

1 **Assessment and Prediction of Clinical Outcomes for ICU-Admitted Patients Diagnosed with** 2 **Hepatitis: Integrating Sociodemographic and Comorbidity Data.**

3 Dimple Sushma Alluri¹, and Felix M. Pabon-Rodriguez^{1,2*}

4 ¹Department of Biomedical Engineering and Informatics, Indiana University, Luddy School of Informatics,
5 Computing, and Engineering, Indianapolis, Indiana, United States

6 ²Department of Biostatistics and Health Data Science, Indiana University, School of Medicine,
7 Indianapolis, Indiana, United States

8 **Correspondence**

9 *Felix M. Pabon-Rodriguez

10 Email: fpabonrodriguez@gmail.com

11

12

13 **Abstract**

14 Hepatitis, a leading global health challenge, contributes to over 1.3 million deaths annually, with hepatitis
15 B and C accounting for the majority of these fatalities. Intensive care unit (ICU) management of patients is
16 particularly challenging due to the complex clinical care and resource demands. This study focuses on
17 predicting Length of Stay (LoS) and discharge outcomes for ICU-admitted hepatitis patients using
18 machine learning models. Despite advancements in ICU predictive analytics, limited research has
19 specifically addressed hepatitis patients, creating a gap in optimizing care for this population. Leveraging
20 data from the MIMIC-IV database, which includes around 94,500 ICU patient records, this study uses
21 sociodemographic details, clinical characteristics, and resource utilization metrics to develop predictive
22 models. Using Random Forest, Logistic Regression, Gradient Boosting Machines, and Generalized
23 Additive Model with Negative Binomial Regression, these models identified medications, procedures,
24 comorbidities, age, and race as key predictors. Total LoS emerged as a pivotal factor in predicting
25 discharge outcomes and location. These findings provide actionable insights to improve resource
26 allocation, enhance clinical decision-making, and inform future ICU management strategies for hepatitis
27 patients.

28 **KEYWORDS:** Hepatitis, predictive models, resource utilization, health disparity, SDOH.

29

30

31 **1. INTRODUCTION**

32 Hepatitis, an inflammatory liver disease, remains a significant global health challenge, claiming an
33 increasing number of lives each year. According to the World Health Organization (WHO) 2024 Global
34 Hepatitis Report, viral hepatitis is the second leading infectious cause of death globally, responsible for
35 1.3 million deaths annually. This number has risen from 1.1 million in 2019, with 83% of these deaths
36 attributed to hepatitis B and 17% to hepatitis C (WHO, 2024). Every day, approximately 3,500 people die
37 due to hepatitis B and C infections worldwide. Despite advancements in diagnostic tools and treatment
38 options, testing and treatment coverage rates have plateaued, signaling a growing public health crisis
39 (WHO, 2024). In the United States, the most common forms of viral hepatitis are hepatitis A, B, and C,
40 each impacting the liver differently and predominantly affecting distinct populations (CDC, 2024). Hepatitis
41 B and C pose severe health risks, often leading to chronic conditions such as cirrhosis and liver cancer.
42 These diseases are also the primary contributors to liver-related mortality globally (WHO, 2023).
43 Managing intensive care unit (ICU) - admitted patients with hepatitis is particularly challenging due to the
44 treatment complexities and the resource-intensive nature of care required for severe cases.

45

46 One critical issue in ICU settings is the unpredictability of a patient's length of stay (LoS), a key metric
47 influencing hospital resource management, patient care quality, and healthcare efficiency. Prolonged ICU
48 stays are closely linked to increased hospital costs and heightened resource strain (Peres et al., 2020).

49 Research also indicates that longer ICU stays correlate with increased long-term mortality rates,
50 underscoring the need for precise LoS predictions (Moitra et al., 2015). Beyond LoS, discharge outcomes
51 and locations serve as vital metrics for assessing patient recovery, readmission risks, and the burden on
52 healthcare systems (Hickman, 2018). Sociodemographic and clinical factors—including race, gender,
53 marital status, insurance type, age, and type of hepatitis—significantly influence hospitalization outcomes,
54 particularly for ICU-admitted hepatitis patients. These factors also play a crucial role in determining
55 hospital stay duration and discharge outcomes (Arnab Kumar Ghosh et al., 2021; Hayes et al., 2016;
56 Dubin et al., 2024; Eskandari et al., 2022; Ng et al., 2024).

57 Recent years have seen the growing application of machine learning (ML) in healthcare, particularly for
58 predictive modeling. ML algorithms provide several advantages over traditional statistical approaches,
59 including their ability to analyze large, complex datasets and identify subtle patterns that conventional
60 techniques may overlook (An et al., 2023). Predictive models using ML have been successfully applied in
61 ICU settings to improve forecasts related to patient outcomes, treatment responses, and resource needs
62 (Choi et al., 2022; Levin et al., 2020). However, there is a noticeable lack of ML models tailored to
63 hepatitis patients and their unique clinical profiles.

64 The objective of this research is to explore and develop ML models to predict the length of stay, discharge
65 location, and discharge outcomes for ICU-admitted hepatitis patients. These models will incorporate
66 sociodemographic and clinical variables from ICU admission records as key predictors. By addressing
67 this gap, the study aims to provide actionable insights for healthcare providers, enhancing resource
68 planning and improving hepatitis patient outcomes in ICU settings.

69 **2. MATERIALS AND METHODS**

70 **2.1. Data**

71 The data set used for this study was sourced from the Medical Information Mart for Intensive Care
72 (MIMIC-IV), version 3.0, a comprehensive and deidentified repository of patient records (Johnson et al.,
73 2024; Johnson et al., 2023; Goldberger et al., 2000). MIMIC-IV comprises data collected from patients
74 admitted to the emergency department (ED) or intensive care units at the Beth Israel Deaconess Medical
75 Center (BIDMC) in Boston, MA. It includes over 364,000 unique patient records, with a total of 546,000
76 hospital admissions and nearly 94,500 ICU stays, providing a rich source of information for developing
77 predictive models in healthcare. MIMIC-IV is organized into two main data modules: *hosp* and *icu*. The
78 *hosp* module captures information from the hospital-wide electronic health record (EHR), detailing
79 hospitalizations, patient demographics, laboratory results, medication administration, billing data, and
80 more. In contrast, the *icu* module contains highly detailed clinical data from the ICU, sourced from the
81 MetaVision clinical information system, including treatment plans, and monitoring data (Johnson et al.,
82 2024; Johnson et al., 2023; Goldberger et al., 2000).

83
84 Data collection spans from 2008 to 2022, and the records are meticulously de-identified following Health
85 Insurance Portability and Accountability Act (HIPAA) guidelines. Patient identifiers, such as names and
86 social security numbers were removed to protect privacy. Additionally, dates of death are available up to
87 one-year post-discharge and are derived from hospital and state records. Ethical approval for data
88 collection and the creation of this research resource was obtained from the Institutional Review Board at
89 the BIDMC. The dataset was made accessible through a data use agreement, ensuring compliance with
90 all ethical and legal guidelines for data usage. With this agreement in place, the database is available for
91 researchers and academics.

92 **2.2. Variables**

93 The independent variables for this study include socio-demographic, clinical, and resource utilization
94 data. Socio-demographic variables cover gender, age, race, marital status, and type of insurance.
95 Admission characteristics include the type of admission and admission location. Clinical variables
96 included the number of comorbidities, as well as indicators for Hepatitis A, B, C, D, and E, and Hepatic
97

98 Coma. Resource utilization variables include counts of medications, procedures, and drugs dispensed
99 before ICU admission. Additional ICU-specific measures capture the volume of fluids prescribed and the
100 number of procedures performed within the ICU. This study aims to predict three primary outcomes: *Total*
101 *Length of Stay*, measured as the days spent in the ICU or hospital; *Discharge Location*, detailing the
102 destination upon discharge, such as home, hospice, or skilled nursing; and *Discharge Outcome*, a binary
103 variable indicating discharge status (death or alive).

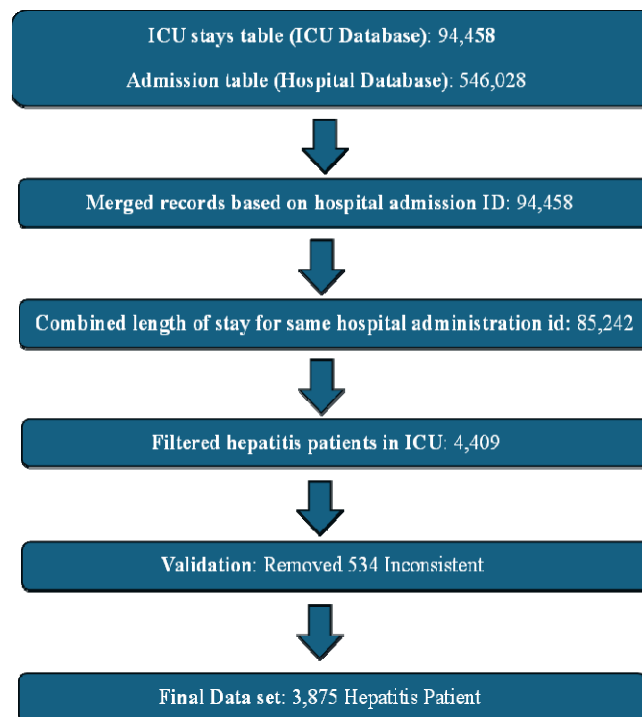
104

105 2.3. Data Processing and Preparation

106 Given the extensive and diverse nature of the MIMIC-IV dataset, specific steps were undertaken using
107 MySQL to refine and tailor the data for this study on hepatitis patients admitted to the ICU (Oracle, 2021).
108 Figure 1 illustrates the data pre-processing workflow, detailing the sequential steps taken to create the
109 final dataset.

110

111



112

113

114 **Figure 1.** Data preprocessing workflow for extracting and refining hepatitis patient records from the
115 MIMIC-IV database.

116

117 a) **International Classification of Diseases (ICD) Code Filtering and Data Integration:** To isolate
118 hepatitis-related hospital admissions, ICD codes specific to viral hepatitis were identified and
119 filtered from a comprehensive table of diagnostic codes. These filtered ICD codes were then
120 linked with a table containing hospital admission identifiers to extract records corresponding to
121 hepatitis diagnoses. Subsequently, the table was merged with another table containing patient
122 demographic information and total length of stay, resulting in a hepatitis-specific patient cohort.

123 b) **Feature Engineering:** To enrich the dataset, several additional variables were calculated. Pre-
124 ICU and post-ICU medication counts, as well as procedure counts, were derived by aggregating
125 respective records. For patients with multiple hepatitis diagnoses, individual indicator columns for
126 Hepatitis A, B, C, D, and E were created, facilitating detailed analysis of co-infections. A
127 comorbidity count was also computed by excluding hepatitis-related ICD codes and counting the
128 remaining diagnoses.

- 129 c) **Data Cleaning and Validation:** The combined dataset initially consisted of 4,409 patient records.
130 A rigorous data validation process identified 534 records with inconsistencies between discharge
131 outcomes and discharge locations. These records were excluded to ensure data integrity,
132 resulting in a final dataset of 3,875 hepatitis patient records. Patients included in this analysis
133 were those admitted to the ICU at Beth Israel Deaconess Medical Center and diagnosed with any
134 form of hepatitis.
- 135 d) **Data Transformation:** The data transformation phase involved preparing categorical variables for
136 ML by converting them into factor data types and recategorizing the race variable to address
137 sparsity and improve interpretability. Subcategories with similar characteristics were grouped
138 under broader categories, such as combining all Asian subcategories (e.g., "Asian - Chinese,"
139 "Asian - Indian") into "Asian" and all Black subcategories (e.g., "Black/African," "Black/Caribbean
140 Island") into "Black." Hispanic/Latino and White subcategories were similarly merged, while
141 categories with low counts, such as "South American," "Native Hawaiian or Other Pacific
142 Islander," and "Multiple Race/Ethnicity," were consolidated into an "Other" category.
- 143 e) **Handling Missing Values:** To address missing values, a thorough check was conducted across
144 all variables, revealing the following counts of missing values per variable: marital status (296),
145 insurance type (64), number of procedures in the ICU (237), number of pre-ICU procedures
146 (642), and medications (248), total length of stay (1), and discharge location (18). Multiple
147 Imputation by Chained Equations (MICE) was employed to impute the missing values (Azur et al.,
148 2011). Predictive Mean Matching (PMM) was used for numeric columns, while Polytomous
149 Logistic Regression (Polyreg) was applied to categorical variables with more than two levels. The
150 imputation process was repeated five times, and the results were pooled to produce a single,
151 complete dataset for analysis.
- 152 f) **Handling Imbalanced Data:** The discharge outcome variable, which indicates whether patients
153 were alive or deceased upon discharge, displayed a significant class imbalance, with 3,353
154 records for the alive category and 522 for death. To address this, the Synthetic Minority
155 Oversampling Technique (SMOTE) was used to balance the classes (Pradipta et al., 2021). This
156 resulted in a revised dataset with 2,088 records labeled as alive and 1,566 as death, ensuring a
157 more equitable representation of both classes for model training.
- 158 g) **Training and Test Data:** After data cleaning and preprocessing, the records were partitioned for
159 subsequent modeling tasks. An 80:20 split was used to separate the data into training and testing
160 sets, a crucial step for evaluating model performance.
- 161 I. For the *Discharge Outcome*, after applying SMOTE to address class imbalance, a total of
162 3,654 records were obtained. Of these, 2,923 records (80%) were used for training, and
163 731 records (20%) were allocated for testing. The training set was used to train the
164 models, while the testing set was reserved to evaluate predictive performance.
- 165 II. For *Length of Stay* and *Discharge Location*, 3,875 records were available. The same
166 partitioning strategy was applied, with 3,100 records (80%) used for training and 775
167 records (20%) for testing. This partition ensured that models could be appropriately
168 trained and assessed for their predictive ability across different outcomes.

170 2.4. Machine Learning Models

171 For predicting the discharge outcome (death or alive), the discharge location variable was not used as
172 predictor for the analysis. This decision was made because certain discharge location categories
173 overlapped with death records, providing redundant information. The following models were considered:

- 174 a) **Model 1: Logistic Regression:** It is a statistical model used for binary classification, estimating
175 the probability of one of two possible outcomes (in this case alive or death) based on a set of
176 predictor variables. This model uses the logistic function to transform predictions into probabilities
177 between 0 and 1 (Stoltzfus, 2011). Logistic Regression was chosen for its interpretability,
178 providing clear insights into how each clinical and socio-demographic feature influences
179 discharge outcome. This characteristic is especially valuable in this research, where

180 understanding the relative impact of individual predictors on patient survival is crucial for making
181 informed inferences about patient care.

182 b) **Model 2: Random Forest Classification:** This model was applied to enhance the accuracy of
183 predicting discharge outcome. The model constructs multiple decision trees using random data
184 samples, with each tree voting for a class. The final prediction is determined by the majority vote
185 across all trees. This ensemble approach captures complex, non-linear relationships among
186 variables, offering strong predictive performance and minimizing the risk of overfitting (Wrld,
187 2024). Although it does not offer the interpretive insights of Logistic Regression, Random Forest's
188 capability to handle heterogeneous data makes it well-suited for maximizing predictive power in
189 clinical scenarios (Couronné et al. 2018).

190 To predict LoS for ICU-admitted hepatitis patients, we considered two scenarios. In the first case, we
191 considered LoS as a count variable (Fernandez & Vatcheva, 2022), and as a continuous variable (Xu et
192 al., 2022) in our second case, allowing for a comprehensive exploration of different modeling approaches.

193 a) **Model 1: Generalized Additive Model (GAM) with Negative Binomial (NB) Distribution:** For
194 the discrete case, a GAM with a NB distribution was utilized to address overdispersion, where the
195 variance in the count data exceeds the mean (Zeileis et al., 2015). By incorporating smooth
196 functions, the GAM framework allowed for modeling potential non-linear relationships between
197 predictors and LoS. This approach provided the necessary flexibility to capture complex patterns
198 in the data, ensuring a robust and accurate analysis (Hastie & Tibshirani, 2014).

199 b) **Model 2: Random Forest Regression:** For the continuous case, Random Forest Regression
200 was employed. This ensemble method constructs multiple decision trees and averages their
201 predictions to capture non-linear relationships and interactions among predictors (Wrld, 2024).
202 Using Random Forest provided an alternative view, enabling a performance comparison between
203 the count-based and continuous modeling approaches.

204 To predict the discharge location of ICU-admitted hepatitis patients, models were trained after excluding
205 discharge outcome variable which is not relevant to this analysis.

206 a) **Model 1: Gradient Boosting Model (GBM):** Gradient Boosting is an ensemble learning
207 technique that builds multiple decision trees sequentially, with each new tree correcting errors
208 made by the previous ones. This model focuses on reducing bias and improving predictive
209 performance through boosting. GBM was chosen for its strong ability to handle complex, non-
210 linear relationships and provide high accuracy (Zhang et al. 2019), which is critical in predicting
211 discharge locations, a multi-class categorical outcome.

212 b) **Model 2: Multinomial Logistic Regression:** Multinomial Logistic Regression is an extension of
213 logistic regression used when the dependent variable has more than two categories (Kwak &
214 Clayton-Matthews, 2002). It models the probability of each class as a function of the predictors,
215 providing a straightforward and interpretable approach to multi-class classification. This model
216 was included to establish baseline and offer a simpler, interpretable model for understanding the
217 impact of predictor variables on different discharge locations.

218 219 **2.5. Model Evaluation**

220 The models were evaluated using appropriate performance metrics based on the type of outcome
221 variable:

222 • **Regression Models:** For LoS, regression models were evaluated using Root Mean Squared
223 Error (RMSE), R-squared (R^2), and Mean Absolute Error (MAE). RMSE measures the average
224 magnitude of error between predicted and actual values, penalizing larger deviations more
225 heavily, making it particularly sensitive to outliers. R^2 quantifies the proportion of variation in the
226 dependent variable explained by the model, offering an overall measure of fit. MAE, unlike

- 227 RMSE, calculates the average of the absolute errors, providing an intuitive measure of model
228 accuracy by treating all errors equally (Zwanenburg, 2022).
- 229 • **Multiclass Classification Models:** These were assessed using Accuracy, the Kappa statistic,
230 the Brier Score, and the Area Under the Receiver Operating Characteristic Curve (ROC AUC).
231 Accuracy reflects the percentage of correctly classified instances, while the Kappa statistic
232 accounts for agreement by chance, providing a more robust evaluation of model performance.
233 The Brier Score measures the mean squared difference between predicted probabilities and the
234 actual class, with lower scores indicating better calibration. ROC AUC quantifies a model's ability
235 to distinguish between classes, with values closer to 1 denoting superior performance
236 (Zwanenburg, 2022).
 - 237 • **Binary Classification Models:** For binary outcomes, Logistic Regression and Random Forest
238 Classification models were evaluated using Brier Score, Accuracy, Kappa, ROC AUC, Sensitivity,
239 and Specificity. Sensitivity, also known as recall, indicates the proportion of true positives
240 correctly identified by the model, while Specificity measures the proportion of true negatives
241 correctly identified. These metrics provide insights into the model's ability to balance false
242 positives and false negatives. The Brier Score evaluates the accuracy of probabilistic predictions,
243 with lower values reflecting better performance. ROC AUC measures the discriminatory power of
244 the model across thresholds, while Accuracy and Kappa offer overall measures of correctness,
245 with Kappa adjusting for chance agreement (Zwanenburg, 2022).

246 2.6. Resampling Techniques

247 To ensure the robustness and generalizability of the models, 10-fold cross-validation was employed. This
248 approach divides the dataset into 10 subsets (or folds). The model is trained on 9 folds and validated on
249 the remaining fold in each iteration. This process is repeated 10 times, with each fold serving as the
250 validation set once, and the average performance across all iterations is reported. The choice of K=10 is
251 supported by experimental analysis and observations, as it provides a balance between computational
252 efficiency and reliability (Verma et al., 2024). K-fold cross-validation is critical for mitigating the risk of
253 overfitting, as it ensures that the model is tested on diverse subsets of the data, not just a single train-test
254 split. This approach provides a more reliable estimate of the model's performance on unseen data by
255 leveraging all data points for both training and validation at different stages. Particularly for complex
256 datasets, such as the one used in this study, 10-fold cross-validation enhances the reliability and
257 generalizability of the results, making it a robust validation strategy for machine learning applications
258 (Wilimitis & Walsh, 2023).

260 2.7. Software

261 Data analysis was conducted in RStudio (R Core Team, 2024; RStudio Team, 2022) using various
262 packages to streamline data manipulation, model development, and visualization. MySQL (Oracle, 2024)
263 was used to extract and manage data from the MIMIC-IV database. In RStudio, data processing was
264 performed using the *tidyverse* (Wickham et al., 2019), while missing data were imputed with *mice* (Buuren
265 & Groothuis-Oudshoorn, 2011), and class imbalance was addressed using *DMwR* (Torgo, 2024). Model
266 development and evaluation were facilitated by *MachineShop* (Smith, 2024), enabling a streamlined
267 approach to implementing and comparing various predictive models. Statistical modeling was conducted
268 using *MASS* (Venables & Ripley, 2002), which supported the fitting of Negative Binomial regression
269 models, and *mgcv* (Wood, 2011), which was utilized for constructing Generalized Additive Models (GAMs)
270 to capture nonlinear relationships. Random Forest modeling was implemented using the *randomForest*
271 package (Liaw & Wiener, 2002). Data visualizations were created with *ggplot2* (Wickham, 2016), while
272 descriptive statistics and baseline characteristics were summarized using *tableone* (Yoshida & Bartel,
273 2022).

274
275
276
277
278

279 **3. RESULTS**

280 **3.1. Demographic and Descriptive Statistics**

281 The study included a total of 3,875 patients admitted to the ICU with hepatitis-related conditions. The
 282 gender distribution was predominantly male, accounting for 68.2% (n = 2,644) of the sample, while
 283 females constituted 31.8% (n = 1,231). The average age of the participants was 53.26 years (SD =
 284 12.77). The racial composition was primarily White (57.3%, n = 2,220), followed by Black (16.9%, n =
 285 653), Hispanic/Latino (6.3%, n = 245), and Asian (6.0%, n = 234). Other racial groups, such as
 286 Indian/Alaska Native and Portuguese, made up a smaller proportion of the population, each representing
 287 less than 1% of the sample. These demographic characteristics, along with details on marital status,
 288 insurance type, clinical and admission characteristics, are summarized in **Table 1**. The distribution of
 289 hepatitis types is shown in **Figure 2**, which reveals that Hepatitis C was the most prevalent condition,
 290 affecting 82.2% (n = 3,186) of the total sample. Hepatitis B was the second most common condition at
 291 19.3% (n = 746), while Hepatitis A (1.0%, n = 37), Hepatitis D (9.2%, n = 357), and Hepatitis E (0.1%, n =
 292 3) were less frequent. Additionally, 8.0% (n = 311) of patients experienced hepatic coma.

293 *Table 1. Baseline demographic, clinical, and admission characteristics of study participants*

| | |
|------------------------------|---------------|
| N | 3,875 |
| GENDER = Male (%) | 2644 (68.2) |
| AGE (MEAN (SD)) | 53.26 (12.77) |
| RACE (%) | |
| INDIAN/ALASKA NATIVE | 10 (0.3) |
| ASIAN | 234 (6.0) |
| BLACK | 653 (16.9) |
| HISPANIC/LATINO | 245 (6.3) |
| OTHER | 120 (3.1) |
| PORTUGUESE | 29 (0.7) |
| UNKNOWN | 364 (9.4) |
| WHITE | 2220 (57.3) |
| MARITAL STATUS (%) | |
| DIVORCED | 398 (10.3) |
| MARRIED | 1106 (28.5) |
| NA | 296 (7.6) |
| SINGLE | 1873 (48.3) |
| WIDOWED | 202 (5.2) |
| INSURANCE (%) | |
| MEDICAID | 1542 (39.8) |
| MEDICARE | 1472 (38.0) |
| NA | 64 (1.7) |
| OTHER | 123 (3.2) |
| PRIVATE | 674 (17.4) |
| TYPE OF ADMISSION (%) | |
| DIRECT EMERGENCY | 230 (5.9) |
| DIRECT OBSERVATION | 10 (0.3) |
| ELECTIVE | 50 (1.3) |

| | |
|---|-----------------------|
| EU OBSERVATION | 5 (0.1) |
| EMERGENCY WARD | 2296 (59.3) |
| OBSERVATION ADMIT | 548 (14.1) |
| SURGICAL SAME DAY ADMISSION | 155 (4.0) |
| URGENT | 581 (15.0) |
| ADMISSION LOCATION (%) | |
| EMERGENCY ADMISSIONS | 1998 (51.6) |
| INTERNAL TRANSFERS | 45 (1.2) |
| REFERRALS | 927 (23.9) |
| TRANSFER FROM OTHER FACILITY | 898 (23.2) |
| UNKNOWN | 7 (0.2) |
| No. of COMORBIDITIES (MEDIAN [IQR]) | 18.00 [12.00, 24.00] |
| HEPATITIS_A = 1 (%) | 37 (1.0) |
| HEPATITIS_B = 1 (%) | 746 (19.3) |
| HEPATITIS_C = 1 (%) | 3186 (82.2) |
| HEPATITIS_D = 1 (%) | 357 (9.2) |
| HEPATITIS_E = 1 (%) | 3 (0.1) |
| HEPATIC_COMA = 1 (%) | 311 (8.0) |
| PRE-ICU MEDICATIONS COUNT (MEDIAN [IQR]) | 82.00 [36.00, 140.25] |
| PRE-ICU PROCEDURES COUNT (MEDIAN [IQR]) | 3.00 [2.00, 6.00] |
| DISPENSED MEDICATIONS COUNT (MEDIAN [IQR]) | 65.50 [33.00, 118.50] |
| ICU MEDICATIONS COUNT (MEDIAN [IQR]) | 54.00 [20.00, 154.00] |
| ICU PROCEDURES COUNT (MEDIAN [IQR]) | 7.00 [3.00, 14.00] |
| TOTAL LENGTH OF STAY (MEDIAN [IQR]) | 2.00 [1.00, 4.00] |
| DISCHARGE_LOCATION (%) | |
| AGAINST MEDICAL ADVICE | 170 (4.4) |
| DEATH | 522 (13.5) |
| HOME-BASED CARE | 1997 (51.5) |
| HOSPICE CARE | 100 (2.6) |
| NA | 18 (0.5) |
| OTHER HEALTHCARE FACILITY | 110 (2.8) |
| PSYCHIATRIC CARE | 83 (2.1) |
| REHABILITATION | 191 (4.9) |
| SKILLED NURSING/LONG-TERM CARE | 684 (17.7) |
| DISCHARGE OUTCOME = 1 (%) | 522 (13.5) |

294
 295 The distribution of hepatitis types across racial groups is depicted in **Figure 3**. Hepatitis C was the most
 296 prevalent condition across all racial groups, with White patients contributing the largest proportion. Among
 297 Black and Hispanic/Latino patients, Hepatitis C was also the most common condition but showed slightly
 298 smaller proportions compared to White patients, while Hepatitis B was notably more prevalent among
 299 Asian patients compared to other racial groups. In addition, Hepatitis B disproportionately affects non-

300 White individuals compared to Whites. Hepatitis A, D, and E were observed in smaller proportions across
 301 all racial groups, with Hepatitis E being particularly rare. Hepatic Coma was most common among White
 302 patients, with smaller contributions from other racial groups.

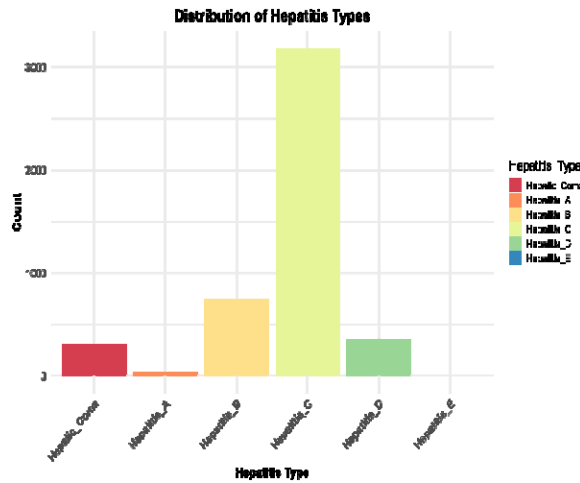


Figure 2. Distribution of hepatitis types among study participants

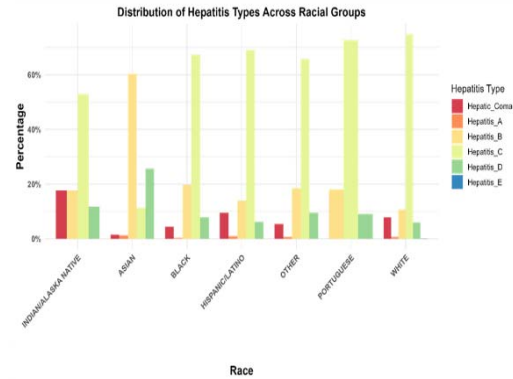


Figure 3. Distribution of hepatitis across racial groups

303 The distribution of total length of stay for hepatitis patients was highly skewed as commonly seen in
 304 hospital records, with most patients having a short hospital stay (**Figure 4**). The median total length of
 305 stay was 2 days (IQR: 1–4 days) as stated in **Table 1**, with some patients having an extended stay of over
 306 30 days. Variations in length of stay were also observed based on admission location (**Figure 5**). For
 307 instance, patients admitted from emergency admissions had shorter median stays compared to those
 308 transferred from other facilities. Notably, patients transferred from other facilities exhibited a wider range
 309 of lengths of stay, including outliers with significantly extended hospitalizations. These trends highlight
 310 how admission pathways may influence the duration of ICU stays.

311

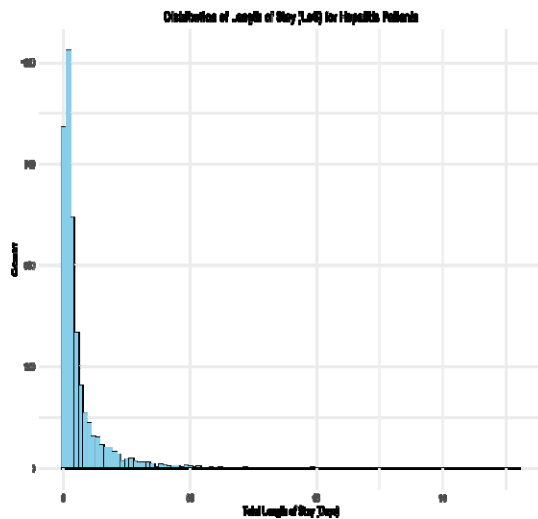


Figure 4. Distribution of total LoS for hepatitis patients

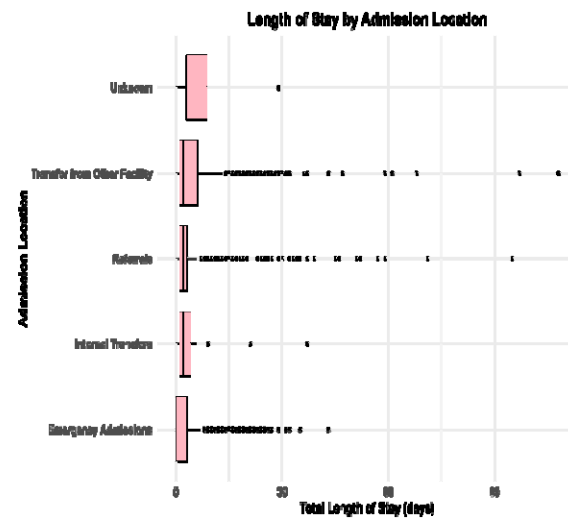


Figure 5. Length of stay by admission location

312 The boxplot in **Figure 6** illustrates the relationship between the total length of stay and discharge
313 outcome. Patients with longer lengths of stay were more likely to have died, as indicated by the higher
314 median and wider range of LoS in the "Died" group compared to the "Survived" group. **Figure 7** highlights
315 the distribution of discharge locations among patients. The majority (51.5%, n = 1,997) were discharged
316 to home-based care. Skilled nursing or long-term care facilities accounted for 17.7% (n = 684) of
317 discharges, while 13.5% (n = 522) of patients died. Other notable discharge destinations included
318 rehabilitation centers (4.9%, n = 191), hospice care (2.6%, n = 100), psychiatric care (2.1%, n = 83), and
319 other healthcare facilities (2.8%, n = 110). A small proportion of patients left the hospital against medical
320 advice (4.4%, n = 170).
321
322
323

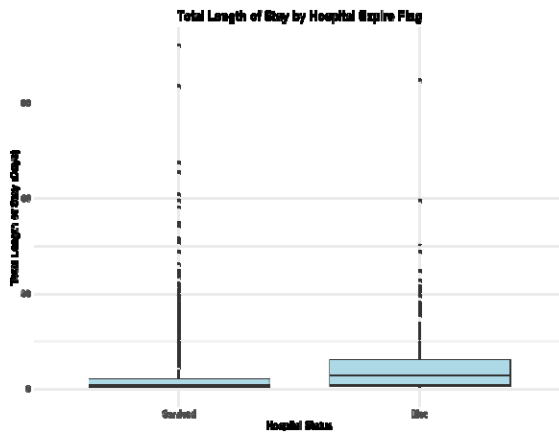


Figure 6. Discharge outcome based on total LoS

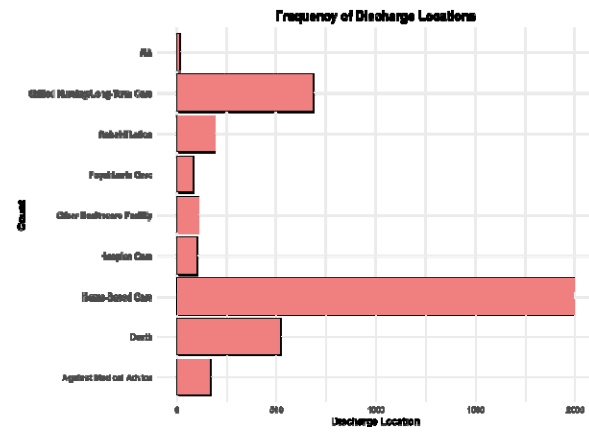


Figure 7. Distribution of discharge locations among patients

324
325 **3.2. Model Results and Performance**
326 Discharge Outcome (Died/Alive)

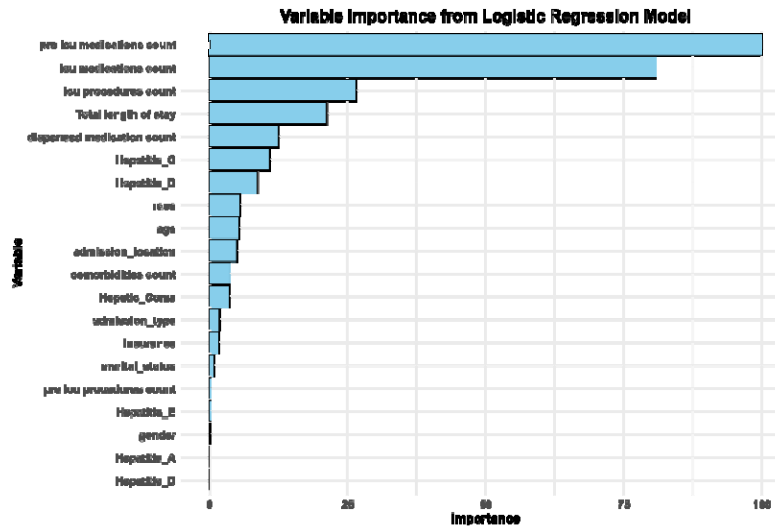
327 The Logistic Regression and Random Forest models were developed and evaluated to predict discharge
328 outcomes (died/alive) for ICU-admitted hepatitis patients. These models were tested on test data and
329 through cross-validation. Performance metrics for both models are presented in **Table 2**. On the test data,
330 Logistic Regression yielded a Brier Score of 0.1302, an accuracy of 0.818, a Kappa of 0.622, and an
331 ROC AUC of 0.887. Sensitivity and specificity were 0.73 and 0.87, respectively. The Random Forest
332 model demonstrated a Brier Score of 0.0875, an accuracy of 0.87, a Kappa of 0.742, and an ROC AUC of
333 0.95. Sensitivity and specificity for the Random Forest model were 0.78 and 0.94, respectively. Cross-
334 validation results are also shown in **Table 2**, indicating the mean and standard deviation for each metric
335 across repeated samples. The mean Brier Score for Random Forest was 0.089 (± 0.010), while Logistic
336 Regression had a mean Brier Score of 0.129 (± 0.011). The Random Forest model had a mean accuracy
337 of 0.87 (± 0.017) and a mean Kappa of 0.729 (± 0.037), compared to Logistic Regression's mean accuracy
338 of 0.82 (± 0.02) and mean Kappa of 0.629 (± 0.043). Moreover, statistical comparisons were conducted to
339 evaluate differences in performance metrics between the Logistic Regression and Random Forest
340 models. The results indicated statistically significant differences ($p < 0.001$) for all key metrics, including
341 Brier Score, accuracy, Kappa, ROC AUC, sensitivity, and specificity. These findings demonstrated that the
342 Random Forest model consistently outperformed Logistic Regression across all evaluated criteria.
343 Figures 8 and 9 display the variable importance for the Logistic Regression and Random Forest models,
344 respectively. The top predictors in the Logistic Regression model (**Figure 8**) were pre-ICU medications
345 count, ICU medications count, and ICU procedures count, with sociodemographic factors such as race

346 and age also ranking among the top 10 variables. For the Random Forest model (**Figure 9**), ICU
 347 medications count, ICU procedures count, and total LoS were identified as the most influential variables,
 348 with race and age again appearing among the top 10 predictors. These findings highlight the consistent
 349 importance of race and age as predictors, emphasizing the need to consider both clinical and
 350 sociodemographic factors when analyzing discharge outcomes.

351 *Table 2. Performance metrics of Logistic Regression and Random Forest models on test and cross-validated data*

| Metrics | Logistic Regression | | Random Forest | |
|-------------|---------------------|------------------|---------------|------------------|
| | Test data | Cross-validation | Test data | Cross-validation |
| Brier Score | 0.1301607 | 0.129 ± 0.011 | 0.08772361 | 0.089 ± 0.010 |
| Accuracy | 0.8166895 | 0.821 ± 0.020 | 0.87551300 | 0.870 ± 0.017 |
| Kappa | 0.6220590 | 0.629 ± 0.043 | 0.74208369 | 0.729 ± 0.037 |
| ROC AUC | 0.8865039 | 0.889 ± 0.017 | 0.94999924 | 0.948 ± 0.012 |
| Sensitivity | 0.7350158 | 0.736 ± 0.040 | 0.78233438 | 0.780 ± 0.033 |
| Specificity | 0.8792271 | 0.884 ± 0.019 | 0.94685990 | 0.937 ± 0.016 |

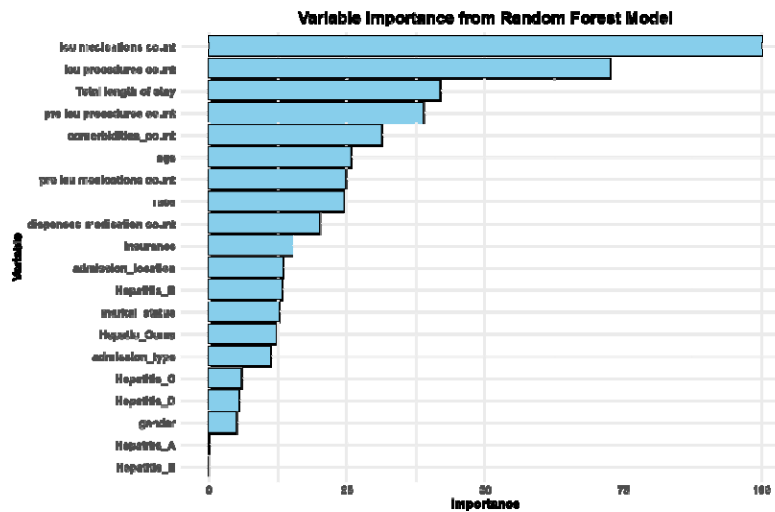
352



353

354

Figure 8: Variable importance for Logistic Regression

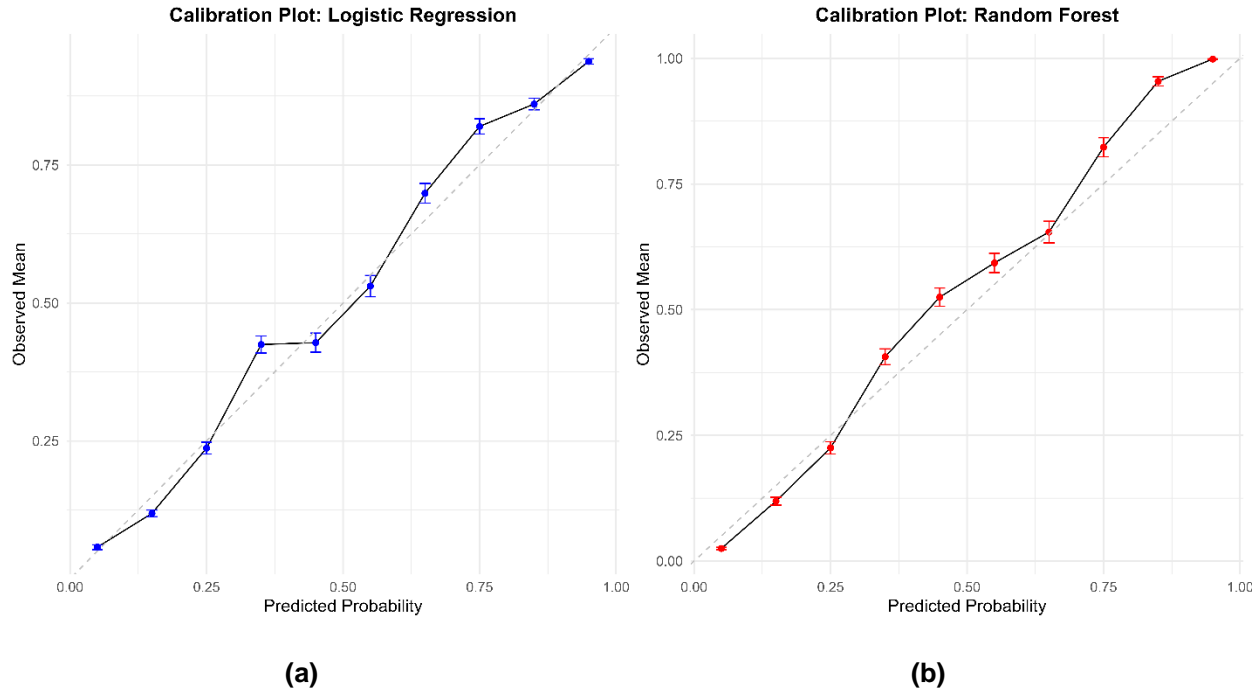


355

356

Figure 9. Variable importance for Random Forest classification

357 The calibration curves shown in **Figure 10** depict the relationship between predicted probabilities and
 358 observed outcomes. **Figure 10a** illustrates the calibration curve for the Logistic Regression model, which
 359 showed reasonable alignment, with minor deviations at the extremes. While **Figure 10b** represents the
 360 calibration curve for the Random Forest model, which demonstrated closer alignment to the diagonal line,
 361 indicating better calibration.



363 **Figure 10.** Calibration curves: (a) Logistic Regression, (b) Random Forest.

364 **Table 3** summarizes the confusion matrix results for both models. The Logistic Regression model showed
 365 an agreement of 0.443 for observed class 0 and 0.2974 for class 1, with sensitivity and specificity values
 366 of 0.773 and 0.696, respectively. The Random Forest model had higher agreement, with 0.483 for class 0
 367 and 0.326 for class 1, a sensitivity of 0.844, and a specificity of 0.763.

368 *Table 3. Confusion matrix metrics for Logistic Regression and Random Forest models*

| Metric | Logistic Regression | | Random Forest | |
|-----------------|---------------------|-----------|---------------|-----------|
| | Class 0 | Class 1 | Class 0 | Class 1 |
| Observed | 0.5726993 | 0.4273007 | 0.5726993 | 0.4273007 |
| Predicted | 0.5729394 | 0.4270606 | 0.5847769 | 0.4152231 |
| Agreement | 0.4431088 | 0.2974701 | 0.4836014 | 0.326125 |
| Sensitivity | 0.7737198 | 0.6961610 | 0.8444246 | 0.7632218 |
| Specificity | 0.6961610 | 0.7737198 | 0.7632218 | 0.8444246 |
| PPV (Precision) | 0.7733956 | 0.6965524 | 0.8269844 | 0.7854216 |
| NPV | 0.6965524 | 0.7733956 | 0.7854216 | 0.8269844 |

369 **Figures 11** and **12** illustrate the confusion matrix results for the **Logistic Regression** and **Random**
 370 **Forest** models, respectively. In **Figure 12**, the Random Forest model demonstrates a stronger alignment
 371 along the diagonal, indicating higher agreement between predicted and observed classes. The color
 372 gradients in both **Figures 11** and **12** represent the predicted probabilities, with darker shades reflecting
 373 higher probabilities of correct classification. Overall, the Random Forest model exhibits improved

374 sensitivity and specificity compared to the Logistic Regression model, aligning with the numerical results
 375 presented in **Table 3**.

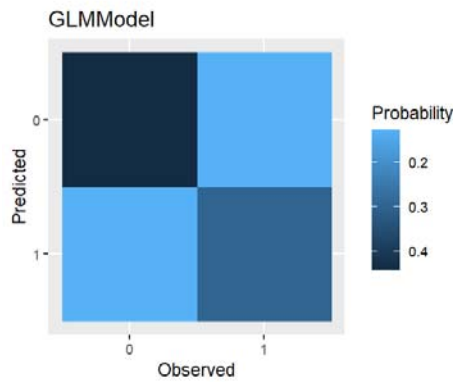


Figure 11. Confusion matrix for Logistic Regression

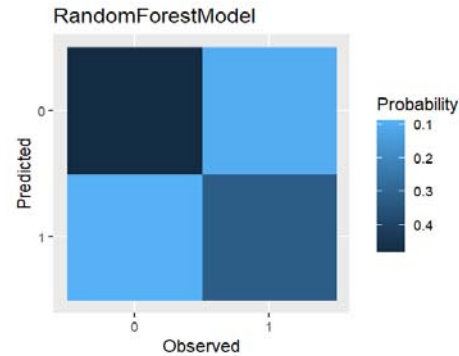


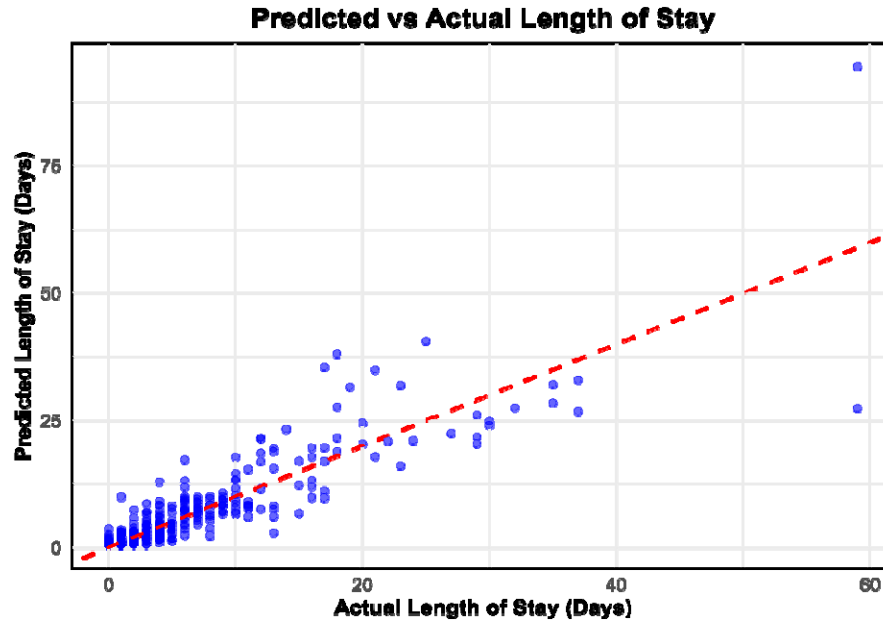
Figure 12. Confusion matrix for Random Forest

376 Total Length of Stay

377 While the previous models focused on classification outcome, the Generalized Additive Model (GAM) with
 378 Negative Binomial distribution and the Random Forest Regression model were utilized to predict total
 379 length of stay as discrete and continuous outcome respectively in the test data and through cross-
 380 validation. The Generalized Additive Model (GAM) with a Negative Binomial family was applied to predict
 381 LoS, as the mean in the dataset for LoS was 3.69, with a variance of 42.54, indicating overdispersion. On
 382 the test data, the model achieved a Root Mean Squared Error (RMSE) of 2.9619, a Mean Absolute Error
 383 (MAE) of 1.4237, and an R-squared (R^2) of 0.7594, demonstrating its ability to explain 75.94% of the
 384 variance in the LoS (**Table 4**). The deviance for the model was 2772.78, and the Akaike Information
 385 Criterion (AIC) was 10410.0, reflecting the model's goodness-of-fit and complexity. Cross-validation was
 386 conducted with 10 folds to assess the model's robustness. The cross-validated RMSE was 5.5212, while
 387 the mean R-squared (R^2) was 0.7601, and the mean MAE was 1.5634. These cross-validation results
 388 highlight slightly higher variability and prediction errors compared to the test data, likely due to the
 389 presence of extreme values or overdispersion in the data. Furthermore, as shown in **Figure 13**, the
 390 majority of predictions are closely aligned with actual values, deviations are observed at higher LoS,
 391 indicating potential challenges in predicting extended stays.

392 *Table 4. Performance metrics of General Additive Model with Negative Binomial*

| Metric | Test Data | Cross-Validation |
|-------------------------|-----------|------------------|
| Root Mean Squared Error | 2.9619 | 5.5212 |
| Mean Absolute Error | 1.4237 | 1.5634 |
| R-squared | 0.7594 | 0.7601 |



393

394 **Figure 13:** Predicted vs. Actual Length of Stay for the Generalized Additive Model (GAM)

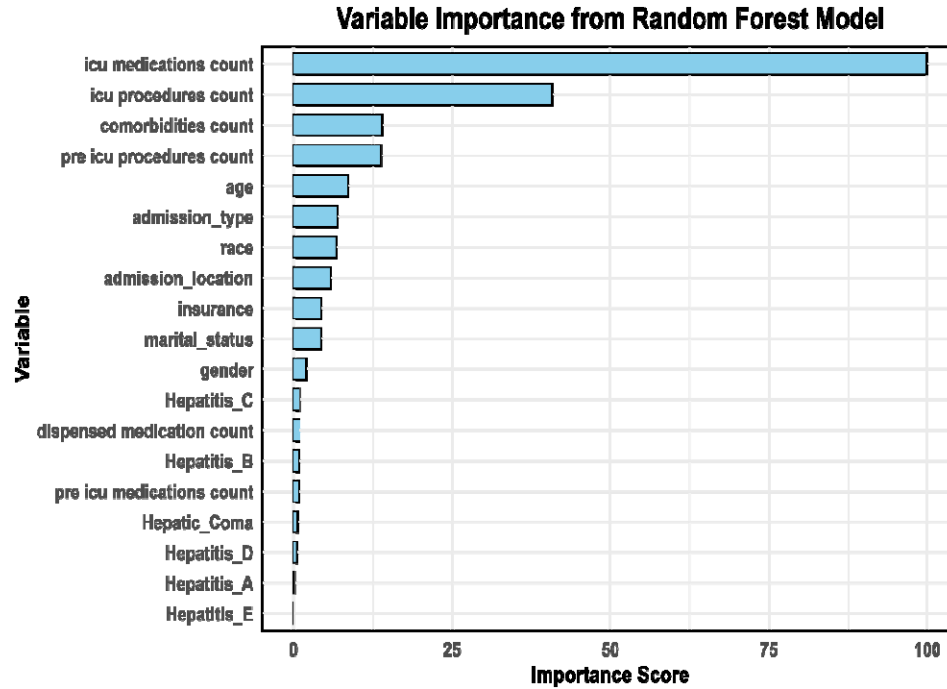
395

396 *Table 5. Performance Metrics of Random Forest Regression Models*

| Metric | Test Data | Log Transformed Test Data | Cross-Validation |
|-------------------------|-----------|---------------------------|-----------------------|
| Root Mean Squared Error | 3.271 | 0.3144 | 0.295 (± 0.020) |
| Mean Absolute Error | 1.360 | 0.2377 | 0.229 (± 0.014) |
| R-squared | 0.755 | 0.821 | 0.838 (± 0.021) |

397 The Random Forest model was subsequently applied after log-transforming the total Length of Stay
 398 (LoS), which resulted in improved performance, achieving an RMSE of 0.3144, an R-squared (R^2) of
 399 0.821, and an MAE of 0.2377 on the test data. Cross-validation further confirmed the model's
 400 consistency, with a mean RMSE of 0.295 (± 0.020), a mean R^2 of 0.838 (± 0.021), and a mean MAE of
 401 0.229 (± 0.014). However, when back transformed to the original scale, the performance metrics were:
 402 RMSE of 3.271, R^2 of 0.755, and MAE of 1.360 (**Table 5**).

403 **Figure 14** illustrates the variable importance of the Random Forest Regression model for predicting
 404 length of stay. The most influential predictors were the number of ICU medications and procedures,
 405 comorbidities count, and pre-ICU procedure count. While clinical factors were most prominent,
 406 sociodemographic variables such as age and race also emerged as notable predictors, highlighting their
 407 potential impact on length of stay alongside clinical characteristics. The scatter plots in **Figure 15**
 408 illustrate the relationship between predicted and actual Length of Stay values for the Random Forest
 409 Regression model. **Figure 15a** presents the scatter plot for the log-transformed data, where a strong
 410 linear relationship is observed, with predictions closely aligning with the actual values. In contrast, **Figure**
 411 **15b** displays the scatter plot for the original Length of Stay data, where greater dispersion is evident,
 412 particularly for higher values.

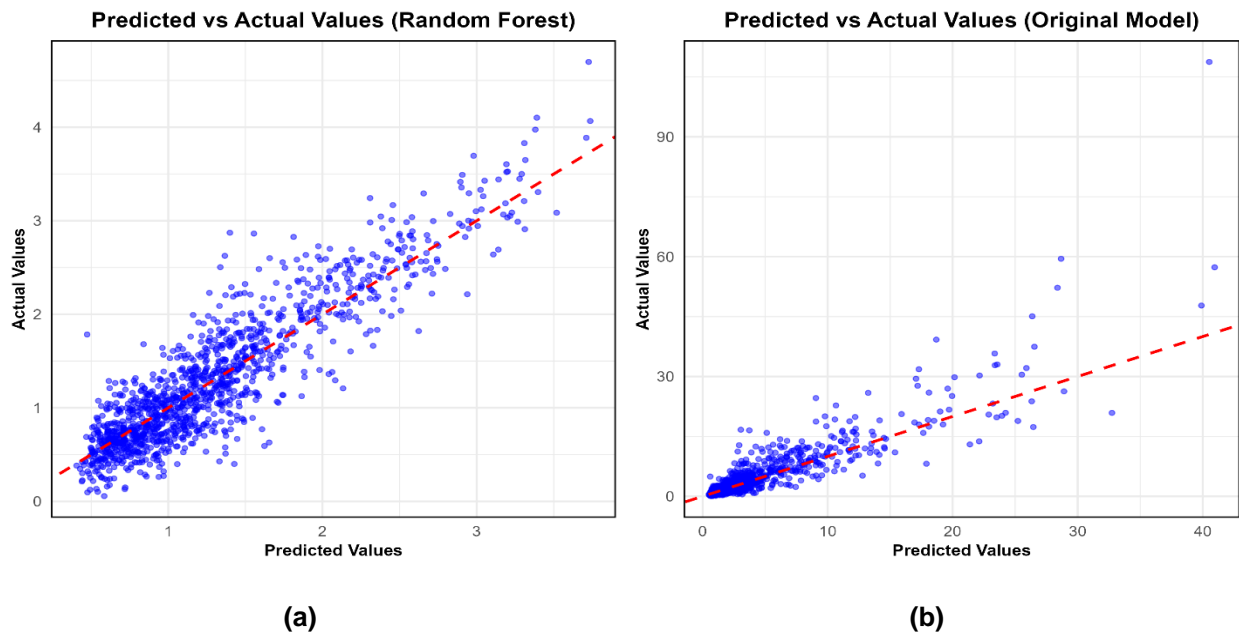


413

Figure 14. Variable Importance for Random Forest Regression Model

414

415



416

Figure 15: Predicted vs. Actual Length of Stay for the Random Forest Regression Model: (a) Log-Transformed Data and (b) Normal Data.

417

418

419

420

421

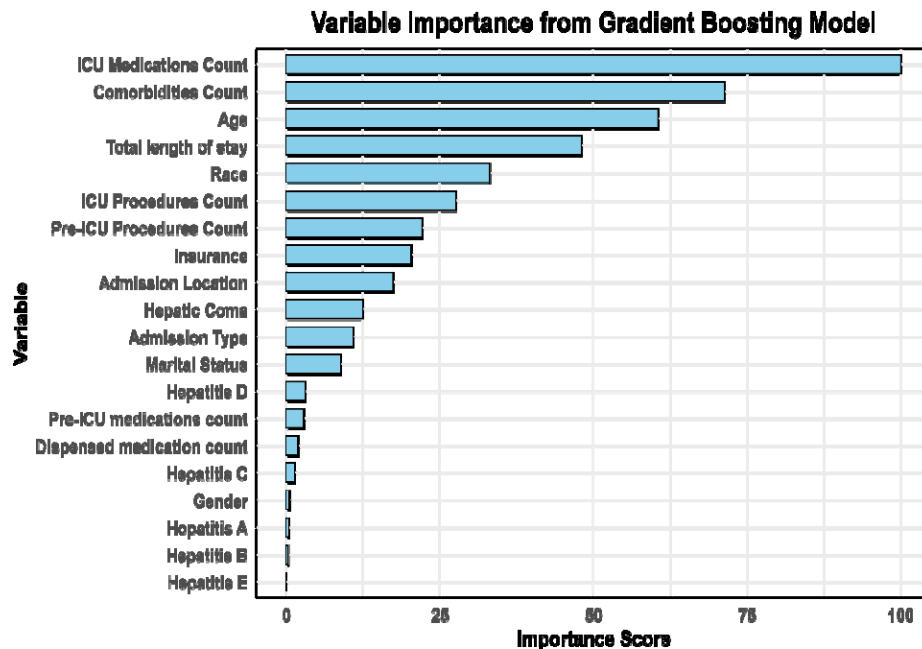
422 Discharge Location

423 Building on the evaluation of model performance for previous outcomes, the focus now turns to predicting
 424 discharge location. The Gradient Boosting Model (GBM) was evaluated on test data, achieving a Brier
 425 Score of 0.589, an accuracy of 0.558, a Kappa of 0.2110, and an ROC AUC of 0.740 (Table 6). Cross-
 426 validation analysis showed a mean Brier Score of 0.5950 (± 0.0206), a mean accuracy of 0.559 (± 0.227),
 427 and a mean Kappa of 0.2077 (± 0.0391), indicating stable performance across folds. The Multinomial
 428 Regression Model produced an accuracy of 0.5639 (95% CI: 52.81% - 59.91%) on the test data, with a
 429 Kappa value of 0.2037 (Table 6). The decay tuning parameter during cross-validation yielded consistent
 430 results, with accuracy values around 0.5658 (± 0.153) and Kappa values near 0.2335 (± 0.30). Significant
 431 predictors ($p < 0.05$) identified in this model included gender, marital status, insurance type, type of
 432 admission, number of comorbidities, and number of ICU procedures and medications. Figure 16
 433 illustrates the variable importance of the Gradient Boosting Model, highlighting the number of ICU
 434 medications, comorbidities, age, total length of stay and race as the most influential predictors.

435 Table 6. Performance metrics of Gradient Boosting and Multinomial Regression model on test and cross-validated
 436 data

| Metric | Gradient Boosting Model | | Multinomial Regression | |
|-------------|-------------------------|-------------------------|------------------------|------------------------|
| | Test Data | Cross Validation | Test Data | Cross Validation |
| Accuracy | 55.87% | 0.5591 (± 0.227) | 0.5639 | 0.5658 (± 0.153) |
| Kappa | 0.2110 | 0.2077 (± 0.0391) | 0.2037 | 0.2335 (± 0.30) |
| Brier Score | 0.5898 | 0.5950 (± 0.0206) | - | - |
| ROC AUC | 0.7407 | 0.7319 (± 0.0282) | - | - |

437



438

439 Figure 16. Variable importance for Gradient Boost model

440

441

442

443 4. DISCUSSION

444 Despite significant advancements in critical care for hepatitis-related conditions, this study highlights the
445 complexities associated with ICU LoS, discharge outcomes (died/alive), and discharge locations.
446 Leveraging ML models, we successfully predicted these outcomes and identified key predictors, including
447 ICU procedures, medication counts, comorbidities, and sociodemographic factors such as age and race.
448 Notably, race consistently emerged as one of the top predictors across all models, underscoring its critical
449 role in influencing health outcomes for hepatitis patients. This finding is particularly significant as it aligns
450 with observed racial disparities in the raw data, suggesting systemic inequities. Additionally, Hepatitis C
451 emerged as the most prevalent condition in our cohort, aligning with global epidemiological trends that
452 underscore its significant burden on critical care systems (Brunner & Bruggmann, 2021; Petruzzello et
453 al., 2016). These findings emphasize the potential of predictive analytics in not only identifying key drivers
454 of health outcomes but also in addressing disparities, optimizing resource allocation, improving patient
455 outcomes, and reducing healthcare costs.

456
457 This study investigated the prediction of ICU LoS for hepatitis patients using ML models, including the
458 Generalized Additive Model (GAM) with a Negative Binomial distribution and Random Forest regression.
459 Random Forest demonstrated superior predictive accuracy, particularly after applying a log transformation
460 to stabilize variance. The log-transformed Random Forest model achieved an R^2 of 0.821 and an RMSE
461 of 0.3144 on test data, effectively capturing complex interactions among predictors. By comparison, GAM
462 with a Negative Binomial distribution achieved an R^2 of 0.7594 and RMSE of 2.9619, effectively
463 addressing overdispersion but yielding higher errors. The key predictors for LoS included ICU medication
464 counts, ICU procedure counts, comorbidities, pre-ICU procedures, and patient age. These findings align
465 with prior studies, such as Gwynn et al. (2019), which demonstrated that higher procedural and
466 pharmacological interventions often correlate with disease severity and resource intensity, leading to
467 longer ICU stays. Similarly, Tola Getachew Bekele et al. (2024) emphasized that prolonged ICU stays are
468 associated with complications, readmissions, and sedative use, underscoring the importance of clinical
469 management strategies to optimize LoS. Although the models performed well for shorter stays, predicting
470 longer LoS posed challenges due to increased variability and extreme values, as reflected in scatter plot
471 (Figures 13). Xu et al. (2022) observed similar challenges, reporting that ML models often struggle to
472 predict longer stays due to the inherent variability of extended ICU durations. This study reflects these
473 findings, with greater error dispersion observed in the upper range of LoS predictions. The performance of
474 Random Forest in predicting LoS aligns with Iwase et al. (2022), who reported high predictive accuracy
475 for ICU stays using Random Forest. Ensemble methods such as Random Forest excel in capturing
476 nonlinear interactions and handling heterogeneous predictors, as further evidenced by Alghatani et al.
477 (2021), who achieved 65% accuracy for LoS predictions in ICU datasets. While GAM with Negative
478 Binomial distribution effectively addressed overdispersion, its higher cross-validation errors indicate that
479 variability in extended LoS remains a challenge. Future research could incorporate additional clinical
480 variables or hierarchical modeling techniques to enhance predictions for longer LoS.

481
482 The prediction of discharge outcomes (alive vs. death) was conducted using Logistic Regression and
483 Random Forest models. Random Forest outperformed Logistic Regression across all metrics, achieving
484 an accuracy of 87.5%, a Kappa of 0.742, and an ROC AUC of 0.95 on test data. Logistic Regression
485 achieved slightly lower metrics, with an accuracy of 81.67%, a Kappa of 0.622, and an ROC AUC of 0.89.
486 These results highlight the robustness of ensemble methods like Random Forest in handling complex
487 interactions among clinical and demographic predictors. Statistically significant predictors for discharge
488 outcomes included ICU medication counts, ICU procedure counts, and total LoS, all of which reflect
489 disease severity and the intensity of care and key variables in predicting LoS. Prolonged LoS has been
490 linked to increased mortality, consistent with Moitra et al. (2016), who found that each additional ICU Day
491 beyond seven days significantly raises mortality risk. Similar findings were reported by Lingsma et al.
492 (2018) and Tola Getachew Bekele et al. (2024), who highlighted the correlation between extended LoS
493 and mortality, emphasizing its role as both a predictor and consequence of severe illness.

494

495 Sociodemographic factors also played a crucial role in predicting discharge outcomes. Race, for instance,
496 emerged as a significant predictor in our study, aligning with Olanipekun et al. (2021), who found
497 disparities in ICU mortality among demographic groups. Age was another critical variable, consistent with
498 Saadatmand et al. (2022), who identified age and Lactate dehydrogenase (LDH) levels as significant
499 determinants of mortality. While LDH levels were unavailable in our dataset, the inclusion of demographic
500 and clinical predictors supported robust outcome predictions.

501
502 The prediction of discharge location presented unique challenges, primarily due to the small dataset of
503 3,875 records and eight discharge categories. Gradient Boosting Model (GBM) and Multinomial
504 Regression yielded moderate performance, with accuracies of 55.87% and 56.39% and Kappa values of
505 0.2110 and 0.2037, respectively. These results highlight the complexities of multi-class prediction in
506 critical care settings with limited data. The Synthetic Minority Oversampling Technique (SMOTE) was
507 applied to address class imbalance but resulted in data loss, which limited the model's performance. Abad
508 et al. (2021) addressed similar challenges using hierarchical classifiers and SMOTE to manage
509 imbalanced multi-class datasets, suggesting that advanced techniques and larger datasets are critical for
510 improving predictive accuracy. Additionally, as noted by Abad et al. (2021), discharge decisions are
511 influenced by subjective factors such as caregiver preferences, resource availability, and social support,
512 which may not be fully captured by clinical and demographic data alone. Top predictors for discharge
513 location included gender, marital status, insurance type, admission type, number of comorbidities, and
514 ICU procedures and medications. These findings align with Mickle and Deb (2022), who highlighted the
515 importance of clinical and patient characteristics in discharge planning. The association of ICU
516 procedures and medications with discharge location reflects the intensity of care required, further
517 emphasizing their relevance in critical care contexts. Future studies should explore hierarchical
518 classifiers, minimize the number of discharge categories, and incorporate additional variables such as
519 functional status and physiological scores to enhance predictive performance. As highlighted by Mickle
520 and Deb (2022), integrating larger datasets and advanced modeling techniques can provide more
521 clinically actionable insights.

522
523 This study has several limitations that warrant further consideration. First, the reliance on the MIMIC-IV
524 database, which reflects a single institution's patient population and care practices, may limit the
525 generalization of the findings to other healthcare settings with differing demographics or protocols.
526 Second, the imbalanced distribution of outcomes, particularly for rare discharge destinations (e.g.,
527 hospice, against medical advice), may have reduced model stability and predictive accuracy. While ML
528 models like Random Forest demonstrated strong performance, their limited interpretability compared to
529 simpler models like Logistic Regression poses challenges for clinical adoption. Future research should
530 focus on validating these findings across diverse datasets to enhance generalizability. Addressing
531 outcome imbalances through advanced techniques without data loss will be critical for better prediction of
532 rare outcomes like hospice discharge or leaving against medical advice. Moreover, future studies could
533 explore fitting a Generalized Additive Model with a Zero-Inflated Negative Binomial (ZINB) distribution to
534 account for overdispersion and zero inflation in LoS data. Expanding this work to include real-time
535 predictive tools integrated into Clinical Decision Support Systems (CDSS) could provide dynamic insights
536 for resource allocation and care optimization.

537 538 **5. CONCLUSION**

539 This study highlights the potential of machine learning in advancing critical care for hepatitis patients. The
540 models identified key predictors of LoS, discharge outcomes, and discharge locations, with clinical factors
541 such as pre-ICU and ICU procedures and medication counts playing a significant role in determining
542 these outcomes. Additionally, sociodemographic factors, including age and race, were consistently
543 identified as important predictors, highlighting the presence of racial disparities among ICU-admitted
544 patients diagnosed with hepatitis. These findings highlight the importance of addressing such disparities
545 while using predictive analytics to optimize resource allocation, improve patient outcomes, and guide
546 targeted interventions. Despite its strengths, the study also revealed challenges such as variability in

547 extended LoS and limited data for multi-class predictions, underscoring the need for further
548 advancements in predictive modeling strategies.

549

550 **6. CONFLICT OF INTEREST**

551 The authors declare that there is no conflict of interest.

552

553 **7. FUNDING STATEMENT**

554 This study did not receive any funding.

555

556 **8. ETHICS STATEMENT**

557 The Indiana University Human Research Protection Program (HRPP) staff determined the analysis done
558 in this study was not human subject research and did not require further Institutional Review Board (IRB)
559 review before conducting the study. The study was conducted under the data use agreement governing
560 MIMIC-IV, ensuring compliance with ethical standards for research involving human subjects. The
561 collection of patient information and creation of the research resource was reviewed by the Institutional
562 Review Board at the Beth Israel Deaconess Medical Center, who granted a waiver of informed consent
563 and approved the data sharing initiative.

564

565 **9. DATA AVAILABILITY STATEMENT**

566 The datasets used for this study are available at [MIMIC-IV v3.0](https://mimic-iv.org/) only for credentialed users following the
567 PhysioNet Credentialed Health Data Use Agreement (DUA).

568

569 **References**

- 570 [1] Abad, Z. S. H., Maslove, D. M., & Lee, J. (2021). Predicting Discharge Destination of Critically Ill
571 Patients Using Machine Learning. *IEEE journal of biomedical and health informatics*, 25(3), 827–
572 837. <https://doi.org/10.1109/JBHI.2020.2995836>
- 573 [2] Alghatani, K., Ammar, N., Rezgui, A., & Shaban-Nejad, A. (2021). Predicting intensive care unit
574 length of stay and mortality using patient vital signs: Machine learning model development and
575 validation. *JMIR Medical Informatics*, 9(5), e21347. <https://doi.org/10.2196/21347>
- 576 [3] An, Q., Rahman, S., Zhou, J., & Kang, J. J. (2023). A Comprehensive Review on Machine
577 Learning in Healthcare Industry: Classification, Restrictions, Opportunities and Challenges.
578 *Sensors*, 23(9), 4178. <https://doi.org/10.3390/s23094178>
- 579 [4] Arnab Kumar Ghosh, Geisler, B. P., & Ibrahim, S. A. (2021). Racial/ethnic and socioeconomic
580 variations in hospital length of stay. *Medicine*, 100(20), e25976–e25976.
581 <https://doi.org/10.1097/md.00000000000025976>
- 582 [5] Azur, M. J., Stuart, E. A., Frangakis, C., & Leaf, P. J. (2011). Multiple imputation by chained
583 equations: what is it and how does it work? *International Journal of Methods in Psychiatric
584 Research*, 20(1), 40–49. <https://doi.org/10.1002/mpr.329>
- 585 [6] Brunner, N., & Bruggmann, P. (2021). Trends of the global Hepatitis C disease burden: Strategies
586 to achieve elimination. *Journal of Preventive Medicine and Public Health*, 54(4).
587 <https://doi.org/10.3961/jpmph.21.151>
- 588 [7] Buuren, S. van, & Groothuis-Oudshoorn, K. (2011). mice: Multivariate Imputation by Chained
589 Equations in R. *Journal of Statistical Software*, 45(3). <https://doi.org/10.18637/jss.v045.i03>
- 590 [8] Centers for Disease Control and Prevention. (2024, April 3). *2022 Viral Hepatitis Surveillance
591 Report | CDC*. <https://www.cdc.gov/hepatitis/statistics/2022surveillance/index.htm>
- 592 [9] Choi, M. H., Kim, D., Choi, E. J., Jung, Y. J., Choi, Y. J., Cho, J. H., & Jeong, S. H. (2022).
593 Mortality prediction of patients in intensive care units using machine learning algorithms based on
594 electronic health records. *Scientific Reports*, 12(1). <https://doi.org/10.1038/s41598-022-11226-4>

- 595 [10] Couronné, R., Probst, P., & Boulesteix, A. L. (2018). Random forest versus logistic regression: A
596 large-scale benchmark experiment. *BMC Bioinformatics*, 19(1), 270.
597 <https://doi.org/10.1186/s12859-018-2264-5>
- 598 [11] Dubin, J. A., Bains, S. S., Rubén Monárrez, Gilmor, R., Swartz, G. N., Katanbaf, R. M., Mont, M.
599 A., Nace, J., & Delanois, R. E. (2024). The Effect of Insurance Type on Length of Stay Following
600 Total Knee Arthroplasty. *The Journal of Arthroplasty*. <https://doi.org/10.1016/j.arth.2024.07.009>
- 601 [12] Eskandari, M., Alizadeh Bahmani, A. H., Mardani-Fard, H. A., Karimzadeh, I., Omidifar, N., &
602 Peymani, P. (2022). Evaluation of factors that influenced the length of hospital stay using data
603 mining techniques. *BMC Medical Informatics and Decision Making*, 22(1).
604 <https://doi.org/10.1186/s12911-022-02027-w>
- 605 [13] Fernandez, G. A., & Vatcheva, K. P. (2022). A comparison of statistical methods for modeling
606 count data with an application to hospital length of stay. *BMC Medical Research Methodology*,
607 22(1). <https://doi.org/10.1186/s12874-022-01685-8>
- 608 [14] Goldberger, A., Amaral, L., Glass, L., Hausdorff, J., Ivanov, P. C., Mark, R., ... & Stanley, H. E.
609 (2000). PhysioBank, PhysioToolkit, and PhysioNet: Components of a new research resource for
610 complex physiologic signals. *Circulation [Online]*. 101 (23), pp. e215–e220.
- 611 [15] Gwynn, M. E., Poisson, M. O., Waller, J. L., & Newsome, A. S. (2019). Development and
612 validation of a medication regimen complexity scoring tool for critically ill patients. *American*
613 *Journal of Health-System Pharmacy*, 76(Supplement_2), S34–S40.
614 <https://doi.org/10.1093/ajhp/zxy054>
- 615 [16] Hastie, T. and Tibshirani, R. (2014). Generalized Additive Models. In *Wiley StatsRef: Statistics*
616 *Reference Online* (eds N. Balakrishnan, T. Colton, B. Everitt, W. Piegorisch, F. Ruggeri and J.L.
617 Teugels). <https://doi.org/10.1002/9781118445112.stat03141>
- 618 [17] Hayes, R. M., Carter, P. R., Gollop, N. D., Reynolds, J., Uppal, H., Sarma, J., Chandran, S., &
619 Potluri, R. (2016). 108 The Impact of Marital Status on Mortality and Length of Stay in Patients
620 Admitted with Myocardial Infarction: Abstract 108 Table 1. *Heart*, 102(Suppl 6), A77.1-A77.
621 <https://doi.org/10.1136/heartjnl-2016-309890.108>
- 622 [18] Hickman, D. (2018, January 3). *Choosing location after discharge wisely*. The Hospitalist.
623 [https://www.the-hospitalist.org/hospitalist/article/155424/transitions-care/choosing-location-after-](https://www.the-hospitalist.org/hospitalist/article/155424/transitions-care/choosing-location-after-discharge-wisely)
624 [discharge-wisely](https://www.the-hospitalist.org/hospitalist/article/155424/transitions-care/choosing-location-after-discharge-wisely)
- 625 [19] Ismail, N., & Jemain, A. (2007). Handling Overdispersion with Negative Binomial and Generalized
626 Poisson Regression Models.
627 https://www.casact.org/sites/default/files/database/forum_07wforum_07w109.pdf
- 628 [20] Johnson, A., Bulgarelli, L., Pollard, T., Gow, B., Moody, B., Horng, S., Celi, L. A., & Mark, R.
629 (2024). MIMIC-IV (version 3.0). PhysioNet. <https://doi.org/10.13026/hxp0-hg59>.
- 630 [21] Johnson, A.E.W., Bulgarelli, L., Shen, L. Et al. MIMIC-IV, a freely accessible electronic health
631 record dataset. *Sci Data* 10, 1 (2023). <https://doi.org/10.1038/s41597-022-01899-x>
- 632 [22] Kwak, C., & Clayton-Matthews, A. (2002). Multinomial logistic regression. *Nursing research*,
633 51(6), 404–410. <https://doi.org/10.1097/00006199-200211000-00009>
- 634 [23] Levin, S., Barnes, S., Toerper, M., Debraine, A., DeAngelo, A., Hamrock, E., Hinson, J., Hoyer, E.,
635 Dungarani, T., & Howell, E. (2020). Machine-learning-based hospital discharge predictions can
636 support multidisciplinary rounds and decrease hospital length-of-stay. *BMJ Innovations*, 7(2),
637 414–421. <https://doi.org/10.1136/bmjinnov-2020-000420>
- 638 [24] Liaw, A., & Wiener, M. (2002). Classification and regression by randomForest. *R News*, 2(3), 18–
639 22. <https://CRAN.R-project.org/doc/Rnews/>
- 640 [25] Lingsma, H. F., Bottle, A., Middleton, S., Kievit, J., Steyerberg, E. W., & Marang-van de Mheen, P.
641 J. (2018). Evaluation of hospital outcomes: The relation between length-of-stay, readmission, and
642 mortality in a large international administrative database. *BMC Health Services Research*, 18(1).
643 <https://doi.org/10.1186/s12913-018-2916-1>

- 644 [26] Mickle, C. F., & Deb, D. (2022). Early prediction of patient discharge disposition in acute
645 neurological care using machine learning. *BMC Health Services Research*, 22(1).
646 <https://doi.org/10.1186/s12913-022-08615-w>
- 647 [27] Moitra, V. K., Guerra, C., Linde-Zwirble, W. T., & Wunsch, H. (2015). Relationship between ICU
648 length of stay and long-term mortality for elderly ICU survivors. *Critical Care Medicine*, 44(4), 1.
649 <https://doi.org/10.1097/ccm.0000000000001480>
- 650 [28] Ng, M., Patrizia Maria Carrieri, Lindila Awendila, Maria Eugenia Socías, Knight, R., & Ti, L.
651 (2024). Hepatitis C virus infection and hospital-related outcomes: A systematic review. *Canadian*
652 *Journal of Gastroenterology & Hepatology*, 2024, 1–23. <https://doi.org/10.1155/2024/3325609>
- 653 [29] Olanipekun, T., Abe, T., Sobukonla, T., Tamizharasu, J., Gamo, L., Kuete, N. T., Bakinde, N.,
654 Westney, G., & Snyder, R. H. (2021). Association between race and risk of ICU mortality in
655 mechanically ventilated COVID-19 patients at a safety net hospital. *Journal of the National*
656 *Medical Association*. <https://doi.org/10.1016/j.jnma.2021.09.003>
- 657 [30] Oracle Corporation. (2023). MySQL: The world's most popular open-source database (Version
658 8.0) [Computer software]. Oracle. <https://www.mysql.com/>
- 659 [31] Oracle. (2021). *What is MySQL?* <https://www.oracle.com/mysql/what-is-mysql/>
- 660 [32] Peres, I. T., Hamacher, S., Oliveira, F. L. C., Thomé, A. M. T., & Bozza, F. A. (2020). What factors
661 predict length of stay in the intensive care unit? Systematic review and meta-analysis. *Journal of*
662 *Critical Care*, 60, 183–194. <https://doi.org/10.1016/j.jcrc.2020.08.003>
- 663 [33] Petruzzello, A., Marigliano, S., Loquercio, G., Cozzolino, A., & Cacciapuoti, C. (2016). Global
664 epidemiology of hepatitis C virus infection: An up-date of the distribution and circulation of
665 hepatitis C virus genotypes. *World Journal of Gastroenterology*, 22(34), 7824.
666 <https://doi.org/10.3748/wjg.v22.i34.7824>
- 667 [34] Pradipta, G. A., Wardoyo, R., Musdholifah, A., Sanjaya, I. N. H., & Ismail, M. (2021). SMOTE for
668 handling imbalanced data problem: A review. *Proceedings of the 2021 Sixth International*
669 *Conference on Informatics and Computing (ICIC)*, 1–8.
670 <https://doi.org/10.1109/ICIC54025.2021.9632912>
- 671 [35] R Core Team. (2023). R: A language and environment for statistical computing (Version 4.4.1)
672 [Computer software]. R Foundation for Statistical Computing. <https://www.R-project.org/>
- 673 [36] RStudio Team. (2023). RStudio: Integrated development for R (Version 2024.4.2.764) [Computer
674 software]. RStudio, PBC. <https://www.rstudio.com/>
- 675 [37] Smith, B. J. (2024). MachineShop: Machine learning models and tools (Version 3.8.0) [R
676 package]. CRAN. <https://cran.r-project.org/package=MachineShop>
- 677 [38] Stoltzfus, J. C. (2011). Logistic regression: A brief primer. *Academic Emergency Medicine*, 18(10),
678 1099–1104. <https://doi.org/10.1111/j.1553-2712.2011.01185.x>
- 679 [39] Tola Getachew Bekele, Melaku, B., Lemlem Beza Demisse, Legese Fekede Abza, & Awol Seid
680 Assen. (2024). Outcomes and factors associated with prolonged stays among patients admitted
681 to adult intensive care unit in a resource-limited setting: A multicenter chart review. *Scientific*
682 *Reports*, 14(1). <https://doi.org/10.1038/s41598-024-64911-x>
- 683 [40] Torgo, L. (2010). *Data mining with R: Learning with case studies*. Chapman and Hall/CRC.
684 Retrieved from <http://www.dcc.fc.up.pt/~ltorgo/DataMiningWithR>
- 685 [41] Venables, W. N., & Ripley, B. D. (2002). *Modern applied statistics with S* (4th ed.). Springer.
686 <https://doi.org/10.1007/978-0-387-21706-2>
- 687 [42] Verma, V.K., Saxena, K., Banodha, U. (2024). Analysis Effect of K Values Used in K Fold Cross
688 Validation for Enhancing Performance of Machine Learning Model with Decision Tree. In: Garg,
689 D., Rodrigues, J.J.P.C., Gupta, S.K., Cheng, X., Sarao, P., Patel, G.S. (eds) *Advanced*
690 *Computing. IACC 2023. Communications in Computer and Information Science*, vol 2053.
691 Springer, Cham. https://doi.org/10.1007/978-3-031-56700-1_30

- 692 [43] Wickham, H. (2016). *ggplot2: Elegant graphics for data analysis*. Springer-Verlag.
693 <https://ggplot2.tidyverse.org>
- 694 [44] Wickham, H., Averick, M., Bryan, J., Chang, W., McGowan, L. D., François, R., Grolemund, G.,
695 Hayes, A., Henry, L., Hester, J., Kuhn, M., Pedersen, T. L., Miller, E., Bache, S. M., Müller, K.,
696 Ooms, J., Robinson, D., Seidel, D. P., Spinu, V., Takahashi, K., Vaughan, D., Wilke, C., Woo, K.,
697 & Yutani, H. (2019). Welcome to the tidyverse. *Journal of Open Source Software*, 4(43), Article
698 1686. <https://doi.org/10.21105/joss.01686>
- 699 [45] Wilimitis, D., & Walsh, C. G. (2023). Practical Considerations and Applied Examples of Cross-
700 Validation for Model Development and Evaluation in Health Care: Tutorial. *JMIR AI*, 2(1), e49023.
701 <https://doi.org/10.2196/49023>
- 702 [46] Wood, S. N. (2011). Fast stable restricted maximum likelihood and marginal likelihood estimation
703 of semiparametric generalized linear models. *Journal of the Royal Statistical Society: Series B*
704 (Statistical Methodology), 73(1), 3–36. <https://doi.org/10.1111/j.1467-9868.2010.00749.x>
- 705 [47] World Health Organization. (2023, October 27). *Hepatitis*. [Wwww.who.int](http://www.who.int).
706 https://www.who.int/health-topics/hepatitis/viral-hepatitis-in-the-world-map#tab=tab_1
- 707 [48] World Health Organization. (2024, April 9). *WHO sounds alarm on viral hepatitis infections*
708 *claiming 3500 lives each day*. [https://www.who.int/news/item/09-04-2024-who-sounds-alarm-on-](https://www.who.int/news/item/09-04-2024-who-sounds-alarm-on-viral-hepatitis-infections-claiming-3500-lives-each-day)
709 [viral-hepatitis-infections-claiming-3500-lives-each-day](https://www.who.int/news/item/09-04-2024-who-sounds-alarm-on-viral-hepatitis-infections-claiming-3500-lives-each-day)
- 710 [49] Wrld, A. (2024, September 10). *Random Forest: The Ultimate Guide to Regression and*
711 *Classification*. Medium. [https://medium.com/@adeevmardia/random-forest-the-ultimate-guide-to-](https://medium.com/@adeevmardia/random-forest-the-ultimate-guide-to-regression-and-classification-33506d6cf865)
712 [regression-and-classification-33506d6cf865](https://medium.com/@adeevmardia/random-forest-the-ultimate-guide-to-regression-and-classification-33506d6cf865)
- 713 [50] Xu, Z., Zhao, C., Scales, C. D., Henao, R., & Goldstein, B. A. (2022). Predicting in-hospital length
714 of stay: a two-stage modeling approach to account for highly skewed data. *BMC Medical*
715 *Informatics and Decision Making*, 22(1). <https://doi.org/10.1186/s12911-022-01855-0>
- 716 [51] Yoshida, K., & Bartel, A. (2022). *tableone: Create 'Table 1' to describe baseline characteristics*
717 *with or without propensity score weights (Version 0.13.2) [R package]*. CRAN. [https://CRAN.R-](https://CRAN.R-project.org/package=tableone)
718 [project.org/package=tableone](https://CRAN.R-project.org/package=tableone)
- 719 [52] Zeileis, A., Kleiber, C., & Jackman, S. (2015). *Regression Models for Count Data in R*. The
720 Comprehensive Archive R Network. [https://cran.r-](https://cran.r-project.org/web/packages/pscl/vignettes/countreg.pdf)
721 [project.org/web/packages/pscl/vignettes/countreg.pdf](https://cran.r-project.org/web/packages/pscl/vignettes/countreg.pdf)
- 722 [53] Zhang, Z., Zhao, Y., Canes, A., Steinberg, D., Lyashevskaya, O., & written on behalf of AME Big-
723 Data Clinical Trial Collaborative Group (2019). Predictive analytics with gradient boosting in
724 clinical medicine. *Annals of translational medicine*, 7(7), 152.
725 <https://doi.org/10.21037/atm.2019.03.29>
- 726 [54] Zwanenburg, A. (2022, December 16). *Performance metrics*. R-Project.org. [https://cran.r-](https://cran.r-project.org/web/packages/familiar/vignettes/performance_metrics_precompiled.html)
727 [project.org/web/packages/familiar/vignettes/performance_metrics_precompiled.html](https://cran.r-project.org/web/packages/familiar/vignettes/performance_metrics_precompiled.html)