset contains 125 MRI sample ges and metadata of 29 patients. Each sample instance includes the Fluid- enuated-inversion recovery (T2-FLAIR), T2-weighted MRI preoperative ges, T1 weighted (T1), post-contrast T1-weighted (T1-ce), and correspond- ground truth. In addition, the dataset is skull-stripped, aligned to the netical anatomical structure, and re-sampled to an isotropic resolution. The nentation class labels, as defined in the BraTS-2020, are label-0 for back- and words, label 1 for pagentia and non enhanced tumor words (NCP) are
pre ima atte ima idei seg: gro	e training BraTS 2020 [8–10] dataset includes 369 3D MRI samples for segmentation and metadata (resection status information, Age, and sur- al days). Out of this, 236 patients' metadata are provided for the SD diction task. The validation BraTS 2020 dataset contains 125 MRI sample ges and metadata of 29 patients. Each sample instance includes the Fluid- enuated-inversion recovery (T2-FLAIR), T2-weighted MRI preoperative ges, T1 weighted (T1), post-contrast T1-weighted (T1-ce), and correspond- ground truth. In addition, the dataset is skull-stripped, aligned to the niccal anatomical structure, and re-sampled to an isotropic resolution. The nentation class labels, as defined in the BraTS-2020, are label-0 for back- ind voxels, label-1 for necrotic and non-enhanced tumor voxels-(NCR or $\Gamma$ ) label 2 for adama words ED, and label 4 for anhancing tumor words.

In this work, the prediction model was evaluated through the BraTS evalua-

tion platform [29]. In addition, we have used the BraTS-2020 top-performing

ET.

Results and Discussion

#### 

models as benchmarks to compare our results. Finally, this section discusses the results of the proposed end-to-end model for its performance and interpretability.

### 3.1 Correlation study of dominant features

To gain a better insight, we have plotted the correlation matrix of the features as shown in Figure 2 (Refer to Supplementary Table A3 for the annotation of the features). The plot shows that most features are highly uncorrelated, which signifies that they have captured distinct properties of phenotypes. Furthermore, the histogram in Figure 2 validates that most of the selected features correlate from -0.13 to +0.17, suggesting they are uncorrelated, and it justifies the merits of our chosen features.



Figure 2 Correlation matrix of feature-set obtained through PI method. The histogram plot on the right-hand side depicts the range and the count for all the correlations in the heat map. (Refer Supplementary Table A3 for features annotation)

### 3.2 Survival days (SD) prediction results

The comparison of our survival days (SD) prediction results with top-ranking methods of BraTS 2020 are shown in Table 2. A robust method must perform

(4)

well on multiple performance metrics apart from accuracy, as each quantifies the models on different parameters. Hence, we compared the proposed model with benchmark models [7] and reported the improvement as computed using equation 4. Here the percentage of improvement  $\phi$  for each performance metrics x for our proposed model P given by:

$$\phi(x) = \frac{Proposed\_model(P) - Top\_ranking\_model(S)}{Top\_ranking\_model(S)} \times 100$$

With this, the survival prediction result of the proposed method shows a 33.33% improvement in accuracy. There is a 19.13% improvement in MSE, which measures the variance around the fitted regression and indicates the deviation of model prediction from the actual one. However, it is sensitive to outliers [49]. In the case of median SE, there is a 60.80% improvement, which uses the median value of the residuals and is unaffected by the outliers. All these results obtained using various metrics indicate the robustness of the prediction [49]. We can see a 2.62% improvement in stdSE and a 181.03% improvement in SpearmanR coefficient often used to measure the relation between the therapy response and the survival days [48]. As shown in Table 2, our model has performed consistently in all the standard metrics used for SD prediction.

Table 2SD performance comparisons with top-ranking models on the training and<br/>validation datasets BraTS 2020. The numbers of other models are obtained from the<br/>validation leader-board [29]. NA: Not Available

Dataset	Method	Accuracy	MSE	medianSE	$\mathbf{stdSE}$	$\mathbf{SpearmanR}$
	Mckinley et al. [7]	NA	NA	NA	NA	NA
<b>T</b>	Asenjo and Solís et al.[15]	0.822	55499.71	11351.02	147319.00	0.833
Training	Bommineni et al.[27]	NA	NA	NA	NA	NA
	Ali et al.[50]	0.641	62305.61	05745.64	200788.00	0.632
	Proposed Method	0.538	60668.61	16037.10	125873.00	0.754
	Mckinley et al.[7]	0.414	098704.66	36100.00	152176.00	0.253
Validation	Asenjo and Solís et al.[15]	0.520	122515.80	70305.26	157674.00	0.130
validation	Bommineni et al.[27]	0.379	093859.54	67348.26	102092.00	0.280
	Ali et al.[50]	0.483	105079.40	37004.93	146376.00	0.134
	Proposed Method	0.552	79826.24	14148.89	148288.00	0.711

### 3.3 Interpretability of SD prediction model

This section presents a detailed analysis of the influence of features on SD prediction. The SHAP results provide local and global impact details, whereas PDP helps analyze features' global impact.

### 3.3.1 SHAP analysis results

SHAP depicts the importance of features in predicting a sample by calculating SHAP values. It shows the contribution of features to the expected prediction



**Figure 3** SHAP summary plot of dominant features: Each blue dot represents a single patient. The X-axis shows the Shapley (SHAP) value of a specific patient, which signifies the feature's effect on the survival days for the particular patient. The absolute higher SHAP value indicates a higher impact on survival days, where a sign indicates increasing or decreasing average survival days. The Y-axis shows features arranged according to importance (high to low) calculated through the average of absolute SHAP-value. On the opposite side, feature value shows high (red) or low (blue) values.

among all the feature combinations. The SHAP value shows how much a single feature affects the forecast, whereas the signs indicate whether the impact is positive or negative on the prediction outcome. Figures 3 and 5 show the SHAP summary and waterfall plots, respectively. The SHAP-summary plot helps us to visualize the global (generalize) and local (as it plots for every sample) impact of features on the model. In contrast, the waterfall plot allows us to visualize and study the features' impact on an individual sample. It will enable us to explore the role of features and their value on the particular prediction, where we can minutely examine each feature behavior for any desirable sample. In the SHAP-summary plot, X-axis displays the SHAP value, which signifies the impact of features on the target feature (Here, the target feature is the *survival days*). The greater the value (absolute), the more significant the effect on the target component, whereas the sign (+/-) indicates whether that impact is positive or negative. In the Y-axis, features are listed in the order

of importance (from top to bottom). Each point on the summary plot represents a sample, and the point's color represents the value of the corresponding instance. Here, blue denotes a low feature value, and red a high one.

From Figure 3, we observe that *Wavelet-LLL\_firstorder\_InterquartileRange* (WIR) feature has the highest importance. It is a first-order radiomic feature extracted using the wavelet low pass filter and depicts the distribution of specific pixel values. WIR measures the pixel intensity between the 25% to 75%percentile range. From the plot, we can observe that the samples with intermediate or high feature values (purple and red color) of WIR contribute positively. to the prediction, which has a maximum positive SHAP value. In other words, the intermediate or high feature value of the WIR feature increases the SD of patients. It is also apparent that there is an aggregation of large samples (with blue color) within the SHAP value range of -15 to -25 (refer to the WIR feature row listed on the Y-axis). It signifies that the majority of the samples fall into this SHAP range. The samples within this range are responsible for reducing patients' SD. Also, it shows that tumor intensity (pixel value) information falls within this range, reducing the SD. It signifies that the intensity of pixels of a tumor in an MRI plays a significant role. Both Aboussaleh et al. and Bae et al. mention this fact [51-53].

The 2nd most crucial feature is Age. From Figure 3, it is clear that samples with the lower feature value of Age have positive SHAP values. In other words, the lower Age increases the SD of patients. This observation aligns with medical science inference, i.e., the Age of GBM patients is crucial in determining SD, i.e., the lower the Age, the more the survivability [54].

The 3rd most crucial feature is the  $cent_wb_x$  shown in Figure 3. It is a 26 location-based feature representing the centroid coordinate of a whole tumor 27 along the X-axis of an MRI image (a physical coordinate). The plot shows that this feature negatively impacts prediction with intermediate and higher 29 feature values. That means the higher feature value is responsible for reducing 30 the survival days of patients. Here, the X-axis represents the axial view [55], 31 and higher feature values represent the physical coordinates of the central part 32 of the brain. This plot signifies that tumors in the brain's central and latter-33 mid parts will reduce patients' survival days [56]. Similar resemblance can be 34 observed for  $cent_at_x$  and  $cent_nec_x$  features, which are centroid of active 35 tumor and centroid of necrosis, respectively. 36

The 4th most important log-sigma-2-0feature is the mm3D\_firstorder\_Kurtosis (LFK). It is a first-order radiomic feature extracted using a LOG filter, which signifies the distribution of voxels without considering their spatial relations [57]. This feature measures the tailedness (outliers) of data distribution. From the plot, we can observe that low kurtosis values are increasing SD, and higher kurtosis values are reducing the SD of patients. Most samples fall within the SHAP- value range of -20 to 60.

The 5th most crucial feature is log-sigma-2-0-mm-3D\_glcm\_Correlation. It is a second-order radiomic feature extracted using a LOG filter, which measures the inter-relationship of intensity between neighbor voxels [57]. The plot shows

14 15

1

2 3

4

5

6

7

8

9

10

11

12

13

16

19

20

21

22

23

17 18



28

37

38

39

40

41

42

43

44

45

that higher feature values are responsible for increasing SD, and lower feature values reduce SD. In other words, the higher correlation between voxels value increases SD, and low correlation values reduce SD.

The 6th most important feature is wavelet- $HHH_firstorder_Kurtosis$ . It is a first-order radiomic feature extracted using a wavelet filter that uses highpass filters in the series of z, y, and x directions. The distribution of voxels is independent of their spatial relations, similar to the 4th most important feature. Here, the plot shows that lower feature values are responsible for increasing SD (for more information, see Figure 4 for SHAP-distribution plot).

In summary, comparing all the features, we can say that the range of SHAP values for all the features is -40 to +40 (in X-axis). Also, with the decreasing of feature importance, the range of SHAP value decreases. That means the features with a low SHAP range have a lower impact on the SD.

Note: Most samples and their SHAP-value can also be verified through Figure 4, which shows the respective features' SHAP-value and feature value distributions.



**Figure 4** SHAP value and its distribution for the top six dominant features, say (a) WIR, (b) Age, (c) cent\_wb\_x, (d) log-sigma-2.0-mm-3D\_firstorder\_Kurtosis, (e) log-sigma-1-0mm-3D\_glcm\_Correlation, (f) wavelet-HHH\_firstorder\_Kurtosis. X-axis displays the feature value of respective features whereas Y-axis displays the SHAP-value of the respective feature. Each blue dot indicates a specific sample. And bars show histograms of feature value. The rest of the features SHAP value distribution can be found in Supplementary Figure A3.

Again the SHAP-waterfall plot provides the visual interpretation of features contribution for a single prediction. Figure 5 is a SHAP-waterfall plot for a single sample. The average SD is shown on the X-axis, and the features are

arranged on the Y-axis in descending order according to their SHAP values (from top to bottom). From this plot we can analyze, how much the features are impacting negatively (blue) or positively (red) and thereby shift the prediction from the expected outcome E[f(x)]. The expected outcome is the average of all the outcomes for all the samples. We observe that for our example sample (for which the plot is generated), the model output is f(x) = 331.732. The expected output is E[f(x)] = 478.91. This deviation in the model outcome can be understood by quantifying the influences of each of the features.

The SHAP value of each feature in Figure 5 depicted this quantification. By adding all of the SHAP values from each feature, it is possible to determine how much each feature (N) contributed to the model output. This is given by  $f(x) = E(f(x)) + \sum_{N} SHAP$ . Here the  $\sum_{N} SHAP$  represents the sum of the SHAP value of all the features. From this analysis also we can see that the feature Age is having a higher impact on the model outcome. For this sample, the Age value is 71.37 and it reduces the average survival days by 39.33 days (- (minus) value indicates a reduction in SD). Similarly,  $cent_wb_x$  value is 164.651, which is also reducing SD by 28.22 days. Whereas mapping Age,  $cent_wb_x$  features to SHAP-summary Plot (Figure 5) or SHAP-distribution plot Figure A3 which shows global impacts. We can extract similar observances of reducing SD for these features. For e.g., visualizing Age, cent\_wb\_x feature on SHAP-summary Plot, which shows a higher value of these features are reducing SD. Similarly, visualizing Age, cent\_wb\_x feature on SHAP-distribution plot Figure A3 also shows a reduction in SD. This proves that features show the same behavior both globally and locally. 

Further, more information was derived by combining SHAP summary (Figure 3), SHAP-distribution plot (Figure 4), and PDP of the top 6 dominant features (Figure 6) which is explained in Section 3.3.2.

#### 3.3.2 PDP analysis

A PDP shows a marginal effect between desirable and target features (Survival Days) [30]. It shows how a dependent variable changes when an explanatory variable changes, provided all other variables remain constant. If changing the value of a particular feature creates more variation in the average survival days indicates that the feature is crucial. In this analysis, we consider the top six features according to their importance (with respect to their absolute SHAP value). These are WIR, Age, cent\_wb\_x, LFK, log-sigma-1-0-mm-3D\_glcm\_Correlation, and wavelet-HHH\_firstorder\_Kurtosis. The PDPs of the dominant six features are shown in Figure 6 (and plots of the rest of the features are shown in Figure A2 Supplementary section). The PDPs were arranged in the order of importance (higher to lower) obtained through the SHAP summary plot (Figure 3).

Furthermore, visualizing the PDPs, we found the marginal impacts are in line with the order of importance of features that supplements SHAP-value analysis. The detailed analysis of the top six features is: The marginal effect of WIR feature on the survival days is shown in Figure 6. The trend shows

#### AUTHOR SUBMITTED MANUSCRIPT - MLST-100882.R1



**Figure 5** SHAP-waterfall plot of dominant features. This plot shows the features' impact on a particular instance. Each row in the Figure shows how the negative (blue) or positive (red) contribution of each feature shifts the value from the expected output (E[f(x)]) to the model predicted value (f(x)). f(x) = E(f(x))+ shapley value of each feature. Value besides every feature row shows their respective observed value.

value sharply increases within the range 100 to 300, reduces within the 300-350 range, and remains saturated within 350-800 intensity value. It indicates that intensity heterogeneity is very high in the range from 100 to 300, which causes a sharp increase in marginal impact. This suggests that intra-tumor tissues are highly heterogeneous. Comparing PDP (Figure 6(a)) with SHAP (Figure 3) and its distribution plots (Figure A3 (a)), we can conclude that this intensity range between 100-300 is decreasing SD (testified through a decrease in SHAP value). Hence we can conclude that tumor pixel intensity within this range is detrimental to a patient's survival. Also, as mentioned in this study, higher tumor heterogeneity is associated with increased malignancy [58]. This also complies with the other studies, which suggest that wavelet filters help capture enhanced texture features [58, 59].

Similarly, from the PDP of the Age feature (Figure 6(b)), we can observe the trend of marginal effect, which shows a maximum deviation in marginal effects for lesser Age patients, signifies maximum impacts on SD. Further, the marginal effect reduces with the increasing Age of patients. Comparing SHAP

46 47 48

26

27

28

29 30 31

32

33 34

35

36

37

38 39

40

41

42

43

- 49
- 50 51
- 52





**Figure 6** PDP analysis of the top six features. This plot depicted the marginal effect of the (a) WIR, (b) Age, (c) cent\_wb\_x, (d) log-sigma-2.0-mm-3D\_firstorder\_Kurtosis, (e) log-sigma-1-0-mm-3D\_glcm\_Correlation, (f) wavelet-HHH\_firstorder\_Kurtosis on the survival days. Here, X-axis shows values of the respective feature, whereas Y-axis shows the average rate of change (marginal-impact) respective feature value creates on the target feature. The vertical bar on the X-axis shows most of the samples' distribution. The rest of the features for PDP can be found in Supplementary Figure A2.

summary and SHAP-distribution plots, we can observe that after 60 years of Age, there is a decrease in SD (as there is a decrease in SHAP-value beyond this range). Whereas, the PD plot for the  $cent\_wb\_x$  feature (shown in Figure 6 (c)) is the physical coordinate of the whole tumor. The plot shows that the marginal effect is more significant if the centroid is within the range of value, i.e., 75-112 (approx.), and less significant for the 113-160 range of value. Also, comparing these ranges to the SHAP distribution plot, we can observe the former range of values is increasing the SD and the latter is reducing the SD, which signifies tumor lesions in the central or latter part are detrimental to patients.

At the same time, the log-sigma-2.0-mm-3D\_firstorder\_Kurtosis feature is a radiomic first-order statistical information that measures the peakedness of data distribution. For a normal distribution, kurtosis (k) is 3. If k > 3, the dataset tends to have significant outliers. If k < 3, the dataset has fewer or no outliers. PD plot (Figure 6(d)) shows for k = 3; there is a higher marginal impact on SD. Comparing PDP with SHAP and SHAP-distribution plots, we can conclude that, for most samples, k is 3, and it increases SD. At the same

2 3

4

5

6

time, there are enough samples with k > 3, decreases SD. It signifies that there are considerable amounts of outliers or intra-heterogeneity among samples. As mentioned by Steven et al. [60], diffusion kurtosis imaging works on a similar principle of capturing non-normal distribution behavior, which signifies tissue heterogeneity. It is observed that the survival days are positively skewed [61].

The  $log-sigma-1-0-mm-3D_{-}$  glcm\_Correlation is a radiomics feature that 7 calculates the joint likelihood of occurrence of the given pixel pairs with the 8 specified intensity value. At the same time, the grav-level co-occurrence ma-9 trix explores spatial relationships between pixels at a specific distance and 10 direction. From the PDP (Figure 6(e)), we can observe that if the pixel pairs 11 correlation is more than or equal to 0.6 value, it impacts the SD more. Similarly, 12 observing the correlation threshold of 0.6 in the SHAP and SHAP-distribution 13 plot, one can observe that it impacts positively (having a positive SHAP value). 14 Also, it is mentioned by Sanghani et al. in their study [62], which shows texture 15 features played a crucial role in SD prediction. 16

Further, wavelet-HHH\_firstorder\_Kurtosis is a first-order statistical feature 17 like log-sigma-2.0-mm-3D\_firstorderKurtosis (the 4th most feature), but it is 18 extracted using wavelet filters. Comparing their PDPs (Figure 6(d) and (f)) 19 shows both capture kurtosis information but in different dimensions. From the 20 PDP (shown in Figure 6(f)), we can observe that the kurtosis value is steeply 21 increasing between 0 to 100 and then stagnant for the rest of the values. 22 Further, observing SHAP- distribution, we can conclude that most samples are 23 in this range, and samples near 1-10 values are decreasing the SD, and the 24 rest are increasing the SD. However, some samples are sparsely distributed, 25 signifying that they are outliers. 26

All these signify the importance of these features in determining the SD. 27 With the above analysis, we find the Age, WIR, cent\_wb\_x, LFK, log-sigma-28  $1-0-mm-3D_glcm_Correlation$ , and wavelet-HHH\_firstorder\_Kurtosis plays a 29 crucial role in determining a patient's SD. Similarly, we can analyze other re-30 maining features. Finally, we agree that the WIR feature could tell us about 31 tumor heterogeneity associated with high malignancy. Again, the Age feature 32 showed us the trend of survivability, where the survival chances decrease with 33 the patient's increasing age (this is further validated by the Kaplan-Meier 34 (KM) [34] plot as shown in the Supplementary Figure A4). At the same time, 35 the centroid of tumors enabled us to locate tumors in the central or latter-36 central part, which are detrimental for patients. All these analyses using the 37 SHAP and PDPs are analogous to medical findings and related studies. This 38 signifies the model's reliability and validates the explainability methods such 39 as SHAP and PDP. 40

45 46 47

- 48
- 49 50
- 51 52

### 4 Limitations of the proposed approach and future prospect

The SHAP and PDP techniques are the post-hoc methods that interpret the model after the completion of training. However, for the further understanding of a model, the study of the intrinsic characteristics may help to an extent. Functional imaging like PET, fMRI, and MRS can provide insights into GBM by capturing molecular or physiological information not captured by normal MRI or CT Scans. The methods like the Neural Ordinary Differential Equation model (NODE) can provide the learning behaviour of a model, especially to understand the spatiotemporal deep feature extraction of a segmentation model [63]. Further, the diffusion imaging modalities such as diffusion kurtosis Imaging [64] may help us to understand the underline biological and pathological characteristics of GBM. However, these kind of functional imaging are more complex to analyze, has a high variability across imaging sessions, are more susceptible to noise, and are also expensive. In short, they face several challenges for routine GBM prognosis [65, 66]. Still we believe, integration of these modalities with conventional MRI techniques will enhance the understanding of GBM with added model transparency and interpretability.

### 5 Conclusion

We have proposed an end-to-end approach for the SD prediction task. We have identified the 29 most dominant features that help predict SD accurately. Again, we validate the optimality of these features using correlation and histogram plots. The trained model performs better on multiple performance metrics. Also, it predicts a more accurate SD than the top-ranking method of the BraTS-2020 competition. Further, we also explore the interpretability of the model to understand its decision globally and locally using post-hoc methods, i.e., SHAP and PDP. Observing these plots, we found that first-order statistical features, Age, location-based and texture features play a crucial role in prediction. Also, these interpretability methods can provide valuable insights into the model that can give human-understandable inferences. The inferences obtained for six dominant features using these interpretability methods were in line with medical facts. We also find that WIR, Age, and location-based features influence the most in predicting survival days. We further verify this conclusion using the KM estimation method on the metadata available with the BraTS dataset. Thus, the model is robust in predicting brain tumor patients' survivability. In addition, the interpretability methods can help us to understand model behavior at multiple levels. This will ultimately help to develop trust between medical experts and ML models and incorporate it into clinical practices.

50 51

1

2 3

4

5

6 7

8 9

10

11

12

13 14

15

16 17

18 19

20 21 22

23 24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

2 3

4

5

6

7

8

9

10

11 12

13 14

15

16

17

18 19

20

21

22

23

24 25

26

27

28

29

30

31 32

33

34

35

36 37

38

39

40

41

42

43

# **Declarations**

М. Rov acknowledges the seed No. grant ORSP/R&D/PDPU/2019/MR/RO051 of PDEU (for the computing facility), the core research grant No. CRG/2020/000869 of the Science and Engineering Research Board (SERB), Govt. of India and the project grant no  $GUJCOST/STI/2021 - 22/3873 \circ fGUJCOST$ , Govt. of Gujarat, India. M. S. Raval acknowledges grant No. GUJCOST/STI/2021-22/3858 of Gujarat Council of Science and Technology (GUJCOST), Govt. of Gujarat, India for computing facility.

## References

- [1] Hanif F, Muzaffar K, Perveen K, Malhi SM, Simjee SU. Glioblastoma multiforme: a review of its epidemiology and pathogenesis through clinical presentation and treatment. Asian Pacific journal of cancer prevention: APJCP. 2017;18(1):3.
- [2] Ostrom QT, Cioffi G, Waite K, Kruchko C, Barnholtz-Sloan JS. CBTRUS Statistical Report: Primary Brain and Other Central Nervous System Tumors Diagnosed in the United States in 2014–2018. Neuro-Oncology. 2021 Oct;23(Supplement\_3):iii1-iii105. https://doi.org/10.1093/neuonc/ noab200.
  - [3] Rindi G, Klimstra DS, Abedi-Ardekani B, Asa SL, Bosman FT, Brambilla E, et al. A common classification framework for neuroendocrine neoplasms: an International Agency for Research on Cancer (IARC) and World Health Organization (WHO) expert consensus proposal. Modern Pathology. 2018 Aug;31(12):1770-1786. https://doi.org/10.1038/ s41379-018-0110-y.
  - [4] Fernández-Llaneza D, Gondová A, Vince H, Patra A, Zurek M, Konings P, et al. Towards fully automated segmentation of rat cardiac MRI by leveraging deep learning frameworks. Scientific Reports. 2022 Jun;12(1). https://doi.org/10.1038/s41598-022-12378-z.
- [5] Ronneberger O, Fischer P, Brox T. U-net: Convolutional networks for biomedical image segmentation. In: International Conference on Medical image computing and computer-assisted intervention. Springer; 2015. p. 234 - 241.
- [6] McKinlev R. Meier R. Wiest R. Ensembles of densely-connected CNNs with label-uncertainty for brain tumor segmentation. In: International MICCAI Brainlesion Workshop. Springer; 2018. p. 456–465.

- 50 51
- 52

#### AUTHOR SUBMITTED MANUSCRIPT - MLST-100882.R1

- [7] McKinley R, Rebsamen M, Daetwyler K, Meier R, Radojewski P, Wiest R. Uncertainty-driven refinement of tumor-core segmentation using 3Dto-2D networks with label uncertainty. arXiv preprint arXiv:201206436. 2020;.
  - [8] Bakas S, Akbari H, Sotiras A, Bilello M, Rozycki M, Kirby JS, et al. Advancing the cancer genome atlas glioma MRI collections with expert segmentation labels and radiomic features. Scientific data. 2017;4(1):1–13.
- [9] Bakas S, Reyes M, Jakab A, Bauer S, Rempfler M, Crimi A, et al. Identifying the best machine learning algorithms for brain tumor segmentation, progression assessment, and overall survival prediction in the BRATS challenge. arXiv preprint arXiv:181102629. 2018;.
- [10] Menze BH, Jakab A, Bauer S, Kalpathy-Cramer J, Farahani K, Kirby J, et al. The multimodal brain tumor image segmentation benchmark (BRATS). IEEE transactions on medical imaging. 2014;34(10):1993–2024.
- [11] Rajput S, Raval MS. A Review on End-To-End Methods for Brain Tumor Segmentation and Overall Survival Prediction. arXiv preprint arXiv:200601632. 2020;.
- [12] Jia H, Cai W, Huang H, Xia Y. H2NF-Net for Brain Tumor Segmentation using Multimodal MR Imaging: 2nd Place Solution to BraTS Challenge 2020 Segmentation Task. arXiv preprint arXiv:201215318. 2020;.
- [13] Wang Y, Zhang Y, Hou F, Liu Y, Tian J, Zhong C, et al. Modality-Pairing Learning for Brain Tumor Segmentation. arXiv preprint arXiv:201009277. 2020;.
- [14] McKinley R, Rebsamen M, Meier R, Wiest R. Triplanar ensemble of 3d-to-2d cnns with label-uncertainty for brain tumor segmentation. In: International MICCAI Brainlesion Workshop. Springer; 2019. p. 379–387.
- [15] Asenjo JM, Solís AML. MRI Brain Tumor Segmentation Using a 2D-3D U-Net Ensemble. In: Brainlesion: Glioma, Multiple Sclerosis, Stroke and Traumatic Brain Injuries. Springer International Publishing; 2021. p. 354–366.
- [16] Myronenko A. 3D MRI brain tumor segmentation using autoencoder regularization. In: International MICCAI Brainlesion Workshop. Springer; 2018. p. 311-320.
- [17] Crimi A, Bakas S. Brainlesion: Glioma, Multiple Sclerosis, Stroke and Traumatic Brain Injuries: 5th International Workshop, BrainLes 2019, Held in Conjunction with MICCAI 2019, Shenzhen, China, October 17, 2019, Revised Selected Papers, Part I. vol. 11992. Springer Nature; 2020.
- 48 49

1

2

3

4

5 6

7

8

9 10

11

12

13

14 15

16

17

18 19

20

21

22 23

24

25

26

27

28

29 30

31

32

33 34

35

36

37

38 39

40

41

42 43

44

45

46

- 50 51
- 52

1		20	1
1 2 3	[18]	Isensee F, Jaeger PF, Full PM, Vollmuth P, Maier-Hein KH. nnU-Net for Brain Tumor Segmentation. arXiv preprint arXiv:201100848. 2020;.	C
4 5 6 7	[19]	Braunstein V.: Nvidia data scientists take top spots in MICCAI 2021 Brain Tumor Segmentation Challenge. Available from: https://developer. nvidia.com/blog/.	6
8 9 10 11	[20]	Agravat RR, Raval MS. A Survey and Analysis on Automated Glioma Brain Tumor Segmentation and Overall Patient Survival Prediction. Archives of Computational Methods in Engineering. 2021;p. 1–36.	3
12 13 14 15 16	[21]	Karami G, Giuseppe Orlando M, Delli Pizzi A, Caulo M, Del Gratta C. Predicting Overall Survival Time in Glioblastoma Patients Using Gra- dient Boosting Machines Algorithm and Recursive Feature Elimination Technique. Cancers. 2021;13(19):4976.	
17 18 19 20	[22]	Baid U, Rane SU, Talbar S, Gupta S, Thakur MH, Moiyadi A, et al. Overall survival prediction in glioblastoma with radiomic features using machine learning. Frontiers in computational neuroscience. 2020;p. 61.	
21 22 23	[23]	Hermida LC, Gertz EM, Ruppin E. Predicting cancer prognosis and drug response from the tumor microbiome. Nature Communications. 2022 May;13(1). https://doi.org/10.1038/s41467-022-30512-3.	
24 25 26	[24]	Walid MS. Prognostic factors for long-term survival after glioblastoma. The Permanente Journal. 2008;12(4):45.	
27 28 29 30	[25]	Feng X, Dou Q, Tustison N, Meyer C. Brain tumor segmentation with uncertainty estimation and overall survival prediction. In: International MICCAI Brainlesion Workshop. Springer; 2019. p. 304–314.	
31 32 33 34	[26]	Vale-Silva LA, Rohr K. Long-term cancer survival prediction using multimodal deep learning. Scientific Reports. 2021 Jun;11(1). https://doi.org/10.1038/s41598-021-92799-4.	
35 36 37 38	[27]	Bommineni VL. PieceNet: A Redundant UNet Ensemble. In: Brainlesion: Glioma, Multiple Sclerosis, Stroke and Traumatic Brain Injuries. Springer International Publishing; 2021. p. 331–341.	
39 40 41 42 43	[28]	Isensee F, Kickingereder P, Wick W, Bendszus M, Maier-Hein KH. Brain tumor segmentation and radiomics survival prediction: Contribu- tion to the brats 2017 challenge. In: International MICCAI Brainlesion Workshop. Springer; 2017. p. 287–297.	
44 45 46 47	[29]	Spyridon (Spyros) BCS.: Validation Survival Leaderboard 2020. Accessed: 2021-06-12. https://www.cbica.upenn.edu/BraTS20//lboardValidationSurvival.html.	
48 49 50		Y	

#### AUTHOR SUBMITTED MANUSCRIPT - MLST-100882.R1

1

2

3

4 5

6

7 8

9

10

11 12

13

14

15

16 17

18

19 20

21

22

23

24

25 26

27

28

29 30

31

32 33

34

35 36

37

38

39

40 41

42

43

44

24

[30] Friedman JH. Greedy function approximation: A gradient boosting machine. The Annals of Statistics. 2001;29(5):1189 – 1232. https://doi.org/ 10.1214/aos/1013203451.[31] Friedman JH, Popescu BE. Predictive learning via rule ensembles. The Annals of Applied Statistics. 2008;2(3):916–954. [32] Lundberg SM, Lee SI. A unified approach to interpreting model predictions. In: Proceedings of the 31st international conference on neural information processing systems; 2017. p. 4768–4777. [33] Lundberg SM, Erion G, Chen H, DeGrave A, Prutkin JM, Nair B, et al. From local explanations to global understanding with explainable AI for trees. Nature Machine Intelligence. 2020 Jan;2(1):56–67. https://doi.org/ 10.1038/s42256-019-0138-9. [34] Goel MK, Khanna P, Kishore J. Understanding survival analysis: Kaplan-Meier estimate. International journal of Ayurveda research. 2010;1(4):274. [35] Van Griethuysen JJ, Fedorov A, Parmar C, Hosny A, Aucoin N, Narayan V, et al. Computational radiomics system to decode the radiographic phenotype. Cancer research. 2017;77(21):e104-e107. [36] Singh SP, Urooj S. Wavelets: biomedical applications. International Journal of Biomedical Engineering and Technology. 2015;19(1):1–25. [37] Kong H, Akakin HC, Sarma SE. A generalized Laplacian of Gaussian filter for blob detection and its applications. IEEE transactions on cybernetics. 2013;43(6):1719-1733.[38] Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et al. Scikit-learn: Machine Learning in Python. arXiv:1201090. 2012;. [39] MIT MK, Lopuhin K.: permutation\_importance. Available from: https: //eli5.readthedocs.io/en/latest/blackbox/permutation\_importance.html. [40] Fernáandez-Delgado M, Cernadas E, Barro S, Amorim D. Do we need hundreds of classifiers to solve real world classification problems? The journal of machine learning research. 2014;15(1):3133–3181. [41] Puybareau E, Tochon G, Chazalon J, Fabrizio J. Segmentation of Gliomas and Prediction of Patient Overall Survival: A Simple and Fast Procedure. In: Crimi A, Bakas S, Kuijf H, Keyvan F, Reyes M, van Walsum T, editors. Brainlesion: Glioma, Multiple Sclerosis, Stroke and Traumatic Brain Injuries. Cham: Springer International Publishing; 2019. p. 199– 209.

[42]	Agravat RR, Raval MS. Brain tumor segmentation and survival predic- tion. In: International MICCAI Brainlesion Workshop. Springer; 2019. p. 338–348.
[43]	Ishwaran H, Kogalur UB, Blackstone EH, Lauer MS, et al. Random survival forests. Annals of Applied Statistics. 2008;2(3):841–860.
[44]	Rajput S, Agravat R, Roy M, Raval MS. Glioblastoma multiforme patient survival prediction. arXiv preprint arXiv:210110589. 2021;.
[45]	Rozemberczki B, Watson L, Bayer P, Yang HT, Kiss O, Nilsson S, et al. The Shapley Value in Machine Learning. arXiv preprint arXiv:220205594. 2022;.
[46]	Molnar C. Interpretable Machine Learning A Guide for Making Black Box Models Explainable. Leanpub; 2021.
[47]	Pan X, Zhang T, Yang Q, Yang D, Rwigema JC, Qi XS. Survival pre- diction for oral tongue cancer patients via probabilistic genetic algorithm optimized neural network models. The British Journal of Radiology. 2020;93(1112):20190825.
[48]	Molina G, Chawla A, Clancy TE, Wang J.: The correlation between the proportion of patients with pancreatic ductal adenocarcinoma who received neoadjuvant therapy and overall survival between 2004 and 2015. American Society of Clinical Oncology.
[49]	Minoru.: regression - What does the median absolute error metric say about the models? URL:https://stats.stackexchange.com/q/253892 (version: 2017-04-13) accessed: 2021-06-12. Cross Validated.
[50]	Ali MJ, Akram MT, Saleem H, Raza B, Shahid AR. Glioma Segmentation Using Ensemble of 2D/3D U-Nets and Survival Prediction Using Multiple Features Fusion. In: Brainlesion: Glioma, Multiple Sclerosis, Stroke and Traumatic Brain Injuries. Springer International Publishing; 2021. p. 189– 199.
[51]	Aboussaleh I, Riffi J, Mahraz AM, Tairi H. Brain tumor segmenta- tion based on deep learning's feature representation. Journal of Imaging. 2021;7(12):269.
[52]	Bae S, Choi YS, Ahn SS, Chang JH, Kang SG, Kim EH, et al. Ra- diomic MRI phenotyping of glioblastoma: improving survival prediction. Radiology. 2018;289(3):797–806.
[53]	Tessamma T, Ananda Resmi S. Texture Description of low grade and high grade Glioma using Statistical features in Brain MRIs. ACEEE; 2010

#### AUTHOR SUBMITTED MANUSCRIPT - MLST-100882.R1

- [54] ASCO ASoCO.: Brain Tumor: Statistics. Available from: https://www.cancer.net/cancer-types/brain-tumor/statistics.
  - [55] Mahmoudzadeh AP, Kashou NH. Interpolation-based super-resolution reconstruction: effects of slice thickness. Journal of Medical Imaging. 2014;1(3):034007.
  - [56] Fyllingen EH, Bø LE, Reinertsen I, Jakola AS, Sagberg LM, Berntsen EM, et al. Survival of glioblastoma in relation to tumor location: a statistical tumor atlas of a population-based cohort. Acta neurochirurgica. 2021;163(7):1895–1905.
- [57] Rizzo S, Botta F, Raimondi S, Origgi D, Fanciullo C, Morganti AG, et al. Radiomics: the facts and the challenges of image analysis. European radiology experimental. 2018;2(1):1–8.
- [58] Gupta M, Rajagopalan V, Rao BP. Glioma grade classification using wavelet transform-local binary pattern based statistical texture features and geometric measures extracted from MRI. Journal of Experimental & Theoretical Artificial Intelligence. 2019;31(1):57–76.
- [59] Deepa B, Sumithra M, Kumar RM, Suriya M. Weiner filter based hough transform and wavelet feature extraction with neural network for classifying brain tumor. In: 2021 6th International Conference on Inventive Computation Technologies (ICICT). IEEE; 2021. p. 637–641.
- [60] Steven AJ, Zhuo J, Melhem ER. Diffusion kurtosis imaging: an emerging technique for evaluating the microstructural environment of the brain. American journal of roentgenology. 2014;202(1):W26–W33.
- [61] Der G, Everitt BS. Statistical analysis of medical data using SAS. Chapman and Hall/CRC; 2005.
- [62] Sanghani P, Ang BT, King NKK, Ren H. Overall survival prediction in glioblastoma multiforme patients from volumetric, shape and texture features using machine learning. Surgical oncology. 2018;27(4):709–714.
- [63] Yang Z, Hu Z, Ji H, Lafata K, Floyd S, Yin FF, et al. A Neural Ordinary Differential Equation Model for Visualizing Deep Neural Network Behaviors in Multi-Parametric MRI based Glioma Segmentation. arXiv preprint arXiv:220300628. 2022;.
- [64] Li Y, Kim MM, Wahl DR, Lawrence TS, Parmar H, Cao Y. Survival Prediction Analysis in Glioblastoma With Diffusion Kurtosis Imaging. Frontiers in Oncology. 2021;11:690036.

- [65] Horská A, Barker PB. Imaging of brain tumors: MR spectroscopy and metabolic imaging. Neuroimaging Clinics. 2010;20(3):293–310.
- [66] Law M, Yang S, Wang H, Babb JS, Johnson G, Cha S, et al. Glioma grading: sensitivity, specificity, and predictive values of perfusion MR imaging and proton MR spectroscopic imaging compared with conventional MR imaging. American journal of neuroradiology. 2003;24(10):1989–1998.





Figure A2 PDP of Dominant features.

#### AUTHOR SUBMITTED MANUSCRIPT - MLST-100882.R1



Figure A2 PDP of dominant features.



Figure A2 PDPs of Dominants features: X-axis shows values of respective features, and Y-axis shows the average rate of change of feature effect on target feature. The vertical bars on X-axis show data distribution. This captures global trends of desirable features on the target variable by considering all the samples.











**Figure A3** SHAP value and its distribution for the dominant features. The X-axis shows the feature value of the respective feature, whereas Y-axis shows the SHAP-value of respective instances. The shaded region shows the distribution of instances. Each dot is an instance from the training dataset. This can help us to visualize and analyze where the majority of SHAP feature value lies, how individually the instances impact target features (range of impact instance wise), and its distribution. This testifies how these values play a role in defining the important feature. From all the SHAP plots, we can observe the magnitude of the SHAP value reduces with the order of importance of features (high to low).

**Note:** SHAP value calculates how much feature value changes the model's predicted value from the average.



Figure A4 Kaplan-Meier survival plot for survival probability.

The KM estimator measures the percentage of patients that have survived over a certain period after the treatment or surgery. It computes probabilities of the occurrence of events for a duration of time by dividing them into small intervals and re-estimates the probabilities to get the final estimate. The survival probability is computed as follows:

$$S_{t+1} = S_t \times ((N_{t+1} - D_{t+1})/N_{t+1})$$
(A1)

where, N denotes the number of people at risk and D denotes the number of people who died and t is the time interval. The KM survival curve is shown in Figure A4. It is a cumulative measure, and the survival remains the same until another individual encounters the risk. From this plot, we observed that the survival probability of older patients is low. The survivability reduces almost linearly after the age of 50 and almost exponentially after the age of 70 and it is very low beyond the age of 80. This KM analysis on the metadata supports the PDP analysis, which shows the exponential decay of survivability from the age near 70.

### A.2 Supplementary Tables:

Table A1: Dominant features obtained through RFE and their details.

#### Features with their descriptions

Age: age information given with the dataset.

*cent-at-y* : centroid of active tumor across the y axis. *pos-ed-wb-x*: centroid of enhancing tumor w.r.t brain centroid across x

*pos-ea-wo-x*: centroid of enhancing tumor w.r.t brain centroid across x axis.

original-shape-LeastAxisLength: It calculates smallest axis length of the ROI.

**Wavelet-LLH-firstorder-Maximum**: It measures maximum gray intensity within the ROI after applying the Wavelet LLH-band filter.

Wavelet-LLH-Gldm-Dependence Variance: It calculates variance in
the dependence matrix of the image after applying the Wavelet LLH-band
filter.

Wavelet-LLH-Glrlm-LongRunLowGrayLevel-Emphasis: It esti mates the length in the terms of successive pixels run lengths with lower
 gray-level intensity values, after applying Wavelet transform.

Wavelet-LHL-Glcm-Correlation: It assesses the correlation between
 gray-level values and their related voxels in the gray-level co-occurrence
 matrix, after applying the Wavelet LHL band filter.

Wavelet-LHH-Gldm-DependenceNonUniformityNormalized: It
 assesses the similarity between dependencies in an image, after applying
 the Wavelet LHH band filter.

Wavelet-LHH-Gldm-SmallDependenceHighGrayLevelEmphasis:

It assesses the combined distribution of small dependence with higher gray-level values after applying the Wavelet LHH band filter.

gray-level values after applying the Wavelet LHH band filter.
 Wavelet-LHH-Glszm-ZoneEntropy: It evaluates the randomization
 of zone sizes and gray level values in the distribution after applying the
 Wavelet LHH band filter.

Wavelet-HLL-Glem-Imc1: It assesses the correlation between the prob ability distribution of two pixels after applying the Wavelet HLL band
 filter.

Wavelet-HLH-firstorder-Kurtosis: It assesses the peakedness of the
 spread of pixel intensities in the given image after applying the Wavelet
 HLH band filter.

Wavelet-HLH-Gldm-DependenceEntropy: It assesses randomness in the dependencies of an image after applying the Wavelet HLH band filter. Wavelet-HLH-Gldm-SmallDependenceLowGrayLevel-Emphasis:

It assesses the combined distribution of small dependence with higher gray-level values after applying Wavelet HLH band filter.

45 46 47

41

42

43

44

1

2 3 4

5

6 7

8

9

10

11

12

13

14

15

16

- 48 49
- 50
- 51 52

1	
2 3	Wavelet-HHH-Glcm-MaximumProbabilitu: It finds the most fre-
4	quently occurring neighboring pair of intensity values from the grev-level
5	co-occurrence matrix after applying the Wavelet HHH band filter.
6	Wavelet-LLL-Glcm-Correlation: It assesses the association between
7	pairs and their corresponding voxel intensity value after applying the
8	Wavelet LLL band filter.
9	LoG-sigma-1-0-mm-3D-Glcm-Correlation: It assesses the associa-
10	tion between pairs and their respective voxel intensity value after applying
11	LoG filter with sigma value 1.
12	Wavelet-LLH-Ngtdm-Strength: It assesses strength in an image after
13	applying a Wavelet filter using the LLH band.
14	LoG-sigma-5-0-mm-3D-Glrlm-RunLengthNonUniformity-
15	<i>Normalized</i> : It assesses the homogeneity in the gray level run lengths in
16	the image after applying the LoG filter with sigma value 5.
17	LoG-sigma-3-0-mm-3D-Glrlm-RunVariance: It calculates variance
18	in the gray-level run-lengths in the image, after applying LoG filter with
19	sigma value 3.
20	LoG-sigma-2-0-mm-3D-Glcm-ClusterShade: It calculates unifor-
21	mity in the gray-level co-occurrence matrix after applying the LoG filter
22	with sigma value 2.
23	LoG-sigma-5-0-mm3D-firstorder-TotalEnergy: It assesses the local-
24	ized change of the image after applying the LoG filter with sigma value
25	5.
26	LoG-sigma-3-0-mm-3D-Glcm-MaximumProbability: It assesses the
27	occurrences of the most prevalent pairing of neighboring intensity values
28	in the grey-level co-occurrence matrix after applying the LoG filter with
29	sigma value 5.
30	LoG-sigma-2-0-mm-3D-firstorder-90Percentile: It assesses 90th
31	percentile intensity values of an image after applying LoG filter with sigma
32	value 2.
33	LoG-sigma-2-0-mm-3D-firstorder-Skewness: It calculates the asym-
34	metry of the distribution of intensity values that deviates from the mean
35	intensity value after applying the LoG filter with sigma value 2.
36	LoG-sigma-1-0-mm-3D-Glcm-MCC: It assesses the complexity of the
37	texture in the co-occurrence matrix of an image after applying the LoG
38	filter with sigma value 1.
39	Wavelet-LLL-Glszm-SmallAreaEmphasis: It assesses the number of
40	connected voxels with the same gray-level intensity value or the spread of
41	smaller size zones after applying Wavelet LLL band filter.
42	Wavelet-HLL-Glcm-MCC: It assesses the complexity of the texture
43	in the co-occurrence matrix of an image after applying the Wavelet HLL
44	band filter.
45	
46	V ′
47	
48	V
49	₹

~

**Table A2**: Dominant feature set through PI and their weights.The threshold value of the weights is 100.

Weight	Features
309.94	Age : age information given with the dataset.
)761.33	<b>LoG-sigma-1-0-mm-3D-glcm-Correlation</b> : It assesses the association between pairs and its respective voxel intensity value after explained the LoC filter with size solution 1
)722.95	after applying the LoG filter with sigma value 1. <b>Wavelet-HHH-Gldm-Dependence Variance</b> : It calculates variance in the dependence matrix of the image after applying the Wavelet HHH band filter.
0678.10	LoG-sigma-4-0-mm-3D-Glcm-JointEntropy: It calculates the randomness in neighborhood intensity values.
0669.58	<b>LoG-sigma-2-0-mm-3D-firstorder-Kurtosis</b> : It assesses the peakiness of the intensity distribution of a given image after applying the LoG filter with sigma value 2
0558.74	LoG-sigma-2-0-mm-3D-Glrlm- HighGrayLevelRunEmphasis: It assesses the spread of the image's upper grav-level values in the image
0555.77	Wavelet-HLH-Gldm-SmallDependence- LowGrayLevelEmphasis: It assesses the combined spread of small-dependence with lower gray-level values after applying Wavelet HLH band filter.
0509.37	LoG-sigma-3-0-mm-3D-Gldm-LowGrayLevelEmphasis: It calculates the spread of low gray-level values in the image.
0476.60	<i>cent-ncr-x</i> : centroid of necrosis across x-axis.
0464.70	<b>Wavelet-LLL-firstorder-InterquartileRange</b> : It assesses the difference between the 75th and 25th percentile of the image array after applying the Wavelet LLL band filter.
0444.84	LoG-sigma-4-0-mm-3D-Glcm-ClusterShade: It calculates uniformity in the gray level co-occurrence matrix after applying LoG filter with sigma value 4.
0438.99	<b>Wavelet-LHH-firstorder-RootMeanSquared</b> : It assesses the root-mean-square of the intensity value of an image after applying the Wavelet LHH band filter.
0420.35	<b>LoG-sigma-4-0-mm-3D-Glcm-SumAverage</b> : It assesses the relationship between pair occurrences with lower intensity values and pair occurrences with higher intensity values, after applying LoG filter with sigma value 4.
0406.08	Wavelet-HHH-Glrlm-RunLengthNonUniformity: It as-

1 2		
2	0395.63	LoG-sigma-5-0-mm-3D-Glszm-SmallAreaEmphasis: It
5 ⊿		assesses the spread of small size-zones or the number of con-
4 5		nected voxels that have the same gray-level intensity value,
5		after applying LoG filter with sigma value 5.
7	0357.96	Wavelet-LLH-Ngtdm-Coarseness: It assesses the spatial
7 Q		rate of change in the intensity value after applying the Wavelet
0 0		LLH band filter.
10	0357.51	Wavelet-LLH-firstorder-InterquartileRange: It assesses
11		the difference between the 75th and 25th percentile of the image
12		array after applying the Wavelet LLH band filter.
13	0340.11	<i>cent-at-x</i> : centroid of active tumor across x-axis.
14	0314.18	LoG-sigma-4-0-mm-3D-Glszm-
15		LargeAreaLowGrayLevel-Emphasis.
16	0282.22	Wavelet-HHH-firstorder-Kurtosis: It assesses the peaked-
17		ness of the spread of the image's intensity values after applying
18	001500	the Wavelet HHH band filter.
19	0247.80	Wavelet-HHH-Glcm-DifferenceAverage: It assesses the re-
20		intensity values and these with different intensity values after
21		applying the Waydet filter
22	0247.36	<i>cent-wb-r</i> : centroid of whole-tumor brain across y-axis
23	023256	LoC-sigma 2-0-mm 2D-first order Energy: It assesses the
24	0252.50	magnitude of voyel values in an image
25	0229.91	LoG-sigma-1-0-mm-3D-firstorder-Variance: It measures
27	0000_	the distribution spread about the mean intensity value after
28		applying LoG filter with sigma value 1.
29	0226.71	Wavelet-LHH-firstorder-Kurtosis: It assesses the image's
30		peakedness in terms of intensity distribution, applying Wavelet
31		LHH band filter.
32	0217.80	LoG-sigma-2-0-mm-3D-Glszm-
33		LargeAreaHighGrayLevel-Emphasis: It assesses the
34		combined spread of larger size-zones with higher gray-level
35		values, after applying the LoG filter with sigma value 1.
36	0183.47	Wavelet-LLH-firstorder-Range: It assesses the distribution
3/	0191 10	of gray-level values of an image.
38	0151.10	downess in the dependencies of an image after applying Wavelet
<u>40</u>		LHH band filter
41	0118.90	Wavelet-LHL-Glcm-ClusterShade: It calculates uniformity
42	0110100	in the grav level co-occurrence matrix after applying the Wavelet
43		LHL band filter.
44		
45		
46		7
47	X	
48		
49		

1 
 Table A3
 Feature annotation of a correlation matrix.
 2 3 Index Features Name Features Type 4 No. 5 1 Aae Meta-Data 6 2 cent-at-x Image-based 7 3 cent-ncr-x Image-based 8 4 cent-wb-x Image-based 9 5LoG-sigma-1-0-mm-3D-FirstorderVariance Radiomics-based 6 LoG-siama-1-0-mm-3D-Glcm-Correlation Radiomics-based 10 7 LoG-sigma-2-0-mm-3D-FirstorderKurtosis Radiomics-based 11 Radiomics-based 8 LoG-sigma-2-0-mm-3D-Glrlm-12 HighGrayLevelRunEmphasis13 9 LoG-sigma-2-0-mm-3D-Glszm-Radiomics-based 14 LargeAreaHighGrayLevelEmphasis10 LoG-sigma-3-0-mm-3D-Firstorder-Energy Radiomics-based 15 11 LoG-sigma-3-0-mm-3D-Gldm-Radiomics-based 16 LowGrayLevelEmphasis 17 12LoG-sigma-4-0-mm-3D-Glcm-ClusterShade Radiomics-based 18 13 LoG-sigma-4-0-mm-3D-Glcm-JointEntropy Radiomics-based 19 14 LoG-sigma-4-0-mm-3D-Glcm-SumAverage Radiomics-based 20 LoG-sigma-4-0-mm-3D-Glszm-Radiomics-based 1521 LargeAreaLowGrayLevelEmphasis 16 LoG-sigma-5-0-mm-3D-Glszm-Radiomics-based 22 SmallAreaEmphasis 23 Wavelet-HHH-Firstorder-Kurtosis Radiomics-based 1724 18 Wavelet-HHH-Glcm-DifferenceAverage Radiomics-based 25 19 Wavelet-HHH-Gldm-Dependence Variance Radiomics-based 26 20Wavelet-HHH-Glrlm-RunLengthNonUniformity Radiomics-based 27 21Wavelet-HLH-Gldm-SmallDependence-Radiomics-based Low Gray Level Emphasis28 Wavelet-LHH-Firstorder-Kurtosis 22Radiomics-based 29 Wavelet-LHH-Firstorder-RootMeanSquared 23Radiomics-based 30 24 Wavelet-LHH-Gldm-DependenceEntropy Radiomics-based 31 25Wavelet-LHL-Glcm-ClusterShade Radiomics-based 32 26Wavelet-LLH-Firstorder-InterguartileRange Radiomics-based 33 Wavelet-LLH-Firstorder-Range 27Radiomics-based 34 28Wavelet-LLH-Ngtdm-Coarseness Radiomics-based 35 29 Wavelet-LLL-Firstorder-InterquartileRange Radiomics-based 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50

#### AUTHOR SUBMITTED MANUSCRIPT - MLST-100882.R1

**Table A4** Example of calculating SHAP value. Let us consider Feature set  $(F) = \{A, B, D\}$ , and values (contribution) of features are:  $v\{A\} = 8$ ,  $v\{B\} = 10$ ,  $v\{D\} = 9$ ,  $v\{A, B\} = 18$ ,  $v\{A, D\} = 20$ ,  $v\{B, D\} = 22$  and  $v\{A, B, D\} = 25$ .

Possible combinations of feature	Marginal Combination				
	Feature A	Feature B	Feature D		
$\{\mathbf{A}, \mathbf{B}, \mathbf{D}\}$	$v{A} -\phi = 8$	$v{A,B}-v{A} = 10$	$v{A,B,D} - v{A,B} = 7$		
$\{A, D, B\}$	$v{A} -\phi = 8$	$v{A,B,D} -v{A,D}=5$	v{A,D} -v{A}=12		
$\{D, B, A\}$	$v{A,B,D} - v{D,B} = 25-22 = 3$	$v{D,B} -v{B} =12$	v{D} - $\phi = 9$		
$\{B, A, D\}$	$v{A,B} - v{B} = 8$	$v{B} - \phi = 10$	$v{A,B,D} - v{A,B} = 7$		
$\{D, A, B\}$	$v{A,D} -v{D} = 11$	$v{A,B,D} - v{A,D} = 5$	$v{D} - \phi = 9$		
{B, D, A}	v{A,B,D} -v{B,D}=3	$v{B} - \phi = 10$	v{B,D} -v{B}=12		
SHAP value	$(8+8+8+10+11+3) \mid 6 = 6.833$	$(10+5+12+10+5+10) \mid 6 = 8.667$	(7+12+9+7+9+12)   6 = 9.334		