

# Beyond Known Reality: Exploiting Counterfactual Explanations for Medical Research

Toygar Tanyel<sup>1</sup>(✉), Serkan Ayvaz<sup>2</sup>(✉), and Bilgin Keserci<sup>3</sup>

<sup>1</sup> Biomedical Engineering Graduate Program, Istanbul Technical University

<sup>2</sup> Centre for Industrial Software, Maersk Mc-kinney Moeller Institute, University of Southern Denmark

<sup>3</sup> Department of Biomedical Engineering, Yildiz Technical University  
(✉) tanyel23@itu.edu.tr (✉) seay@mmmi.sdu.dk

**Abstract.** The field of explainability in artificial intelligence (AI) has witnessed a growing number of studies and increasing scholarly interest. However, the lack of human-friendly and individual interpretations in explaining the outcomes of machine learning algorithms has significantly hindered the acceptance of these methods by clinicians in their research and clinical practice. To address this issue, our study uses counterfactual explanations to explore the applicability of "what if?" scenarios in medical research. Our aim is to expand our understanding of magnetic resonance imaging (MRI) features used for diagnosing pediatric posterior fossa brain tumors beyond existing boundaries. In our case study, the proposed concept provides a novel way to examine alternative decision-making scenarios that offer personalized and context-specific insights, enabling the validation of predictions and clarification of variations under diverse circumstances. Additionally, we explore the potential use of counterfactuals for data augmentation and evaluate their feasibility as an alternative approach in our medical research case. The results demonstrate the promising potential of using counterfactual explanations to improve AI-driven methods in clinical research.

**Keywords:** counterfactual explanations · machine learning · pediatric brain tumors · explainable AI · magnetic resonance imaging

## 1 Introduction

As we incorporate automated decision-making systems into the real world, explainability and accountability questions become increasingly important [33]. In some fields, such as medicine and healthcare, ignoring or failing to address such a challenge can seriously limit the adoption of computer-based systems that rely on machine learning (ML) and computational intelligence methods for data analysis in real-world applications [3, 11, 54]. Previous research in eXplainable Artificial Intelligence (XAI) has primarily focused on developing techniques to interpret decisions made by black box ML models. For instance, widely used

approaches such as local interpretable model-agnostic explanations (LIME) [42] and shapley additive explanations (SHAP) [29] offer attribution-based explanations to interpret ML models. These methods can assist computer scientists and ML experts in understanding the reasoning behind the predictions made by AI models.

On the other hand, counterfactual explanations [32, 57] are a form of model-agnostic interpretation technique that identifies the minimal changes needed in input features to yield a different output, aligned with a specific desired outcome. This approach may be more interesting to end users, including clinicians and patients, rather than focusing solely on how models arrive at their predictions. Because counterfactual explanations can help to better understand the practical consequences of the ML model’s predictions and take practical actions at the individual level. Patients’ primary concern is not only to learn about their diseases, but also to seek guidance on how to regain their health. For example, counterfactual explanations could potentially be used by clinicians to help inform patients about what small changes they need to make to become healthy again. Understanding the doctor’s or ML model’s decision-making process is less important to patients.

Leveraging counterfactual explanations holds promise in enhancing the data synthesize and for answering causal questions [39, 45], and interpretability of AI models by offering deeper insights into their decision-making processes [36, 53, 58]. Our proposed approach aims to uncover the underlying reasons for the observed relationships between MRI features, going beyond just generating actionable outcomes for individual patients. Through counterfactual explanations, previously unseen decisions within the decision space can be brought to light. Numerous questions can be explored, such as how to determine the modifications required to transform a patient’s diagnosis from one tumor subtype to another. Initially, posing such a question may seem nonsensical and illogical since an individual’s actual tumor type cannot be magically altered. However, considering the challenge of distinguishing these two tumor types in clinical settings, asking such a question can effectively demonstrate which features are more informative in differentiating tumor types. Counterfactual explanations enable us to identify the characteristics that distinguish two patient types with the smallest changes in features. Consequently, a deeper understanding of the interactions between MRI features and tumors can be gained; unveiling previously undisclosed outcomes that may be concealed in existing ML studies.

Furthermore, we have identified a potential contribution to clinical practice whereby a new patient with only MRI data available can have their tumor type estimated using a counterfactual approach, prior to receiving histopathological results. Since there is no prior label available for the patient, they are given an "unknown" label and the counterfactual approach is used for each tumor type, allowing estimation of the tumor type with the lowest distance and smallest change in features. While this approach shares similarities with ML, the crucial distinction lies in retaining information about the reasoning behind the estimated

tumor type and its corresponding feature changes. This, in turn, can enhance our understanding and the use of AI models in clinical practice.

Last but not least, in situations where the acquisition of data is limited or not possible, various data augmentation methods have been developed to enhance the performance of ML and related applications [7, 61, 69]. However, these methods also give rise to additional issues while fulfilling their intended purpose, such as introducing biased shifts in data distribution. To address this issue, we employed counterfactuals generated from different spaces in order to balance the data by maximizing its diversity, and subsequently reported the results for different scenarios.

To summarize, our main contributions include:

- introducing a new perspective on the application of counterfactual explanations in the pediatric posterior fossa brain tumor literature,
- demonstrating how the counterfactual approach, which enables us to provide patient-specific local explanations, can inform us in differentiating tumor types,
- providing a systematic comparative analysis of machine explanations, with the aim of evaluating the feasibility of the outcomes.

## 2 Material & Methods

This study integrates a systematic approach combining clinical data acquisition, MRI feature analysis, machine learning, and counterfactual explanations to investigate tumor classification in pediatric posterior fossa brain tumors.

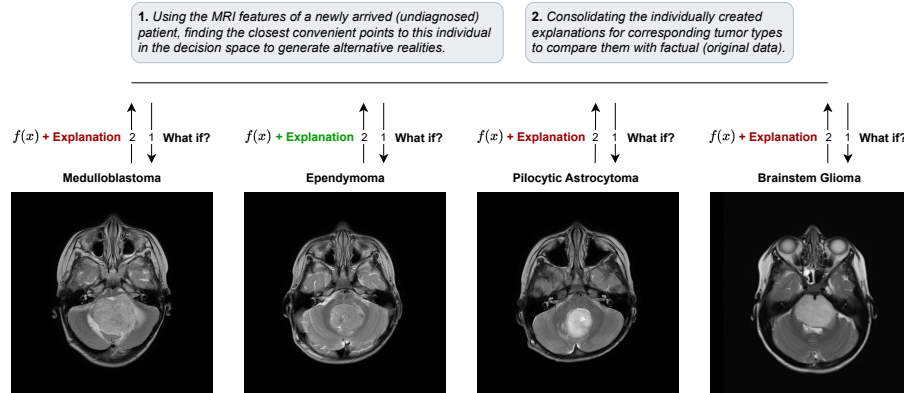


Fig. 1: The figure illustrates a hypothetical clinical scenario showcasing the practical application of counterfactuals. The data consists of features extracted from provided multi-parametric brain MRIs, rather than the raw images themselves.

## 2.1 Ethics Statement and Patient Characteristics

This prospective study (Ref: 632 QĐ-NĐ2 dated 12 May 2019) was carried out in both Radiology and Neurosurgery departments, and was approved by the Institutional Review Board in accordance with the 1964 Helsinki declaration. Written informed consent was obtained from authorized guardians of patients prior to the MRI procedure. Our study comprised a cohort of 112 pediatric patients diagnosed with posterior fossa tumors, including 42 with MB, 25 with PA, 34 with BG, and 11 with EP. All BG patients were confirmed based on full agreement between neuroradiologists and neurosurgeons, whereas the remaining MB, PA, and EP patients underwent either surgery or biopsy for histopathological confirmation.

## 2.2 Data Acquisition and Assessment of MRI Features

For all patients, MRI exams including T1W, T2W, FLAIR, DWI (b values: 0 and 1000) with ADC, and contrast-enhanced T1W (CE-T1) sequences with macro-cyclic gadolinium-based contrast enhancement (0.1 ml/kg Gadovist, Bayer, Germany, or 0.2 ml/kg Dotarem, Guerbet, France) were collected in the supine position using a 1.5 Tesla MRI scanner (Multiva, Philips, Best, the Netherlands).

The Medical Imaging Interaction Toolkit (German Cancer Research Center, Division of Medical Image Computing, Heidelberg, Germany) was utilized for measuring the region of interest (ROI) of posterior fossa tumors and normal-appearing parenchyma and subsequently assessed the following MRI features: signal intensities (SIs) of T2, T1, FLAIR, T1CE, DWI, and ADC. Ratios between the posterior fossa tumor and parenchyma were calculated by dividing the SI of the tumor and the SI of the normal-appearing parenchyma based on T2, T1, FLAIR, T1CE, DWI, and ADC. Additionally, ADC values were quantified for both the posterior fossa tumor and parenchyma on the ADC map using the MR Diffusion tool available in Philips Intellispace Portal, version 11 (Philips, Best, The Netherlands). It is worth noting that, prior to analysis, bias field correction was applied to every image to correct for nonuniform grayscale intensities in the MRI caused by field inhomogeneities.

## 2.3 Standardization

Prior to conducting ML trainings, the dataset was subjected to a standardization process, using Python programming (version 3.9.13) with the Scikit-Learn library (version 1.0.2) module. This technique involved transforming the data to have a mean of zero and a standard deviation of one. To standardize all numerical attributes, the Scikit-Learn StandardScaler function was employed, which subtracted the mean and scaled the values to unit variance, ensuring the data was in a standardized format. To determine the standard score of a sample  $x_i$ , the following formula is used:

$$z = \frac{x_i - \mu}{\sigma}, \quad (1)$$



where,  $\mu$  represents the mean of the training samples, and  $\sigma$  represents their standard deviation.

## 2.4 Distance Calculation

Utilizing counterfactuals as classifiers, the notable discrepancy in MRI feature values, as shown in the example in Table 1, complicates distance calculations. The underlying reason for the distance calculation issue is that some MRI features in fact include ratios that depend on other variables. We addressed this by omitting unchanged values ('-'), rescaling the rest to a consistent scale, and then reintroducing them. The distances were then computed using the Euclidean metric on the counterfactual values of the current factual. Minimizing this distance aids in determining the tumor type by corresponding to the least dissimilarity (Table 2). The Euclidean distance is given by:

$$\text{Distance} = \sqrt{\sum_{i=1}^n (x_i - y_i)^2} \quad (2)$$

Here,  $x_i$  and  $y_i$  are values from the current and baseline rows, respectively. The formula sums the squared differences for each feature, then takes the square root. In this equation,  $n$  signifies the feature count in the dataset.

## 2.5 Statistical Analysis

Using the  $t$ -test from the `scipy` library (version 1.10.1), we considered a two-tailed  $p$ -value of  $< 0.05$  as statistically significant. Our analysis comprised:

1. Assessing the statistical significance of changes in counterfactuals when transitioning the tumor type from  $\mathcal{X}$  to  $\mathcal{Y}$  (dependent  $t$ -test).
2. Comparing the similarity between counterfactuals transitioning from  $\mathcal{X}$  to  $\mathcal{Y}$  and the original patients with tumor type  $\mathcal{Y}$  (Welch's  $t$ -test).

For each transition from tumor type  $\mathcal{X}$  to  $\mathcal{Y}$ , we generated five counterfactuals. In our  $t$ -test evaluations, these counterfactuals were examined in distinct manners. Due to generating a dataset larger than our original sample size, we could not maintain equal dimensions for the dependent analysis. To address this:

- For each counterfactual, we used the data of the corresponding factual patient as a baseline for testing.
- We subsequently tested the significance of the top five most changed feature variables.

For the independent analysis, all counterfactuals were evaluated using the real data of all patients, focusing on the three most altered features. Each feature was evaluated separately.

In summary, the primary distinction between our tests was their focus: one on patients, and the other on features.

## 2.6 Distribution Plotting

To generate individual kernel density estimation (KDE) plots for each feature, we utilized the `kdeplot` function from the Seaborn package (version 0.11.2). By specifying a hue parameter (e.g., Tumor Type), we were able to incorporate a meaningful association using this method. Consequently, we transformed the default marginal plot into a layered KDE plot. This approach tackles the challenge of reconstructing the density function  $f$  using an independent and identically distributed (iid) sample  $x_1, x_2, \dots, x_n$  from the respective probability distribution.

## 2.7 Machine Learning

To decrease overfitting and convergence issue of counterfactuals, especially for EP, we took less patients to implement the task: 25 patients from MB, PA and BG and 11 patients from EP. For testing, to ensure the reliability of our ML models, particularly with a small dataset, we conducted five runs using stratified random sampling based on tumor type with 55% train and 45% test patients.

Using nine ML models, including support vector machine (SVM), adaboost (ADA), logistic regression (LR), random forest classifier (RF), decision tree classifier (DT), gradient boosting classifier (GB), catboost classifier (CB), extreme gradient boosting classifier (XGB) and voting classifier (VOTING), we evaluated the models on the raw data with the outcomes prior to our counterfactual interpretations. CB and XGB were obtained from CatBoost version 1.1.1 and XGBoost version 1.5.1 libraries, respectively, while the other models were obtained from the Scikit-Learn library.

We assessed the performance of the models using precision, recall, and F1 score, which were calculated based on the counts of true positives (TP), true negatives (TN), false positives (FP), and false negatives (FN). In order to ensure an accurate interpretation of the ML results, we opted not to balance the labels. Instead, we employed macro precision, macro recall, and macro F1 score metrics, which take into account the contributions of all labels equally. This approach enabled us to observe the genuine impact of the varying label frequencies, EP in this case.

## 3 Counterfactual Explanations

Given the challenges associated with local approximations, it is worthwhile to explore prior research in the "explanation sciences" to identify potential alternative strategies for generating reliable and practical post-hoc interpretations that benefit the stakeholders affected by algorithmic decisions [33, 44]. To create explanations that are understandable and useful for both experts and non-experts, it is logical to investigate theoretical and empirical studies that shed light on how humans provide and receive explanations [32]. Over the past few

decades, the fields of philosophy of science and epistemology have shown increasing interest in theories related to counterfactual causality and contrastive explanations [24, 27, 44, 63, 64].

In philosophy, counterfactuals serve not only to assess the relationship between a mental state and reality, but also to determine whether a mental state can be considered as knowledge. The problem of identifying knowledge with justified true belief is complicated by various counterexamples, such as Gettier cases (1963) [17]. However, some scholars proposed additional conditions to address these counterexamples. This literature highlighted two significant counterfactual conditions:

**Sensitivity:** If  $\rho$  were false,  $\mathcal{S}$  would not believe that  $\rho$ .

**Safety:** If  $\mathcal{S}$  were to believe that  $\rho$ ,  $\rho$  would not be false.

Both of these conditions express the notion that  $\mathcal{S}$ 's beliefs must be formed in a manner that is sensitive to the truthfulness of  $\rho$ . The counterfactual semantics has influenced from this idea in various ways, including the establishment of their non-equivalence, clarification, and resolution of potential counterexamples [50].

This concept has sparked a fresh wave of counterfactual analyses that employ new methodologies. Hitchcock [21, 22] and Woodward [62], for instance, constructed counterfactual analyses of causation using Bayesian networks (also known as "causal models") and structural equations. The basic idea of the analysis can be summarized as follows: " $\mathcal{X}$  can be considered a cause of  $\mathcal{Y}$  only if there exists a path from  $\mathcal{X}$  to  $\mathcal{Y}$ , and changing the value of  $\mathcal{X}$  alone results in a change in the value of  $\mathcal{Y}$ ".

Ginsberg (1986) [18] initiated his discussion by outlining the potential significance of counterfactuals in artificial intelligence and summarizing the philosophical insights that have been drawn regarding them. Following this, Ginsberg provided a structured explanation of counterfactual implication and analyzed the challenges involved in executing it. Over time, numerous developments in the fields of artificial intelligence and cognitive science, including the Bayesian epistemology approach, have gone beyond what was previously envisioned by Ginsberg regarding the potential application of artificial intelligence and counterfactuals [10, 19, 32, 48, 49, 57]. Furthermore, Verma et al. [55] conducted a comprehensive review of the counterfactual literature, analyzing its utilization in over 350 research papers.

In recent times, there has been a growing interest in the concept of counterfactual explanations, which aim to provide alternative perturbations capable of changing the predictions made by a model. In simple terms, when given an input feature  $x$  and the corresponding output produced by an ML model  $f$ , a counterfactual explanation involves modifying the input to generate a different output  $y$  using the same algorithm. To further explain this concept, Wachter et al. [57] introduce the following formulation in their proposal:

$$c = \arg \min_c \ell(f(c), y) + |x - c| \quad (3)$$

The initial component  $\ell$  of the formulation encourages the counterfactual  $c$  to deviate from the original prediction, aiming for a different outcome. Meanwhile,

the second component ensures that the counterfactual remains in proximity to the original instance, thereby emphasizing the importance of maintaining similarity between the two.

While challenges like the inability to find optimal counterfactual explanations underscore the need for DiCE updates, there are alternative solutions. Dutta et al. [15] and Maragno et al. [30] propose alternative counterfactual algorithms to potentially overcome these challenges. Furthermore, Guidotti’s review [20] presents an extensive list of counterfactual algorithms.

Although the subject of explainable artificial intelligence has been investigated by many researchers in the field of medicine, there are not many studies exploiting counterfactual explanations [5, 56]. Among explainable AI methods, it appears that SHAP, Grad-CAM, and Lime have been extensively studied [4, 52]. Sarp et al. focused on a LIME-based heat map application for interpretation of COVID-19 cases from chest X-ray images [46]. In another study [28], authors applied gradient-weighted class activation mapping (Grad-CAM) to localize decision-making regions. Knapič et al. investigated the use of SHapley Additive exPlanations (SHAP) in medical images and compared its performance with LIME and the Contextual Importance and Utility (CIU) method [25]. In [68], the authors explored the use of different XAI methods including Vanilla gradient, guided backpropagation, integrated gradients, guided integrated gradients, SmoothGrad, Grad-CAM, and guided Grad-CAM to explain brain tumor segmentation. Their work focused on methods for creating visualization maps to better interpret deep learning models. Differently in this study, we propose a counterfactual explanation-based approach that explores alternative decision-making scenarios to provide personalized and context-specific insights based on MRI data.

### 3.1 Generating Counterfactual Explanations

The clear interpretability of counterfactuals helps an individual make decisions about his or her future [1, 2, 12, 43, 58, 60, 65–67]. In situations where a *negative* response is received, understanding how to improve results without resorting to major and unrealistic data alterations becomes important.

We argue that the use of counterfactual explanations effectively leverages the factual insights derived from MRI features. These features act as distinct markers, facilitating the differentiation between various tumors. Such a methodology shines particularly when traditional diagnostics find it challenging to distinguish between tumor types.

The *data manifold* concept, illustrated in Fig. 2, emphasizes the importance of proximity in counterfactual explanations. For counterfactuals to be credible, their features should resemble those of prior classifier observations and be realistic. Counterfactuals with features that diverge from training data or disrupt feature associations are impractical and outside established norms [6]. To ensure that counterfactuals are realistic and align with the training data, we employ constraint-based approaches in our algorithms. For instance, changing parameters such as "age" and "gender" would be highly unreasonable. Therefore, in

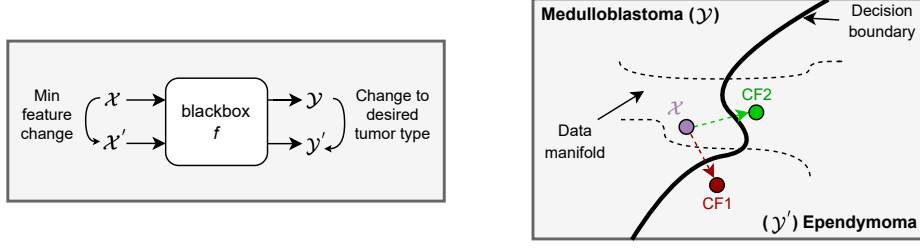


Fig. 2: Generating counterfactual explanations on tumor types. The left illustrates feature manipulation using an ML model with a counterfactual approach. The right delves deeper into the process and counterfactual explanation concept, exemplified by two tumor types.

most scenarios, these parameters are specified in the model to prevent counterfactuals from deviating from reality. Similarly, in our case, while parenchyma features were included during training, they remain invariant during counterfactual generation to preserve tissue characteristic references.

The DiCE library [35] offers a framework for counterfactual generation, viewing it as an optimization task akin to adversarial example discovery. Crucially, the modifications must be diverse, feasible, and implementable. The optimization formula used is:

$$C(x) = \arg \min_{c_1, \dots, c_k} \frac{1}{k} \sum_{i=1}^k \ell(f(c_i), y) + \frac{\lambda_1}{k} \sum_{i=1}^k \text{dist}(c_i, x) - \lambda_2 \text{dpp\_diversity}(c_1, \dots, c_k), \quad (4)$$

with hinge loss as:

$$\ell = \max(0, 1 - z * \text{logit}(f(c))), \quad (5)$$

where  $z$  takes values based on  $y$  and  $\text{logit}(f(c))$  represents unscaled ML output. Diversity is expressed as:

$$\text{dpp\_diversity} = \det(K), \quad (6)$$

where  $K_{i,j} = \frac{1}{1 + \text{dist}(c_i, c_j)}$ , and the distance between counterfactuals is measured.

## 4 Results

### 4.1 What if the counterfactual explanations graciously provide us with additional insights into classification?

Using DiCE’s multi-class training capability, we have established a framework that simultaneously trains for four distinct tumor types. This framework employs counterfactual explanations, serving as classifiers, to determine tumor types based on numeric MRI data. Figure 3 presents an overview of our approach

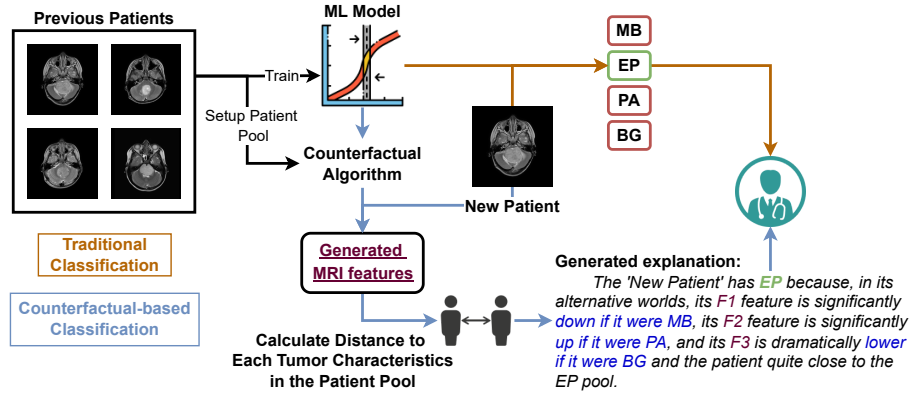


Fig. 3: An overview of our approach demonstrating the generation of counterfactual explanations.

demonstrating the generation of counterfactual explanations. By utilizing all four tumor types, we essentially construct a decision space of reality with our existing patients. As the new patient is guided through this space, attempting to transform into each disease sequentially, the degree of self-modification required for each specific tumor condition will vary. As the required changes decrease, it can be inferred that the patient is closer to that particular tumor type since they necessitate fewer modifications. Similarly, understanding the level of dissimilarity and the contributing features to this dissimilarity has been explored as a critical approach in determining the tumor type. In our previous study [51] the LR model excelled in binary classification and generating counterfactuals, thereby justifying its selection for our current research.

Figure 1 and Table 1 present a scenario involving a patient diagnosed via MRI with an indeterminate tumor type. Using the MRI data, we generate "what-if" scenarios for each tumor type. These scenarios, based on feature distances, help determine the tumor type the patient's data aligns with.

Factual ( $x$ )													
Tumor Type	T2	T2_R	FLAIR	FLAIR_R	DWI	DWI_R	ADC	ADC_R	T1	T1_R	T1CE	T1CE_R	
unknown (EP)	1286	1.529	1311	1.341	1175	1.088	1.009	1.771	473	0.84	892	1.595	
Counterfactual ( $x_{cf}$ )													
Tumor Type	T2	T2_R	FLAIR	FLAIR_R	DWI	DWI_R	ADC	ADC_R	T1	T1_R	T1CE	T1CE_R	
MB	-	-	648	-	-	-	-	-	-	-	1309.5	-	
EP	1423.2	-	-	-	-	-	-	-	-	-	-	-	
PA	2290.2	-	-	-	-	-	2	-	-	-	1492.5	-	
BG	-	-	-	-	544.23	-	2	-	-	-	-	0.781	

Table 1: This table presents the results of our proposed method utilizing counterfactuals (Fig. 1).

Table 1 details an unknown patient classified as EP. For the MB counterfactual, changes in FLAIR\_Tumor and T1CE\_Tumor result in distances of -663 and 417.5, respectively. For EP, T2\_Tumor changes, resulting in a distance of 137.2. In the PA group, changes are observed in T2\_Tumor (from 1286 to 2290.2), ADC\_Tumor (from 1.009 to 2), and T1CE\_Tumor (from 892 to 1492.5). For the BG group, DWI\_Tumor changes from 1175 to 544.23, ADC\_Tumor from 1.009 to 2, and T1CE\_Ratio from 1.595 to 0.781.

In Tables 1 and 2, T represents Tumor and R represents Ratio (Tumor/Parenchyma). The symbol (-) indicates no modification. In Table 1, the patient, with the lowest feature distance to EP, is predicted as EP.

Table 2 displays four new samples after standardization of features and distance calculation. Distances are provided for each patient’s tumor type and their respective counterfactuals. The data represented by (-) matches original values. The distance metric adjusts the distance magnitude between them.

	Tumor Type	T2	T	T2	R	FLAIR	T	FLAIR	R	DWI	T	DWI	R	ADC	T	ADC	R	T1	T	T1	R	T1CE	T	T1CE	R	Distance
Factual ( $x$ )	unknown (MB)	0	-0.5	-0.5		0.5	0	0	-0.5	-1.191	0	0	0	0	-											
Counterfactual ( $x_{cf}$ )	MB	-	-	-	-2	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	<b>2.5</b>
Counterfactual ( $x_{cf}$ )	EP	-	-	2	-	-	-	-	-	0.370	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	2.948
Counterfactual ( $x_{cf}$ )	PA	-	2	-	-	-	-	-	-	0.993	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	3.319
Counterfactual ( $x_{cf}$ )	BG	-	-	-	-	-	-	-	-	2	1.020	-	-	-	-	-	-	-	-	-	-	-	-	-	-	3.337
Factual ( $x$ )	unknown (EP)	-0.583	0	0.5	0	0.5	0	0	-0.816	0	0	0	0	-0.795	0.5	-										
Counterfactual ( $x_{cf}$ )	MB	-	-	-2	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	2.985
Counterfactual ( $x_{cf}$ )	EP	-0.233	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	<b>0.350</b>
Counterfactual ( $x_{cf}$ )	PA	1.982	-	-	-	-	-	-	-	1.225	-	-	-	-	1.550	-	-	-	-	-	-	-	-	-	-	4.031
Counterfactual ( $x_{cf}$ )	BG	-	-	-	-	-	-	-2	-	1.225	-	-	-	-	-	-	-	-	-	-	-	-	-	-2	-	4.082
Factual ( $x$ )	unknown (PA)	0.5	1.223	0	0.5	0.5	0.124	0.816	0	0	0.5	0	0.5	-												
Counterfactual ( $x_{cf}$ )	MB	-	-0.871	-	-2	-	-	-1.225	-	-	-2	-	-	-	-	-	-	-	-	-	-	-	-	-	-	4.588
Counterfactual ( $x_{cf}$ )	EP	-2	-0.869	-	-	-2	-1.749	-1.225	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	4.955
Counterfactual ( $x_{cf}$ )	PA	-	-	-	-	-	1.377	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	<b>1.252</b>
Counterfactual ( $x_{cf}$ )	BG	-	-0.705	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-2	-	-	3.157
Factual ( $x$ )	unknown (BG)	0.811	0.5	-0.373	0.811	0	0	0.831	0.5	-0.5	0.5	-0.602	-0.5	-												
Counterfactual ( $x_{cf}$ )	MB	-1.033	-	-0.847	-1.393	-	-	-1.079	-2	2	-2	-0.166	2	6.109												
Counterfactual ( $x_{cf}$ )	EP	-1.400	-2	1.966	-	-	-	-1.360	-	-	-	-	-	4.627												
Counterfactual ( $x_{cf}$ )	PA	-	-	-	-	-	-	0.778	-	-	-	1.971	-	2.574												
Counterfactual ( $x_{cf}$ )	BG	-	-	-	-1.041	-	-	-	-	-	-	-	-	<b>1.853</b>												

Table 2: Distance results for counterfactuals generated on feature-wise scaled data for four distinct newly arriving patients with varying tumor types.

## 4.2 Revealing Key MRI Features through Counterfactual Explanations

As discussed in Section 3.1, counterfactual explanations can provide insights into feature importance. These explanations allow us to understand the reasoning behind ML model decisions and offer valuable options for restriction. In clinical settings, visible changes in features through counterfactual explanations can be more relevant and meaningful for real-world evaluations and applications.

Considering that we generated five counterfactuals for each patient, we obtained 125 explanations for MB, PA, and BG, and 55 explanations for EP. Table 3 illustrates our reporting method for counterfactual analysis results for a case

scenario (MB to EP). The patient count, the total number of generated counterfactual explanations for them, and the statistical information regarding the frequency of changes observed on which features in these counterfactuals to identify the top 3 influential features are shown. For instance, "FLAIR\_Tumor 71 changes" signifies that out of 125 counterfactuals, 71 of them involved a modification from MB to EP. Therefore, FLAIR\_Tumor creates such a distinction between these two tumors that the model considers altering this feature significantly influential in shifting the decision from one side to the other in the decision space. The greater the repetition of this occurrence, indicated by the magnitude of "changes," the more pronounced the outcome suggesting that even in random selections, optimization is achieved for that particular feature, significantly impacting the decision. Table 4 presents the findings from each tumor

Number of patients: 25			
Number of generated counterfactuals: 125			
FLAIR_Tumor	71 changes	T1_Ratio	6 changes
ADC_Tumor	33 changes	T1CE_Tumor	6 changes
ADC_Ratio	29 changes	T2_Tumor	3 changes
DWI_Ratio	18 changes	T2_Parenchyma	0 changes
FLAIR_Ratio	17 changes	FLAIR_Parenchyma	0 changes
DWI_Tumor	12 changes	DWI_Parenchyma	0 changes
T1_Tumor	10 changes	ADC_Parenchyma	0 changes
T1CE_Ratio	7 changes	T1_Parenchyma	0 changes
T2_Ratio	6 changes	T1CE_Parenchyma	0 changes

Table 3: This example analysis presents the variations in characteristics observed during the generation of counterfactual instances for the transition from MB to EP.

pair to identify feature differences between different tumor types. The observed changes in features align with expected outcomes from clinical studies. MB and EP tumors are distinguished by FLAIR and ADC features. MB and PA typically exhibit differences in T2 and ADC. MB and BG, on the other hand, show variations primarily in ADC, T2, and T1CE. In the case of EP and PA, T2 exhibits the most significant changes, while variations in ADC and T1CE are also observed. The most distinguishing features between EP and BG are T1CE\_Ratio and ADC\_Tumor. As for PA and BG, the T2\_Ratio feature has been identified as a crucial factor in creating differentiation. Additionally, significant variations in T1CE features are frequently observed, further contributing to the dissimilarity between these tumor types.

The results illustrated in Figure 4 provide a clear view of the original data distributions for MB-PA and BG-PA, shedding light on critical MRI features



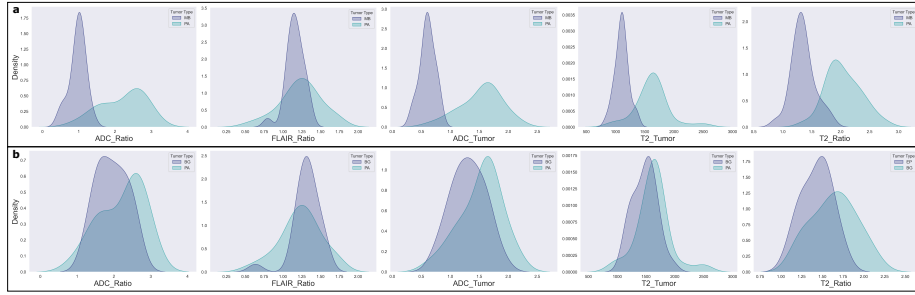


Fig. 4: The original data distributions for MB-PA and BG-PA, highlighting key features for MB-PA and their behavior when applied to BG-PA. The figure demonstrates how the top 5 features vary significantly between MB and PA, while BG and PA distributions remain almost identical, lacking discernible features.

through counterfactual explanations. This figure, along with Table 4, reveals that transitions between MB and PA generally exhibit similar overall distributions. However, there are significant differences in the top 5 features, which aligns with our previous study [51]. In contrast, the distributions for BG and PA are nearly identical, with no distinct features. This lack of variation is also evident in the top 5 features, demonstrating the effectiveness of our algorithm in generating counterfactuals by selecting features that show the most variation, thus achieving maximum impact with minimal changes. Notably, T1CE features and T2\_Ratio stand out as the most distinctive between PA and BG, as previously highlighted in Figure 5 of our earlier research [51].

MB to EP		MB to PA		MB to BG		EP to MB		EP to PA		EP to BG	
Feature	Change	Feature	Change	Feature	Change	Feature	Change	Feature	Change	Feature	Change
FLAIR_Tumor	71	T2_Ratio	87	T2_Tumor	64	T1CE_Tumor	18	T2_Tumor	34	ADC_Tumor	22
ADC_Tumor	33	T2_Tumor	55	ADC_Tumor	52	FLAIR_Ratio	16	T2_Ratio	26	DWI_Ratio	16
ADC_Ratio	29	ADC_Tumor	43	T1CE_Ratio	43	FLAIR_Tumor	13	ADC_Tumor	19	T1CE_Ratio	15
PA to MB		PA to EP		PA to BG		BG to MB		BG to EP		BG to PA	
Feature	Change	Feature	Change	Feature	Change	Feature	Change	Feature	Change	Feature	Change
ADC_Ratio	91	T2_Ratio	95	T2_Ratio	95	FLAIR_Ratio	85	T1CE_Ratio	90	T1CE_Ratio	53
FLAIR_Ratio	88	T2_Tumor	81	T1CE_Tumor	66	ADC_Tumor	71	T1_Tumor	50	T1CE_Tumor	48
ADC_Tumor	76	T1CE_Tumor	45	T1CE_Ratio	54	ADC_Ratio	71	DWI_Ratio	48	T2_Ratio	32

Table 4: The three most important features for each changing reality case.

### 4.3 Statistical Analysis of Generated Counterfactuals

To validate the statistical fidelity of the counterfactual tumors generated from the original tumors, we performed a dependent T-test to evaluate the statistical

difference between the tumor instances as explained in Section 2.5. Although the generated counterfactuals are expected to have similar statistical properties to samples from the same class due to their minimal distance from the original instances, it is important to verify this with statistical analysis, as small borderline feature changes can lead to different classes. The statistical similarity obtained when each tumor is transformed to represent the "what if?" scenario of other tumors. More specifically, when we transform tumor  $x$  to tumor  $x'=y$ , we know that  $x'$  is still dependent on  $x$ . Therefore, we measure how similar  $x'$  is to the original distribution of  $y$  on the feature where  $x$  undergoes the most significant change. The results of statistical hypothesis tests are presented in Table 5.

A high  $p$ -value indicates that we do not reject the difference, implying that the counterfactual explanations we generate sufficiently resemble the original distribution for that particular feature. In other words, for many features the difference between the generated sample and the actual data distribution is statistically insignificant. This is a desirable property for data augmentation as it means that the augmented data samples resemble similar statistical distributions.

MRI Feature	Original	Generated	T-Statistic	P-Value
DWI_Tumor	MB	MB to MB	-0.0605	0.9521
T1CE_Ratio	MB	MB to MB	-0.3643	0.7177
T1_Ratio	MB	MB to MB	-0.0282	0.9776
T1CE_Tumor	MB	EP to MB	-2.4975	0.0156
FLAIR_Ratio	MB	EP to MB	0.4532	0.6524
FLAIR_Tumor	MB	EP to MB	0.7917	0.4331
ADC_Ratio	MB	PA to MB	-0.4017	0.6887
FLAIR_Ratio	MB	PA to MB	11.5615	<0.0001
ADC_Tumor	MB	PA to MB	-4.3706	<0.0001
FLAIR_Ratio	MB	BG to MB	6.7026	<0.0001
ADC_Tumor	MB	BG to MB	-5.2062	<0.0001
ADC_Ratio	MB	BG to MB	-2.7487	0.0071

(a) Difference between original MB and generated MBs.

MRI Feature	Original	Generated	T-Statistic	P-Value
FLAIR_Tumor	EP	MB to EP	-3.2397	0.0061
ADC_Tumor	EP	MB to EP	-2.0273	0.0495
ADC_Ratio	EP	MB to EP	-0.6434	0.5266
FLAIR_Tumor	EP	EP to EP	-1.0603	0.3018
ADC_Tumor	EP	EP to EP	-1.4653	0.1519
DWI_Ratio	EP	EP to EP	-0.4937	0.6278
T2_Ratio	EP	PA to EP	2.956	0.007
T2_Tumor	EP	PA to EP	0.3621	0.72
T1CE_Tumor	EP	PA to EP	-1.1672	0.262
T1CE_Ratio	EP	BG to EP	-2.1967	0.0428
T1_Tumor	EP	BG to EP	-5.9549	<0.0001
DWI_Ratio	EP	BG to EP	0.0059	0.9954

(b) Difference between original EP and generated EPs.

MRI Feature	Original	Generated	T-Statistic	P-Value
T2_Ratio	PA	MB to PA	-2.0667	0.0430
T2_Tumor	PA	MB to PA	-0.1256	0.9004
ADC_Tumor	PA	MB to PA	4.4019	<0.0001
T2_Tumor	PA	EP to PA	-1.3925	0.1694
T2_Ratio	PA	EP to PA	1.6279	0.1091
ADC_Tumor	PA	EP to PA	3.6692	0.0005
ADC_Tumor	PA	PA to PA	0.2781	0.7822
T1_Ratio	PA	PA to PA	-0.8599	0.3948
FLAIR_Tumor	PA	PA to PA	-1.2497	0.2176
T1CE_Ratio	PA	BG to PA	1.711	0.0944
T1CE_Tumor	PA	BG to PA	1.7524	0.0862
T2_Ratio	PA	BG to PA	2.2029	0.0326

(c) Difference between original PA and generated PAs.

MRI Feature	Original	Generated	T-Statistic	P-Value
T2_Tumor	BG	MB to BG	-2.4723	0.0149
ADC_Tumor	BG	MB to BG	5.6539	<0.0001
T1CE_Ratio	BG	MB to BG	-6.6115	<0.0001
ADC_Tumor	BG	EP to BG	2.8215	0.0066
DWI_Ratio	BG	EP to BG	0.806	0.4236
T1CE_Ratio	BG	EP to BG	-4.4199	<0.0001
T2_Ratio	BG	PA to BG	6.78	<0.0001
T1CE_Tumor	BG	PA to BG	-5.3162	<0.0001
T1CE_Ratio	BG	PA to BG	-6.9185	<0.0001
ADC_Tumor	BG	BG to BG	-0.3252	0.7467
DWI_Tumor	BG	BG to BG	-0.8181	0.4176
DWI_Ratio	BG	BG to BG	-0.7461	0.4599

(d) Difference between original BG and generated BGs.

Table 5: The results of hypothesis tests comparing the original data with the generated data.

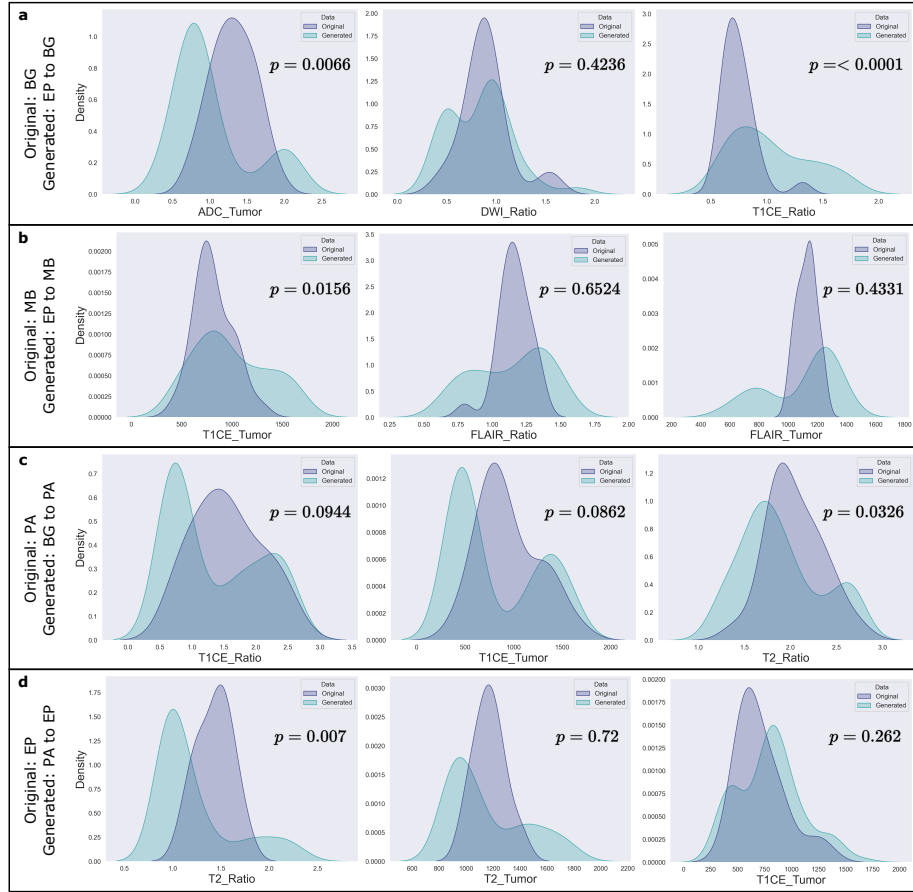


Fig. 5: The distributions of the original data and the generated data.

Apart from the PA to MB transition (e.g.,  $p=0.04763$ ,  $p=0.0307$ ), no significant differences were observed in other tumor transitions. This result can be attributed to both the fundamental optimization principle of minimizing changes during counterfactual generation and the distribution distances shown in Fig. 4. Specifically, Fig. 4a demonstrates a distinct separation in the distributions during the PA to MB transition, requiring a significantly larger change for transformation.

As expected, when attempting self-transformation on each tumor type, the obtained  $p$ -values were notably high. Evaluating at a significance level of 0.05, several features closely aligned with the actual feature distribution of the patients, making them indistinguishable from the ground truth. The following features exhibited this characteristic: FLAIR\_Ratio and FLAIR\_Tumor in the case of transforming EP to MB, ADC\_Ratio when transforming PA to MB,

ADC\_Ratio during the transformation from MB to EP, T2\_Tumor and T1CE\_Tumor in the context of PA to EP transformation, DWI\_Ratio when transforming BG to EP, T2\_Tumor for MB to PA transformation, T2\_Tumor and T2\_Ratio in the case of EP to PA transformation, T1CE\_Ratio and T1CE\_Tumor during BG to PA transformation, and DWI\_Ratio when transforming EP to BG. Fig. 5 presents some of these cases along with their KDE distributions.

#### 4.4 Pushing the Boundaries of Data Augmentation through Alternative Realities

During the construction of counterfactuals, we employed downsampling for MB and BG to align with the number of PA patients (25) during training, considering it appropriate. EP had a count of 11, and we did not increase it. The baseline results for this scenario can be observed in Table 6. For evaluation, the train-test splitting was conducted with a ratio of 45% for the baseline dataset, 35% for EP augmentation, and 25% for EP-PA-BG augmentation.

	Training set	Test set	Precision	Recall	F1 Score
Baseline	47(Real only)	39(Real only)	$73.15 \pm 9.48$	$72.20 \pm 4.78$	$71.28 \pm 5.62$
A	65(56 R, 9 CF)	35(30 R, 5 CF)	$84.83 \pm 4.95$	$83.75 \pm 3.72$	$83.34 \pm 3.65$
B	126(84 R, 42 CF)	42(28 R, 14 CF)	$86.31 \pm 4.57$	$84.64 \pm 4.69$	$84.85 \pm 4.72$
C	124(68 R, 56 CF)	44(Real only)	73.58	72.73	72.04

Table 6: The impact of data augmentation using counterfactuals on classification scores is presented in the table. In the table, CF stands for counterfactuals whereas R is abbreviation for real samples. **(A)** For the first augmentation scenario, only EP counterfactuals are added, resulting in a dataset with 25 samples each for MB, EP, PA, and BG. **(B)** In the second augmentation scenario, counterfactuals for EP, PA, and BG are added to balance the number of samples with original count of MB. Assuming all counterfactual examples represent real data, this scenario results from a dataset comprising 42 samples each for MB, EP, PA, and BG. **(C)** The third scenario involves moving all real samples to the test set, with 11 patients in each category. Consequently, no factual EP samples are left in the training set, and the model is trained accordingly.

To address the data imbalance, we examined the inclusion of generated counterfactuals for data augmentation. For example, by equalizing EP with the other tumor types and incorporating 14 different generated counterfactuals alongside the originals, we excluded EP-to-EP instances. Opting for transitions from various tumor types to maximize variance and generalizability, we achieved an improvement of up to 12.06% as shown in Table 6, case A.

To incorporate the previously set aside MB and BG data, we aligned all tumor types, except themselves, with counterfactuals generated from different

tumor types. BG, PA, and EP were included with MB, and all were evaluated as a group of 42 patients, which was the maximum patient count for one tumor type. When considering the counterfactuals as actual patients, the outcomes align with the results presented in Table 6, case B.

Furthermore, in the case examined in Table 6, case C, 11 patients were included from each tumor type in the test set, resulting in no actual EP patients in the training set. Consequently, during training, we had 31 real samples for MB, 0 real and 31 counterfactual samples for EP, 14 real and 17 counterfactual samples for PA, and 23 real and 8 counterfactual samples for BG. Notably, when evaluating on real samples, the results were intriguing. Despite the absence of real EP patients in the training data, the model successfully identified 5 out of the 11 patients, leading to an overall baseline score that was, on average, 0.76% higher.

In all cases, LR consistently yielded the best performance in terms of prediction accuracy and execution time. The boosting algorithms took significantly more time to converge on the predictions, although they did not improve the prediction performance. Therefore, to simplify the presentation of the results, we only reported the results of the LR classifier in the table.

## 5 Discussion

Spatial heterogeneity in pediatric brain tumors, especially from the posterior fossa, complicates accurate differentiation [9, 13, 14, 16, 26, 34, 37, 40, 41]. Accurate diagnoses are essential since each tumor requires specific treatments impacting patient outcomes. While AI advancements in medical imaging are promising, their black-box nature hinders clinical adoption. Our study introduces a novel approach, leveraging counterfactual explanations to interpret MRI features, aiming to provide clinicians an intuitive tool. This research pioneers feature-based counterfactual investigations in pediatric brain tumors.

The medical literature highlights that individualized care is crucial, aligning with personalized healthcare [8, 23, 38, 47]. The ability to generate hypothetical scenarios for a patient based on MRI features offers a significant advantage in medical diagnosis. By creating what-if scenarios, radiologists are equipped with additional intuitive data, enhancing their decision-making ability. In the proposed approach, counterfactual explanations can be seamlessly integrated into the radiologist’s existing clinical workflow as a decision support system and provide support to the personalized treatment process for patients. This approach could potentially mitigate the need for invasive procedures and provide a clearer perspective on the tumor’s nature based on MRI data alone. Our method produces tailored explanations for each patient, drawing from past cases, facilitating understanding of tumor differentiation based on MRI data.

The concept of counterfactuals, which has been debated in philosophy and psychology for decades, has found its place in the field of artificial intelligence under various names. Though the idea has historical roots, its comprehensive implementation in AI is a more recent phenomenon. Like many preceding stud-

ies, we have adapted this concept for the clinical domain, providing valuable insights for clinicians. Furthermore, our work enriches the literature on medical counterfactuals by offering a unique perspective tailored to specific tasks. Through counterfactuals, we demonstrate alternative possibilities within the decision space and elucidate the rationales behind specific decisions pertaining to pediatric patients.

To the best of our knowledge, no prior studies have addressed counterfactuals regarding posterior fossa tumors. We filled this gap and subjected results to statistical tests, presented in Section 4.3. We explored the utility of counterfactuals both as post-classifiers and indicators of significant MRI features. The LR model was the most effective, hence we used it for counterfactual generation.

We developed a framework aimed at enhancing the utility of counterfactuals beyond what DiCE offers for our case. When faced with a sizable patient pool, utilizing another counterfactual algorithm considering a subsample and substituting excluded patients can aid in statistical testing. In certain cases, the alternate methods could optimize within a more favorable timeframe than DiCE.

Our approach, which utilizes counterfactual explanations in a classifier-like manner, eliminates the need to separate different test sets. Consequently, the performance of our machine learning models significantly exceeds the baseline scores, with only a few patients excluded from the decision space to simulate the scenario of newly arriving patients. In this scenario, all training samples serve as test patients as we explore the decision space. To accomplish this, DiCE provides valuable information about misclassified samples, allowing us to exclude the associated counterfactuals from the statistical analysis through post-processing.

Fig. 1 and Table 1 depict a hypothetical scenario involving a patient with an initially unknown EP tumor. The radiologist examining the MR images was uncertain about whether the tumor was of the MB or EP type. A key challenge in such cases is the lack of additional information, which often necessitates invasive procedures like brain surgery and tissue sampling for histopathological analysis to obtain a definitive diagnosis. To overcome this issue, we generate alternative scenarios based solely on the MRI features. These scenarios provide additional quantitative information to the radiologist, enabling them to assess the response based on the individual’s biological characteristics.

Moreover, Table 2 demonstrates the efficacy of our approach in identifying patients with diverse tumor types that were previously unidentified and not encompassed within the decision space. While ML models can also accomplish this task, our method offers an additional advantage by preserving information regarding tissue characteristics, which in turn reveal similarities or differences among tumors. Additionally, our approach calculates distances to other tumors by transforming the features into a uniform distribution through standard scaling, providing valuable insights about the proximity. This valuable information aids in our comprehension of the differentiation among tumors in the dataset.

Table 3 presents the total count of modifications made to susceptible features, with the exception of the parenchymas that serves as reference points, when generating samples for different patients. The statistical report enables

a human verification of the optimization process, wherein minimal changes are implemented to achieve the desired outcome. It also confirms that the features exhibiting the highest variations during the generation of alternative realities are those with the most distinct distributions between two tumors. To elucidate the analysis of their distributions, we present Fig. 4 as a visual representation. Table 4 presents the top three most variable features extracted from the reports obtained for all tumor matches in Table 3. This provides a more concise overview of all the cases.

Table 5 exhibits a statistical analysis demonstrating the high degree of similarity between the generated data and reality across different data spaces, specifically focusing on the most frequently selected features. A high  $p$ -value indicates that the generated samples cannot be well distinguished, implying the effectiveness of the independent transformation process, which produces significant alternative realities separate from the original space. Fig. 5 illustrates an example of some transformations from Table 5, displaying their corresponding  $p$ -values, as well as the kernel density estimation of the generated data in comparison to the original data.

Our generated counterfactuals offer potential advantages for data augmentation. Traditional methods, like SMOTE [7], fall short in real-world alignment and interpretability. We suggest counterfactuals as a viable alternative, particularly when data is limited. As shown in Fig. 6, while case C cannot be benchmarked directly due to added test patients, the inclusion of more real samples enhances outcomes. This improvement aligns with findings in [51]. However, challenges arise when certain EP patients, which complicate differentiation, are considered. Despite these challenges, the ability to make accurate predictions for many patients without actual EP training data highlights the promise of our approach and suggests directions for future research.

Counterfactual explanations can also address model bias in medical diagnoses [31, 59]. Ensuring fairness and transparency in decision-making processes is vital, suggesting the value of counterfactuals in this direction.

### 5.1 Potential Challenges and Limitations

One of the limitations encountered in this study is the size of MRI data collected. Considering the low incidence of this disease, it is difficult to encounter a large number of patients for each type of pediatric posterior fossa tumor in a single hospital. Although we observe that the current data set provides valid counterfactual explanations and clinically useful results, there is still a need for large research data sets collected at national or international level in this field.

Expanding the dataset’s scope and size may lead to the potential for counterfactual explanations to emerge with new and different latent features. Furthermore, the use of a comprehensive data set that captures a broader range of scenarios for pediatric posterior fossa tumors encountered in clinical practice may play an important role in ensuring a wider applicability and confirming the generalizability of the counterfactual explanations currently uncovered.

There are recognized challenges in applying the DiCE method to diverse datasets, sometimes resulting in extended optimization times and difficulties achieving convergence. Addressing these challenges will be an essential step forward. Exploring alternative methodologies and delving deeper into the vast landscape of counterfactual algorithms might also be beneficial.

To further advance the field, future research should consider incorporating additional advanced MRI protocols, enriching our understanding and diagnostic capabilities regarding pediatric posterior fossa tumors.

## 6 Conclusion

This paper presents a novel interpretability approach in medical research, using pediatric posterior fossa brain tumors as a case study. By generating counterfactual explanations, it delivers personalized insights, validates predicted outcomes, and highlights how predictions vary under different conditions. Although medical regulations and workflow concerns remain hurdles, the use of explainable AI in medicine is poised to grow as its benefits become clearer.

Our method bridges the gap between machine learning and clinical decision-making, potentially leading to better patient outcomes. Further research is needed to integrate counterfactual explanations into clinical practice and evaluate their real-world performance. Larger studies, including different diseases, could produce even more robust “alternative realities” from MRI features. Overall, we believe this approach has the potential to shift the perspective of radiologists and other medical professionals by offering more human-like, actionable insights.

**Acknowledgments.** The authors received no financial support for the research, authorship, and/or publication of this article.

**Disclosure of Interests.** The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

## Institutional Review Board Statement

After obtaining approval from the Institutional Review Board of Children Hospital of 02 with approval number [Ref: 632 QĐ-NĐ2 dated 12 May 2019], we conducted the study in both Radiology and Neurosurgery departments in accordance with the 1964 Helsinki declaration.

## Data & Code Availability

The datasets generated and/or analyzed during the current study are not publicly available due to privacy concerns but are available from the Dr. Keserci upon reasonable request. The source codes of the presented study can be accessed at:

<https://github.com/tanyelai/counterfactual-explanations-for-medical-research>



## References

1. Andini, M., Ciani, E., De Blasio, G., D'Ignazio, A., Salvestrini, V.: Targeting policy-compliers with machine learning: an application to a tax rebate programme in Italy. Bank of Italy Temi di Discussione (Working Paper) No **1158** (2017)
2. Athey, S.: Beyond prediction: Using big data for policy problems. *Science* **355**(6324), 483–485 (2017)
3. Avanzo, M., Wei, L., Stancanella, J., Vallieres, M., Rao, A., Morin, O., Mattonen, S.A., El Naqa, I.: Machine and deep learning methods for radiomics. *Medical physics* **47**(5), e185–e202 (2020)
4. Band, S.S., Yarahmadi, A., Hsu, C.C., Biyari, M., Sookhak, M., Ameri, R., Dehzangi, I., Chronopoulos, A.T., Liang, H.W.: Application of explainable artificial intelligence in medical health: A systematic review of interpretability methods. *Informatics in Medicine Unlocked* p. 101286 (2023)
5. Borys, K., Schmitt, Y.A., Nauta, M., Seifert, C., Krämer, N., Friedrich, C.M., Nensa, F.: Explainable ai in medical imaging: An overview for clinical practitioners—beyond saliency-based xai approaches. *European journal of radiology* p. 110786 (2023)
6. Brown, K.E., Talbert, D., Talbert, S.: The uncertainty of counterfactuals in deep learning. In: *The International FLAIRS Conference Proceedings*. vol. 34 (2021)
7. Chawla, N.V., Bowyer, K.W., Hall, L.O., Kegelmeyer, W.P.: Smote: synthetic minority over-sampling technique. *Journal of artificial intelligence research* **16**, 321–357 (2002)
8. Chawla, N.V., Davis, D.A.: Bringing big data to personalized healthcare: a patient-centered framework. *Journal of general internal medicine* **28**, 660–665 (2013)
9. Chen, D., Lin, S., She, D., Chen, Q., Xing, Z., Zhang, Y., Cao, D.: Apparent diffusion coefficient in the differentiation of common pediatric brain tumors in the posterior fossa: Different region-of-interest selection methods for time efficiency, measurement reproducibility, and diagnostic utility. *Journal of Computer Assisted Tomography* **47**(2), 291 (2023)
10. Chou, Y.L., Moreira, C., Bruza, P., Ouyang, C., Jorge, J.: Counterfactuals and causability in explainable artificial intelligence: Theory, algorithms, and applications. *Information Fusion* **81**, 59–83 (2022)
11. Cinà, G., Röber, T., Goedhart, R., Birbil, I.: Why we do need explainable ai for healthcare (2022)
12. Dai, W., Brisimi, T.S., Adams, W.G., Mela, T., Saligrama, V., Paschalidis, I.C.: Prediction of hospitalization due to heart diseases by supervised learning methods. *International journal of medical informatics* **84**(3), 189–197 (2015)
13. Duc, N.M., Huy, H.Q.: Magnetic resonance imaging features of common posterior fossa brain tumors in children: a preliminary vietnamese study. *Open Access Macedonian Journal of Medical Sciences* **7**(15), 2413 (2019)
14. Duc, N.M., Huy, H.Q., Nadarajan, C., Keserci, B.: The role of predictive model based on quantitative basic magnetic resonance imaging in differentiating medulloblastoma from ependymoma. *Anticancer Research* **40**(5), 2975–2980 (2020)
15. Dutta, S., Long, J., Mishra, S., Tilli, C., Magazzeni, D.: Robust counterfactual explanations for tree-based ensembles. In: *International Conference on Machine Learning*. pp. 5742–5756. PMLR (2022)
16. D'Arco, F., Khan, F., Mankad, K., Ganau, M., Caro-Dominguez, P., Bisdas, S.: Differential diagnosis of posterior fossa tumours in children: new insights. *Pediatric Radiology* **48**, 1955–1963 (2018)

17. Gettier, E.L.: Is justified true belief knowledge? *Analysis* **23**(6), 121–123 (1963). <https://doi.org/10.1093/analys/23.6.121>
18. Ginsberg, M.L.: Counterfactuals. *Artificial intelligence* **30**(1), 35–79 (1986)
19. Griffiths, T.L., Chater, N., Kemp, C., Perfors, A., Tenenbaum, J.B.: Probabilistic models of cognition: Exploring representations and inductive biases. *Trends in cognitive sciences* **14**(8), 357–364 (2010)
20. Guidotti, R.: Counterfactual explanations and how to find them: literature review and benchmarking. *Data Mining and Knowledge Discovery* pp. 1–55 (2022)
21. Hitchcock, C.: The intransitivity of causation revealed in equations and graphs. *The Journal of Philosophy* **98**(6), 273–299 (2001)
22. Hitchcock, C.: Prevention, preemption, and the principle of sufficient reason. *The Philosophical Review* **116**(4), 495–532 (2007)
23. Johnson, K.B., Wei, W.Q., Weeraratne, D., Frisse, M.E., Misulis, K., Rhee, K., Zhao, J., Snowdon, J.L.: Precision medicine, ai, and the future of personalized health care. *Clinical and translational science* **14**(1), 86–93 (2021)
24. Kment, B.: Counterfactuals and explanation. *Mind* **115**(458), 261–310 (2006)
25. Knapič, S., Malhi, A., Saluja, R., Främling, K.: Explainable artificial intelligence for human decision support system in the medical domain. *Machine Learning and Knowledge Extraction* **3**(3), 740–770 (2021)
26. Koob, M., Girard, N.: Cerebral tumors: specific features in children. *Diagnostic and interventional imaging* **95**(10), 965–983 (2014)
27. Lewis, D.K.: Counterfactuals. Cambridge, MA, USA: Blackwell (1973)
28. Lin, Q.H., Niu, Y.W., Sui, J., Zhao, W.D., Zhuo, C., Calhoun, V.D.: Sspnet: An interpretable 3d-cnn for classification of schizophrenia using phase maps of resting-state complex-valued fmri data. *Medical Image Analysis* **79**, 102430 (2022)
29. Lundberg, S.M., Lee, S.I.: A unified approach to interpreting model predictions. *Advances in neural information processing systems* **30** (2017)
30. Maragno, D., Kurtz, J., Röber, T.E., Goedhart, R., Birbil, S.I., den Hertog, D.: Finding regions of counterfactual explanations via robust optimization (2023)
31. Mikołajczyk, A., Grochowski, M., Kwasigroch, A.: Towards explainable classifiers using the counterfactual approach: global explanations for discovering bias in data. *Journal of Artificial Intelligence and Soft Computing Research* **11**(1), 51–67 (2021)
32. Miller, T.: Explanation in artificial intelligence: Insights from the social sciences. *Artificial intelligence* **267**, 1–38 (2019)
33. Mittelstadt, B., Russell, C., Wachter, S.: Explaining explanations in ai. In: *Proceedings of the conference on fairness, accountability, and transparency*. pp. 279–288 (2019)
34. Moharamzad, Y., Sanei Taheri, M., Niaghi, F., Shobeiri, E.: Brainstem glioma: Prediction of histopathologic grade based on conventional mr imaging. *The neuro-radiology journal* **31**(1), 10–17 (2018)
35. Mothilal, R.K., Sharma, A., Tan, C.: Explaining machine learning classifiers through diverse counterfactual explanations. In: *Proceedings of the 2020 conference on fairness, accountability, and transparency*. pp. 607–617 (2020)
36. Nagesh, S., Mishra, N., Naamad, Y., Rehg, J.M., Shah, M.A., Wagner, A.: Explaining a machine learning decision to physicians via counterfactuals. In: *Conference on Health, Inference, and Learning*. pp. 556–577. PMLR (2023)
37. Orphanidou-Vlachou, E., Vlachos, N., Davies, N.P., Arvanitis, T.N., Grundy, R.G., Peet, A.C.: Texture analysis of t1-and t2-weighted mr images and use of probabilistic neural network to discriminate posterior fossa tumours in children. *NMR in Biomedicine* **27**(6), 632–639 (2014)

38. Paranjape, K., Schinkel, M., Nanayakkara, P.: Short keynote paper: Mainstreaming personalized healthcare—transforming healthcare through new era of artificial intelligence. *IEEE journal of biomedical and health informatics* **24**(7), 1860–1863 (2020)
39. Pawlowski, N., Coelho de Castro, D., Glocker, B.: Deep structural causal models for tractable counterfactual inference. *Advances in Neural Information Processing Systems* **33**, 857–869 (2020)
40. Porto, L., Jurcoane, A., Schwabe, D., Hattingen, E.: Conventional magnetic resonance imaging in the differentiation between high and low-grade brain tumours in paediatric patients. *European Journal of Paediatric Neurology* **18**(1), 25–29 (2014)
41. Reddy, N., Ellison, D.W., Soares, B.P., Carson, K.A., Huisman, T.A., Patay, Z.: Pediatric posterior fossa medulloblastoma: the role of diffusion imaging in identifying molecular groups. *Journal of Neuroimaging* **30**(4), 503–511 (2020)
42. Ribeiro, M.T., Singh, S., Guestrin, C.: " why should i trust you?" explaining the predictions of any classifier. In: *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*. pp. 1135–1144 (2016)
43. Rockoff, J.E., Jacob, B.A., Kane, T.J., Staiger, D.O.: Can you recognize an effective teacher when you recruit one? *Education finance and Policy* **6**(1), 43–74 (2011)
44. Ruben, D.H.: *Explaining explanation*. Routledge (2015)
45. Sanchez, P., Kascenas, A., Liu, X., O’Neil, A.Q., Tsaftaris, S.A.: What is healthy? generative counterfactual diffusion for lesion localization. In: *MICCAI Workshop on Deep Generative Models*. pp. 34–44. Springer (2022)
46. Sarp, S., Catak, F.O., Kuzlu, M., Cali, U., Kusetogullari, H., Zhao, Y., Ates, G., Guler, O.: An xai approach for covid-19 detection using transfer learning with x-ray images. *Heliyon* **9**(4) (2023)
47. Shaban-Nejad, A., Michalowski, M., Buckeridge, D.L.: Health intelligence: how artificial intelligence transforms population and personalized health. *NPJ digital medicine* **1**(1), 53 (2018)
48. Spirtes, P., Glymour, C., Scheines, R.: *Discovery algorithms for causally sufficient structures. Causation, prediction, and search* pp. 103–162 (1993)
49. Spirtes, P., Glymour, C., Scheines, R.: *Causation, prediction, and search*. MIT press (2000)
50. Starr, W.: Counterfactuals. In: Zalta, E.N., Nodelman, U. (eds.) *The Stanford Encyclopedia of Philosophy. Metaphysics Research Lab, Stanford University*, Winter 2022 edn. (2022)
51. Tanyel, T., Nadarajan, C., Duc, N.M., Keserci, B.: Deciphering machine learning decisions to distinguish between posterior fossa tumor types using mri features: What do the data tell us? *Cancers* **15**(16) (2023). <https://doi.org/10.3390/cancers15164015>, <https://www.mdpi.com/2072-6694/15/16/4015>
52. Tjoa, E., Guan, C.: A survey on explainable artificial intelligence (xai): Toward medical xai. *IEEE transactions on neural networks and learning systems* **32**(11), 4793–4813 (2020)
53. Todo, W., Selmani, M., Laurent, B., Loubes, J.M.: Counterfactual explanation for multivariate times series using a contrastive variational autoencoder. In: *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. pp. 1–5. IEEE (2023)
54. Vellido, A.: The importance of interpretability and visualization in machine learning for applications in medicine and health care. *Neural computing and applications* **32**(24), 18069–18083 (2020)

55. Verma, S., Boonsanong, V., Hoang, M., Hines, K.E., Dickerson, J.P., Shah, C.: Counterfactual explanations and algorithmic recourses for machine learning: A review. *arXiv preprint arXiv:2010.10596* (2020)
56. de Vries, B.M., Zwezerijnen, G.J., Burchell, G.L., van Velden, F.H., Menke-van der Houven van Oordt, C.W., Boellaard, R.: Explainable artificial intelligence (xai) in radiology and nuclear medicine: a literature review. *Frontiers in medicine* **10**, 1180773 (2023)
57. Wachter, S., Mittelstadt, B., Russell, C.: Counterfactual explanations without opening the black box: Automated decisions and the gdpr. *Harv. JL & Tech.* **31**, 841 (2017)
58. Wang, Z., Samsten, I., Papapetrou, P.: Counterfactual explanations for survival prediction of cardiovascular icu patients. In: *Artificial Intelligence in Medicine: 19th International Conference on Artificial Intelligence in Medicine, AIME 2021, Virtual Event, June 15–18, 2021, Proceedings.* pp. 338–348. Springer (2021)
59. Wang, Z., Zhou, Y., Qiu, M., Haque, I., Brown, L., He, Y., Wang, J., Lo, D., Zhang, W.: Towards fair machine learning software: Understanding and addressing model bias through counterfactual thinking. *arXiv preprint arXiv:2302.08018* (2023)
60. Waters, A., Mikkulainen, R.: Grade: Machine learning support for graduate admissions. *Ai Magazine* **35**(1), 64–64 (2014)
61. Wong, S.C., Gatt, A., Stamatescu, V., McDonnell, M.D.: Understanding data augmentation for classification: when to warp? In: *2016 international conference on digital image computing: techniques and applications (DICTA).* pp. 1–6. IEEE (2016)
62. Woodward, J.: *Making things happen: A theory of causal explanation.* Oxford university press (2005)
63. Woodward, J., Zalta, E.N., et al.: *Scientific explanation.* The Stanford (2017)
64. Woodward, J.: Explanation, invariance, and intervention. *Philosophy of Science* **64**(S4), S26–S41 (1997)
65. Xu, R., Yu, Y., Zhang, C., Ali, M.K., Ho, J.C., Yang, C.: Counterfactual and factual reasoning over hypergraphs for interpretable clinical predictions on ehr. In: *Machine Learning for Health.* pp. 259–278. PMLR (2022)
66. Yang, L., Kenny, E.M., Ng, T.L.J., Yang, Y., Smyth, B., Dong, R.: Generating plausible counterfactual explanations for deep transformers in financial text classification. *arXiv preprint arXiv:2010.12512* (2020)
67. Yang, Z., Liu, Y., Ouyang, C., Ren, L., Wen, W.: Counterfactual can be strong in medical question and answering. *Information Processing & Management* **60**(4), 103408 (2023)
68. Zeineldin, R.A., Karar, M.E., Elshaer, Z., Coburger, .J., Wirtz, C.R., Burgert, O., Mathis-Ullrich, F.: Explainability of deep neural networks for mri analysis of brain tumors. *International journal of computer assisted radiology and surgery* **17**(9), 1673–1683 (2022)
69. Zhang, H., Cisse, M., Dauphin, Y.N., Lopez-Paz, D.: mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412* (2017)