# Lab 1 - ciTBI Data, STAT 214, Spring 2025

2025-02-21

## 1. Introduction

Traumatic brain injury (TBI) is a leading cause of death and disability in children worldwide. In the United States, pediatric head trauma results in approximately 7,400 deaths and more than 600,000 emergency department visits each year[1]. Although computed tomography (CT) scans are the gold standard for diagnosing TBI, overuse of CT carries a risk of radiation-induced malignancies[1]. To reduce unnecessary radiation exposure, clinical decision rules (CDRs) have been developed to identify children at very low risk for clinically important TBI (ciTBI) who may not require CT imaging.

The PECARN TBI dataset, derived from Kuppermann et al. (2009), provides a large-scale study of pediatric head trauma that incorporates key clinical indicators for assessing the risk of ciTBI[1]. In this report, we perform an exploratory data analysis (EDA) on the dataset to understand its structure and assess data quality issues. We first describe the dataset and its variables and then identify and address inconsistencies such as missing values and measurement error. Finally, the goal is to refine data preprocessing techniques to improve clinical decision making in pediatric TBI screening.

## 2. Data

The dataset used in this analysis comes from the Pediatric Emergency Care Applied Research Network (PECARN) Traumatic Brain Injury (TBI) study, which aims to develop clinical decision rules (CDRs) for identifying children at low risk of clinically-important TBI (ciTBI). The dataset includes a variety of clinical and demographic features recorded for pediatric patients presenting with head trauma.

In this report, we focus on the following key variables relevant to TBI risk assessment and clinical decision-making: **GCS Score, Severity of injury mechanism, Presence of loss of consciousness, vomit, headache, Altered mental status (AMS), and Clinically-important traumatic brain injury (ciTBI).**

## 2.1 Data Collection

The data were collected as part of a multi-center prospective cohort study conducted across 25 emergency departments in North America. Pediatric patients under 18 years old who presented with head trauma were enrolled within 24 hours of injury.

Each patient underwent standardized clinical evaluation, and information was recorded on variables such as **Glasgow Coma Scale (GCS) scores, presence of amnesia, seizures, headache, and mechanism of injury.** CT scans were performed at the discretion of emergency department physicians, and follow-up data were collected to determine whether a patient had **ciTBI (defined as requiring neurosurgery, intubation for $> 24$ hours, hospital admission for $> 2$ nights, or resulting in death).**

## 2.2 Data Cleaning

The dataset contained several issues, including a high number of missing values, inconsistencies in variable naming, and unexpected discrepancies in certain categorical values. Notably, some variables had 92 values that did not align with their corresponding 0 values. For example, 'pos_ct' should have 92 values equal to the sum of 92 and 0 values in 'edct', but this was not the case.

1. **Standardizing Variable Names**

   - All variable names were converted to snake_case for consistency.
   - 'ha_verb' was renamed to 'ha' for easier processing.

2. **Handling Missing Values**

   - Drop all data has 'gcs_total'$< 14$
   - 'amnesia_verb' missing values were filled with 0 when 's_fx_bas', 'high_impact_inj_sev', and 'gcs_total' were 0, 0, 15, respectively.
   - 'loc_separate' missing values were set to 0 when 'act_norm', 'high_impact_inj_sev', and 'gcs_total' were 0, 0, 15 or 1, 1, 15.
   - 'seiz' missing values were set to 0 when 's_fx_bas', 'high_impact_inj_sev', 'gcs_total', and 'finding4' were 0, 1, 15, 0.

3. **Re-mapping Categorical Variables**

   - Based on the dataset documentation, values 0, 91, and 92 were reassigned according to their proper definitions.

4. **Further Missing Value Imputation**

   - 'amnesia_verb' missing values were filled with 1 when 's_fx_bas' == 1.
   - 'loc_separate' missing values were set to 1 when 'high_impact_inj_sev', 'finding4', and 's_fx_bas' were 3, 1, 1.

- 'seiz' missing values were set to 1 when 'loc_len', 'finding20', 's_fx_bas', and 'finding4' were 4, 1, 1, 1.
- 'ams' missing values were set to 1 when 'gcs_total' == 14.
- Variables with a small proportion of missing values were handled by removing rows containing missing data rather than imputing them.

5. **Assigning 'high_impact_inj_sev' Based on 'injury_mech'**

   - Following the dataset definition, 'high_impact_inj_sev' values were correctly assigned based on 'injury_mech' categories.
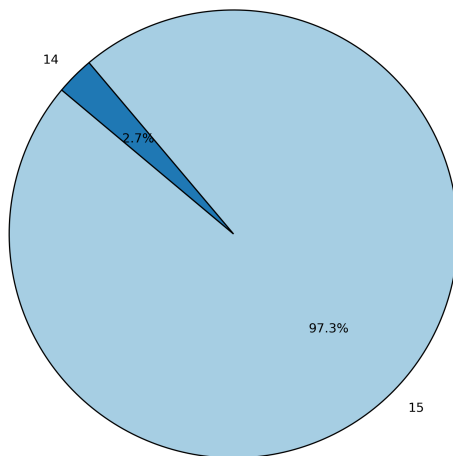
6. **Assigning 'pos_ct' Based on findings**

   - 'pos_ct' was updated based on whether any findings were present, ensuring consistency with the dataset documentation.
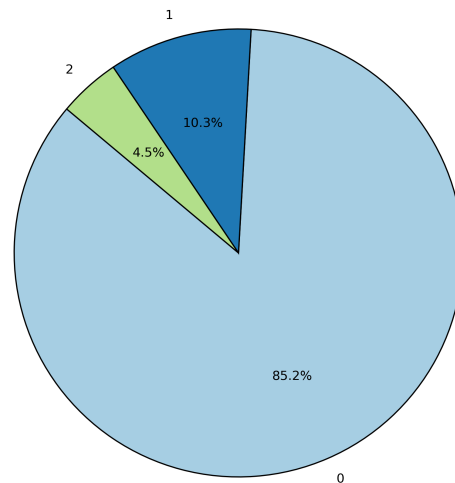
## 2.3 Data Exploration

These graphs provide an overview of the six key variables identified as most relevant in the dataset. **In the visualizations, 0 stands for no, 1 stands for yes, and 2 stands for unclear/suspected.**
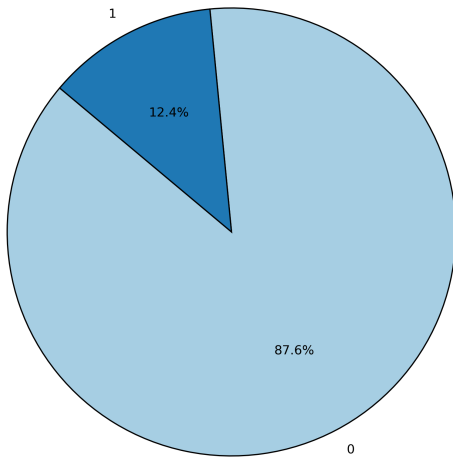


**GCS Scores**          **Loss of Consciousness**

**Pie Chart of ams**

1

12.4%

87.6%

0

**AMS**

**Pie Chart of ha**

1

30.4%

35.5%

34.0%

0

91

**Headache**

**Pie Chart of vomit**

1

13.2%

86.8%

0

**Vomit**

**Pie Chart of high_impact_inj_sev**

1

17.0%

40.5%

42.4%

2

3

**Severity of Injury Mechanism**

# 3. Findings

1. **First Findings**

   From the grouped bar chart below, AMS appears to be a significant factor in determining whether a CT scan is performed. Specifically, when AMS is 0, the majority of cases did not undergo a CT scan, whereas when AMS is 1, most cases received a CT scan. This suggests that AMS status plays a crucial role in clinical decision-making regarding CT utilization.

   **Proportion of CT Done by AMS Status (Normalized)**
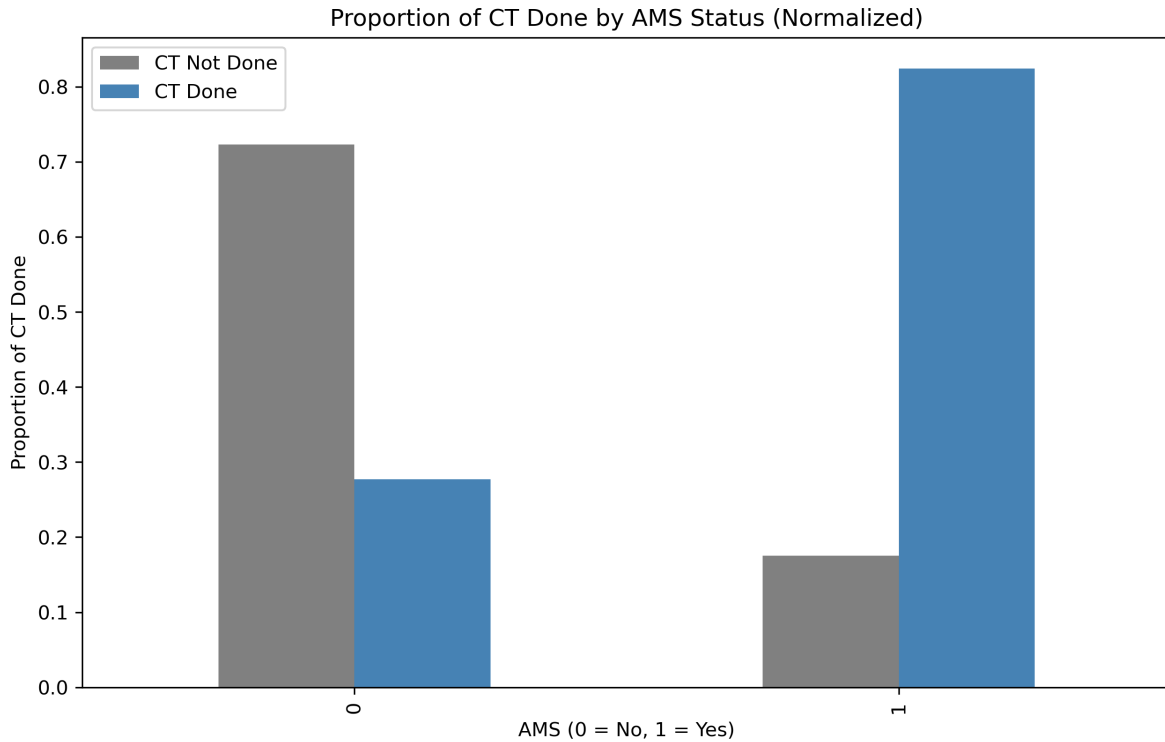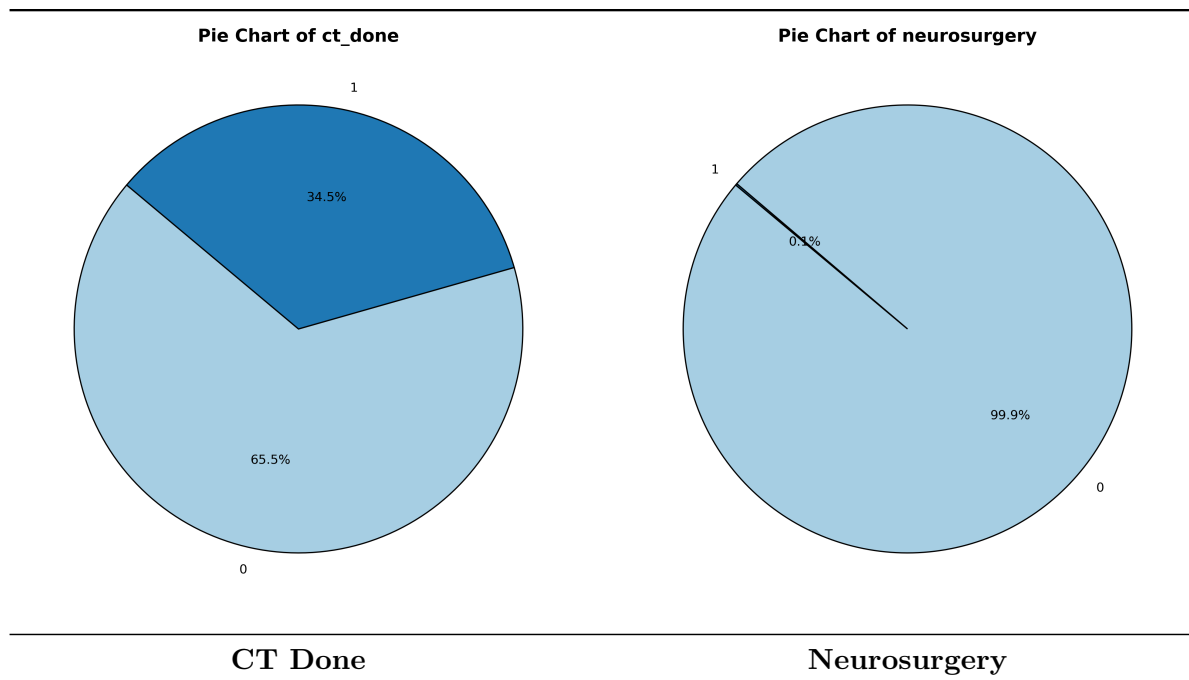
   Figure 1: Grouped Bar Chart of AMS and CT Done

2. **Second Findings**

   Although only 0.1% of patients required neurosurgical intervention, 34.5% underwent a CT scan. This discrepancy really suggests a huge overutilization of CT imaging, where a significant proportion of scans may have been performed on patients who ultimately did not require acute intervention.

**Pie Chart of ct_done**

1

34.5%

65.5%

0

**Pie Chart of neurosurgery**

1

0.1%

99.9%

0

| CT Done | Neurosurgery |
|---------|--------------|

3. **Third Findings**

The figure below indicates a strong relationship between GCS score and ciTBI, suggesting that lower GCS scores are associated with a higher likelihood of clinically important traumatic brain injury.
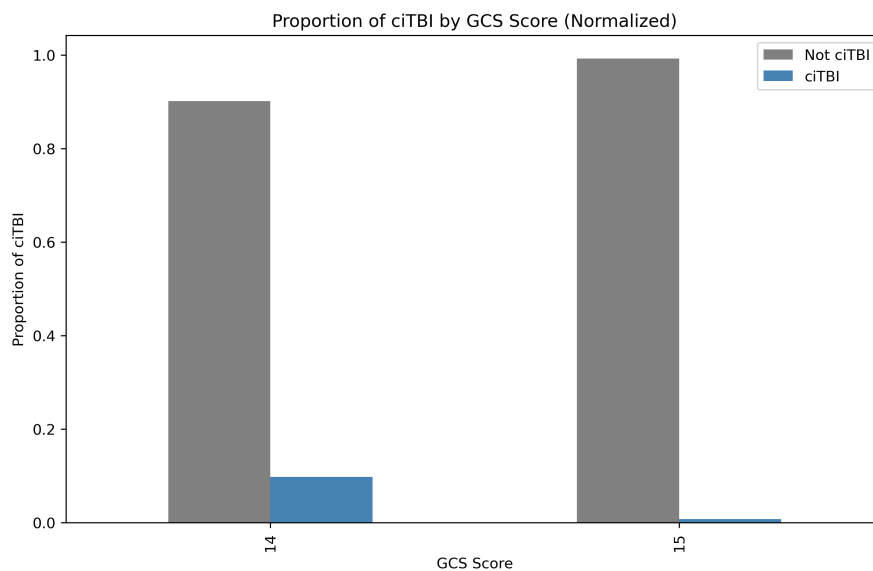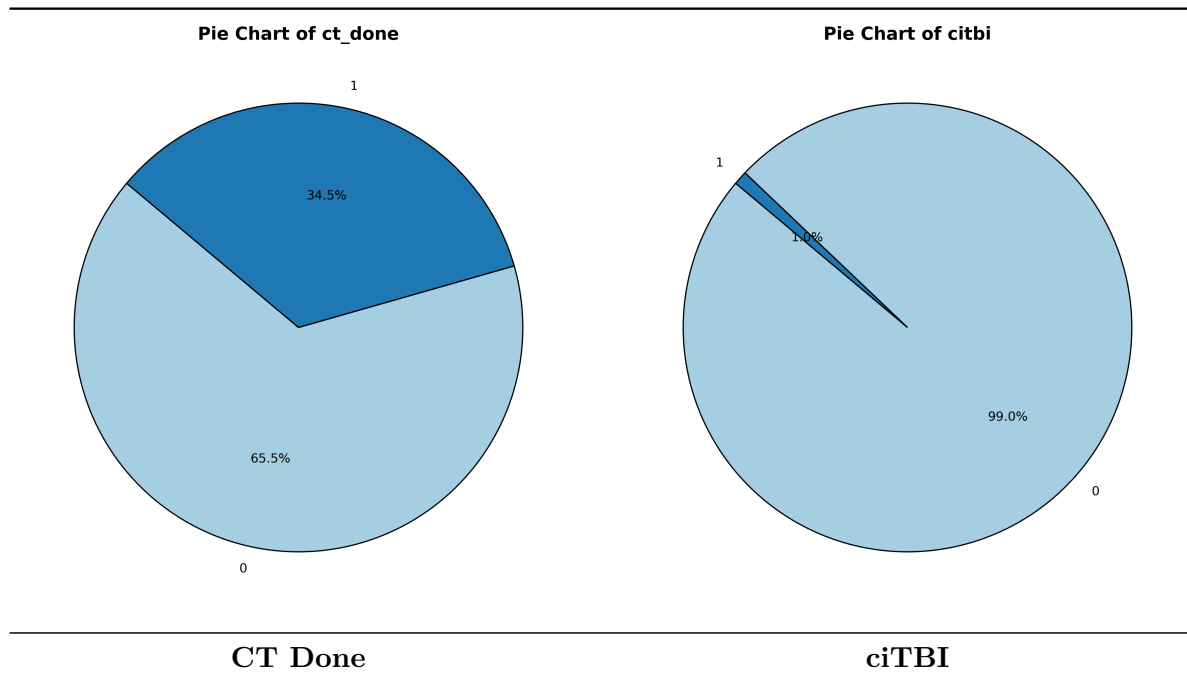
Proportion of ciTBI by GCS Score (Normalized)



Figure 2: Grouped Bar Chart of GCS Score and ciTBI

## 3.4 Reality Check

To validate our cleaned dataset, we compare two key variables—ciTBI occurrence rate and CT scan utilization rate—against the findings reported in Kuppermann et al. [1]. According to this study, clinically-important traumatic brain injuries (ciTBI) occurred in approximately 0.9% of cases, while CT scans were performed in 35.3% of patients. These figures serve as external benchmarks to assess the consistency and reliability of our dataset.

After computing the relative frequencies of ciTBI cases and CT scans performed, we compare them against the reference values. Our dataset shows a ciTBI rate close to 0.9% and a CT scan rate around 35.3%, showing that our data passes the reality check, confirming consistency with Kuppermann et al. [1].



| CT Done | ciTBI |

## 3.5 Stability Check

To assess the stability of our first finding, we compare results using two versions of the dataset. One where missing values were imputed based on domain knowledge, and another where they were imputed in the opposite manner (e.g., values originally filled with 1 were instead filled with 0). The left plot represents the dataset with domain-informed imputation, while the right plot represents the dataset where missing values were left unfilled.

As shown in the comparison, the two distributions exhibit high similarity, suggesting that our imputation strategy does not significantly alter the overall conclusions. This stability check

supports the reliability of our findings, indicating that our results are not overly sensitive to the choice of handling missing values.



**With Subjective Imputation**          **Without Subjective Imputation**

# 4. Modeling

## 4.1 Implementation

We implemented two Decision Tree Classifiers to predict whether a CT scan is necessary (ciTBI = 1). Decision trees were chosen for their interpretability and ability to handle both categorical and numerical features without scaling.

1. **Algorithms Used**

   - Model 1: Features: *gcs_total, loc_separate, ha, ams*
   - Model 2: Features: *gcs_total, loc_separate, vomit, ha_severity, high_impact_inj_sev*
   - Both models use a **Gini-impurity-based Decision Tree Classifier** with class weights to address data imbalance.

2. **Hyperparameter Selection**

   - *criterion*="gini": Selected for computational efficiency.
   - *class_weight* = {0:1, 1:500}: Increased weight for ciTBI=1 to improve recall.
   - *random_state* = 42: Ensured reproducibility.
   - *max_depth* and *min_samples_split*: Tuned using GridSearchCV with cross-validation (cv=5), optimizing for recall.

3. **Reproducibility**

   - Data was preprocessed by removing missing values and encoding categorical variables.

- Train-test split: 80% training, 20% testing (train_test_split).
- Decision trees were trained using scikit-learn's DecisionTreeClassifier.
- Models were evaluated using accuracy, recall, and confidence intervals.
- Decision paths leading to ciTBI=0 were extracted for interpretability.

## 4.2 Interpretability

Our decision tree models are highly interpretable, as they follow a step-by-step decision process based on clinical features. Each split represents a simple, logical rule that aligns with clinical reasoning, making the models easy to understand and apply.

**Model1:**



Figure 3: Model 1

In **Model 2**, several clinical features were transformed into **binary indicators** to simplify decision-making: **GCS Total Score(gcs_total):** Converted into 0(normal, GCS=15) and 1(impaired, GCS=14) to emphasize neurological deterioration. **Loss of Consciousness(loc_separate):** Transformed into 1(no loss of consciousness) and 0(loss of consciousness). **Vomiting(vomit):** Transformed inoto 1(no vmiting) and 0(vomiting). **Headache Severity(ha_severity):** Transformed inoto 1(mild and moderate) and 0(other) to better distinguish risk levels. **High-Impact Injury(high_impact_inj_sev):** Set to 1(low-impact injury) and 0(0ther).



Figure 4: Model 2

## 4.3 Stability

**Model 1(origin data set):**



Figure 5: Model 1(Original Sata Set)

**Model 1(Stability check data set):**



Figure 6: Model 1(Stability Check Data Set)

**Prediction Change Rate: 0.0078**

**Model 2(origin data set):**



Figure 7: Model 2(Original Data Set)

**Model 2(Stability check data set):**

Figure 8: Model 2(Stability Check Data Set)
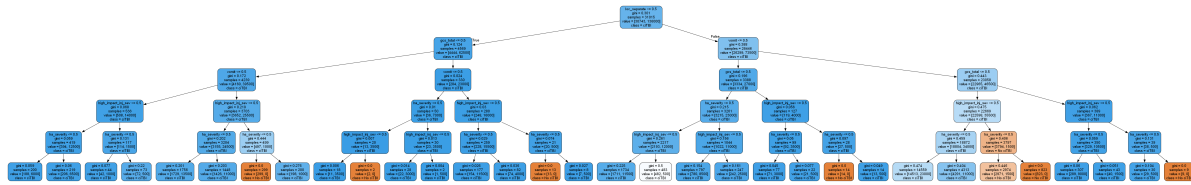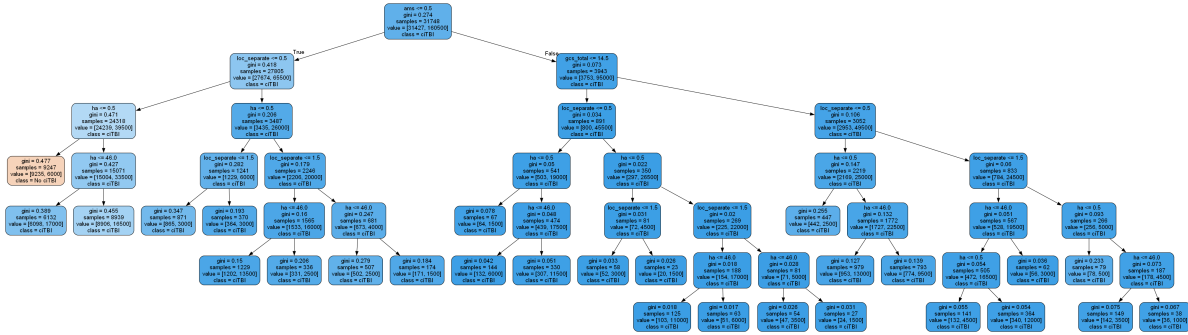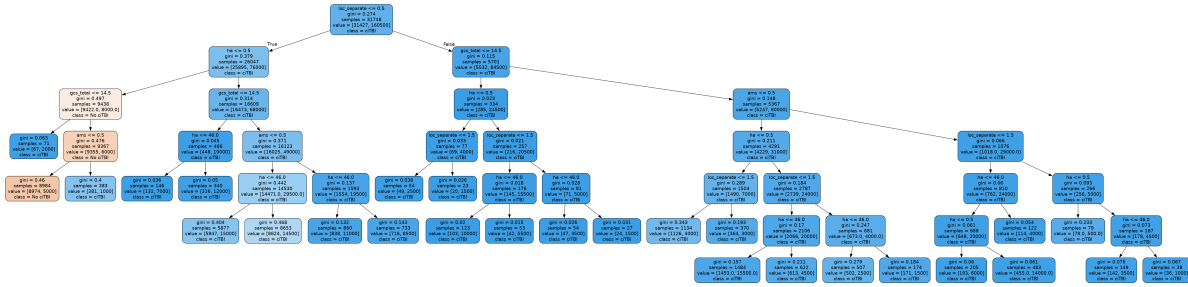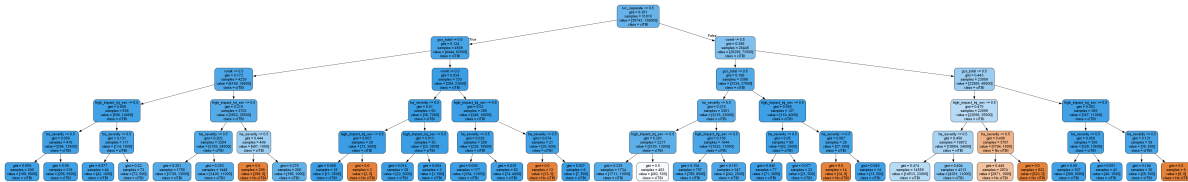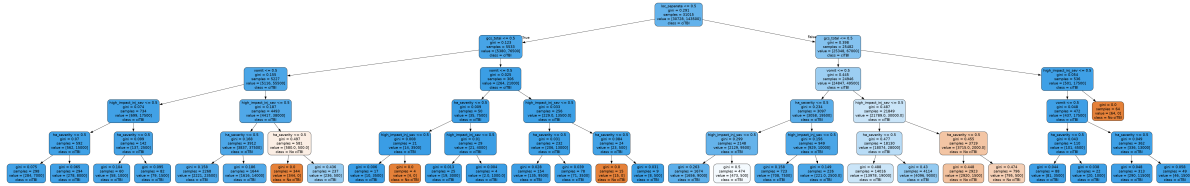
**Prediction Change Rate: 0.0006**

Based on the prediction change rate and the decision tree model showed above, we conclude that both models are stable.

# 5. Discussion

The dataset size did not impose major restrictions, but the limited number of ciTBI=1 cases posed challenges in training a balanced model. To address this, class weighting was applied to ensure the model did not overly favor ciTBI=0. Additionally, maintaining consistency between the original and perturbed datasets was essential for conducting a reliable stability analysis.

This study encompasses the three realms of data science: data and reality, algorithms and models, and future data and reality. While the dataset captures clinical characteristics, it does not perfectly reflect reality due to potential biases, measurement errors, and missing values. The decision tree models were chosen for their interpretability, and hyperparameter tuning, along with stability checks, ensured robustness. However, for these models to be effective in real-world settings, further validation on unseen patient data is necessary to assess their generalizability.

Although the dataset is derived from real clinical observations, it does not provide a one-to-one correspondence with reality. Certain factors influencing clinical decision-making may not be fully captured in structured data. Data visualization plays a crucial role in making findings more interpretable, but it can also oversimplify complex relationships in real-world clinical decision-making. Thus, while the models offer a structured approach to predicting ciTBI, clinical expertise remains essential in ensuring their practical applicability.

# 6. Conclusion

This study explored the PECARN TBI dataset to improve CT scan decision-making in pediatric head trauma. Decision tree models identified key predictors such as altered mental status, GCS score, and injury severity, confirming their clinical relevance. Stability checks

showed that the models remained robust under data perturbation. While the models demonstrated strong interpretability and reliability, further validation on external datasets is needed to ensure broader applicability. Ultimately, this work supports reducing unnecessary CT scans while maintaining high sensitivity for serious TBIs through data-driven decision support.

## 7. Academic honesty statement

I pledge that the work presented in this report is my own, and that all sources I have used, including discussions with classmates, have been properly cited. I understand that academic integrity is essential to maintaining the credibility and reliability of research.

Academic research honesty is necessary because it ensures the validity of knowledge, fosters trust within the academic community, and upholds ethical standards. Misrepresentation of work not only undermines the learning process but also distorts the integrity of scientific inquiry. By adhering to principles of honesty, we contribute to a research environment where ideas are built upon accurately and responsibly.

## 8. Collaborators

I done all the works by myslef in this lab.

## 9. Bibliography

1. Kuppermann, Nathan, et al. "Identification of Children at Very Low Risk of Clinically-Important Brain Injuries after Head Trauma: A Prospective Cohort Study." *The Lancet*, vol. 374, no. 9696, 2009, pp. 1160-1170. https://doi.org/10.1016/S0140-6736(09)61558-0.