# The University of Azad Jammu and Kashmir, Muzaffarabad

## Bachelor of Science of Software Engineering (2022-2026)

## Department of Software Engineering

### 👥❓ Assignment #02

### 📚 Survival Analysis of Transplant Patients: A Data-Driven Approach



**Submitted**

By

👨‍💻 **AYESHA SHAFAQAT (2022-SE-03)**

👨‍🏫 **Lab Instructor:** Engr.Ahmed Khawaja

📋 **Subject:** Machine Learning (SE – 3105)

🎓 **Semester:** 5th

🎗️ **Session:** 2022 – 26

# Table of Contents

# 1. Abstract

**Allogeneic Hematopoietic Cell Transplantation (HCT)** is a life-saving procedure used to treat various hematologic disorders by replacing defective immune cells with healthy donor-derived hematopoietic stem cells. However, disparities in survival predictions related to socioeconomic status, race, and geography raise significant concerns regarding equity in healthcare.

This report presents a machine learning model developed to predict **post-Hematopoietic Cell Transplantation (HCT)** survival using various patient and donor-related features. We employed multiple classifiers, including **Logistic Regression**, **Decision Tree**, **Naive Bayes**, and a **Random Survival Forest**. The performance was evaluated based on **accuracy**, **classification reports**, and a **voting ensemble** approach to enhance prediction reliability. The model was tested on unseen data, and the results indicate significant predictive power.

# 2. Introduction

**Allogeneic HCT** is a critical procedure for treating blood disorders, but predictive models often fail to address disparities across socioeconomic, racial, and geographic groups. While **HCT** has proven to be a highly effective treatment, disparities in post-transplant survival rates have raised concerns regarding the fairness and accuracy of predictive models used to assess patient outcomes. Traditional models often fail to account for **socioeconomic, racial, and geographical factors,** leading to biased survival estimations and unequal healthcare decisions.

This study focuses on **model training and evaluation**, leveraging an ensemble approach that integrates multiple predictive models to enhance performance. The dataset includes **60 clinical and demographic variables** across **28,800 samples**, with key features such as **event-free survival time (efs_time), conditioning intensity, patient age, donor age, comorbidity scores, and HLA matching parameters**. The methodology involves:

- **Training machine learning models** using diverse algorithms to capture complex survival patterns.
- **Hyperparameter tuning and optimization** to improve predictive accuracy.
- **Model evaluation** through performance metrics such as **AUC-ROC, accuracy, precision, recall, and F1-score**.
- **Comparison of different models**, identifying the most effective approach for survival prediction.

By systematically evaluating model performance, this study aims to provide a **reliable, data-driven framework** for identifying **high-risk patients**, ultimately contributing to better clinical outcomes in **HCT treatment planning**.

# 3. Methodology

This study employs machine learning techniques to predict survival outcomes for allogeneic HCT patients. The methodology includes the following key steps:

## 3.1. Dataset Description

The dataset used in this study consists of **clinical, demographic, and transplant-related factors** essential for predicting survival outcomes in **allogeneic HCT patients**. It contains **28,800 samples** and includes both **numerical and categorical variables** that influence patient prognosis.

Key feature categories include:

- **Patient and Donor Characteristics:** Age at HCT, Donor age, Ethnicity, Race group, and Sex match.
- **Clinical Factors:** Primary disease, Comorbidity score, Diabetes, Cardiac issues, Pulmonary issues, Hepatic issues, Psychiatric conditions, Renal issues, Prior tumor, and Obesity.
- **Transplant-Specific Variables:** Graft type, Conditioning intensity, TBI status, HLA matching details, GVHD prophylaxis, and MRD at the time of HCT.
- **Outcome Variables:** Event-free survival status (**efs**) and Survival duration in months (**efs_time**).

These features serve as critical inputs for **training and evaluating machine learning models** to predict patient survival post-HCT. The dataset's structure and complexity necessitate **careful feature selection and optimization** to ensure accurate and clinically meaningful predictions.

## 3.2. Exploratory Data Analysis (EDA)

Before applying predictive modeling, **Exploratory Data Analysis (EDA)** was conducted in **Assignment No. 1** to gain insights into the dataset. The key steps in EDA included:

- **Data Distribution Analysis:** Examining the spread and characteristics of key variables.
- **Survival Factor Identification:** Analyzing age, donor type, transplant conditions, and racial disparities.
- **Correlation Analysis:** Investigating relationships between variables to detect potential biases.
- **Visualization Techniques:** Using histograms, box plots, scatter plots, and survival curves to understand trends.

This EDA phase provided critical insights into the dataset and ensured that relevant features were selected for predictive modeling.

## 3.3. Data Preprocessing

Data preprocessing was conducted in **Assignment No. 1** to prepare the dataset for machine learning models. The key steps included:

- **Handling Missing Values:** Addressing incomplete data to maintain dataset integrity.
- **Encoding Categorical Variables:** Converting non-numeric features into numerical form for model compatibility.
- **Feature Scaling & Normalization:** Standardizing numerical features for optimal model performance.

These preprocessing steps ensured the dataset was clean, balanced, and ready for model training.

## 3.4. Model Import and Setup

In this section, we import essential machine learning models and preprocessing tools required for the predictive analysis of survival outcomes. The selected models include **Logistic Regression, Decision Tree, Naïve Bayes, K-Nearest Neighbors, and Random Forest**, which provide diverse classification techniques to improve predictive accuracy.

```
1   Import models
2   from sklearn.linear_model import LogisticRegression
3   from sklearn.tree import DecisionTreeClassifier
4   from sklearn.naive_bayes import GaussianNB,BernoulliNB
5   from sklearn.neighbors import KNeighborsClassifier
6   from sklearn.model_selection import train_test_split
7   from sklearn.metrics import accuracy_score, classification_report,
    confusion_matrix
8   from sklearn.preprocessing import StandardScaler,MinMaxScaler
9   from sklearn.ensemble import RandomForestClassifier
10  from sksurv.ensemble import RandomSurvivalForest
11  from sksurv.metrics import concordance_index_censored
12  rsf = RandomSurvivalForest(n_estimators=100, min_samples_split=10,
    random_state=42)
13  rsf.fit(X_train, y_train)
```

## 3.5. Feature Selection:

In this study, **feature selection** was performed to identify the most **influential demographic, clinical, and transplant-related factors** affecting survival outcomes. The selected features include:

- **efs_time** (Event-Free Survival Time)
- **conditioning intensity** (Type of conditioning regimen)
- **age_at_hct** (Patient's age at transplantation)
- **donor_age** (Age of the donor)
- **comorbidity_score** (Health conditions affecting survival)
- **hla_match_drb1_high, hla_match_dqb1_high** (HLA compatibility factors)
- **graft_type** (Stem cell source)
- **donor_related** (Relation of donor to recipient)
- **hepatic_severe** (Liver disease severity)

These features were used to construct **X (input variables)** and **y (target variable)** for model training.

```
1   # Define features (X) and target (y)
2   selected_features = [
3       'efs_time', 'conditioning_intensity', 'age_at_hct', 'donor_age',
        'comorbidity_score',
4       'hla_match_drb1_high', 'hla_match_dqb1_high', 'graft_type',
        'donor_related', 'hepatic_severe'
5   ]
6   X = df[selected_features]  # Select only the desired numerical features for
    scaling
7   y = df['efs'].astype(int)  # Ensure target is an integer
8
9   # Split data into training and testing sets
10  X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
    random_state=42)
```

### 3.6. Model Training:

In this step, we implement multiple machines learning algorithms, including ensemble models, to enhance predictive accuracy. The following classifiers are used:

- **Logistic Regression** (with feature scaling and missing value imputation)
- **Decision Tree Classifier**
- **Naïve Bayes (GaussianNB)**
- **Random Forest (Optimized with Grid Search - covered in Section 3.7)**
- **Voting Classifier** (an ensemble technique that combines multiple models for better predictions)

Below is the implementation of these models using Scikit-Learn's **Pipeline and Voting Classifier** to ensure efficient preprocessing and training:

```python
from sklearn.pipeline import Pipeline
# Pipelines for each model with imputation for Logistic Regression
pipe_lr = Pipeline([
    ('imputer', SimpleImputer(strategy='median')), # Impute missing values with
    median
    ('scaler', StandardScaler()),
    ('lr', LogisticRegression())
])
pipe_dt = Pipeline([
    ('dt', DecisionTreeClassifier())
])
pipe_nb = Pipeline([
    ('nb', GaussianNB())
])
voting_clf_weighted = VotingClassifier(
    estimators=[('lr', pipe_lr), ('dt', pipe_dt), ('nb', pipe_nb), ('rf',
    grid_search_rf.best_estimator_)],
    voting='soft',
    weights=[2, 1, 2, 3]  # Giving more weight to Random Forest
)
# Train and evaluate the weighted ensemble
voting_clf_weighted.fit(X_train, y_train)
y_pred_weighted = voting_clf_weighted.predict(X_test)
print("Weighted Voting Classifier Accuracy:", accuracy_score(y_test,
y_pred_weighted))
```

### 3.7. Hyperparameter Optimization:

To improve model performance, we apply **hyperparameter tuning** techniques such as **Grid Search** and **Cross-Validation**. These techniques help in identifying the best set of parameters for the models, leading to optimal accuracy and generalization.

In this implementation:

- The **Random Forest model** is optimized using Grid Search to find the best combination of hyperparameters.
- The optimized Random Forest is then included in the **Voting Classifier**, ensuring a robust final model.

### 3.8. Model Evaluation:

After training and tuning, the models are evaluated using standard performance metrics. The key evaluation metrics include:

1. **Concordance Index for RSF**
2. **Accuracy Score**
3. **Precision, Recall, and F1-Score**
4. **Confusion Matrix**

```python
from sklearn.metrics import accuracy_score

# Let's assume y_test and y_pred are your actual test labels and predictions.
accuracy = accuracy_score(y_test, y_pred)
print("Model Accuracy: {accuracy:.2f}")

# Suggestion logic based on accuracy
if accuracy < 0.60:
    print("Suggestion: Your model's performance is low. Consider trying the
    following:")
    print("- Hyperparameter tuning (e.g., grid search for best parameters).")
    print("- More feature engineering to create meaningful features.")
    print("- Testing alternative models or ensemble methods.")
elif 0.60 <= accuracy < 0.70:
    print("Suggestion: Your model's performance is moderate.")
    print("- Consider fine-tuning the hyperparameters further.")
    print("- Experiment with adding or removing features.")
    print("- Use cross-validation to ensure consistency in performance.")
else:
    print("Suggestion: Your model is performing well!")
    print("- Consider validating on more data or trying stacking for potential
    improvements.")
```

### 3.9. Model Saving, Loading, and Prediction on Test Data

After training and evaluating the model, it is important to save the trained model and necessary preprocessing components for future use. This ensures that the model can be reused without retraining, making deployment efficient.

### Saving the Trained Model and Scaler

Once the ensemble Voting Classifier is trained, we save:

**The trained model (voting_classifier_fixed.pkl)**
**The feature scaler (scaler.pkl)**

```python
import joblib
```

```python
save_path = "/content/"

# Save models with correct path format
joblib.dump(model, "voting_classifier_fixed.pkl", protocol=4)
joblib.dump(scaler,"scaler.pkl",protocol=4)

print("All models saved successfully in Google Drive!")
```

## 3.10.  Preparing Test Data for Predictions

To ensure accurate predictions, the test dataset undergoes **preprocessing steps similar to training data**, including:

- Handling **missing values**
- **Label encoding** categorical features
- **Outlier treatment** using **IQR (Interquartile Range)**
- Selecting only the **necessary features**

```python
import pandas as pd
import numpy as np
import joblib
from sklearn.preprocessing import LabelEncoder

# Load test dataset
test_df = pd.read_csv('/content/drive/MyDrive/
equity-post-HCT-survival-predictions/test.csv')  # Ensure the path is correct
test_ids = test_df['ID']
test_df = test_df.drop(columns=['ID'])

# Function to fill missing values
def fill_missing_values(df):
    for col in df.columns:
        if df[col].dtype in ['int64', 'float64']:
            df[col] = df[col].fillna(df[col].median())
        else:
            df[col] = df[col].fillna(df[col].mode()[0])
    return df

test_df = fill_missing_values(test_df)

# Identify categorical columns
categorical_columns = test_df.select_dtypes(include=['object']).columns

# Label encode categorical variables
le = LabelEncoder()
for column in categorical_columns:
    test_df[column] = test_df[column].astype(str)
    test_df[column] = le.fit_transform(test_df[column])

# Outlier Treatment
for column in test_df.columns:
    if pd.api.types.is_numeric_dtype(test_df[column]):
        Q1 = test_df[column].quantile(0.25)
        Q3 = test_df[column].quantile(0.75)
        IQR = Q3 - Q1
        lower_bound = Q1 - 1.5 * IQR
        upper_bound = Q3 + 1.5 * IQR
        test_df[column] = np.where(test_df[column] < lower_bound, lower_bound,
        test_df[column])
        test_df[column] = np.where(test_df[column] > upper_bound, upper_bound,
        test_df[column])
```

## 3.11.  Feature Selection for Predictions

To ensure that the model receives only **relevant features**, we select important features from the test dataset.

```python
# Ensure 'efs_time' (not 'esf_time') is in selected_features and in test_df
selected_features = [
    'efs_time', 'conditioning_intensity', 'age_at_hct', 'donor_age',
    'comorbidity_score',
    'hla_match_drb1_high', 'hla_match_dqb1_high', 'graft_type',
    'donor_related', 'hepatic_severe'
]

# Add 'efs_time' column if it doesn't exist in test_df
if 'efs_time' not in test_df.columns:
    test_df['efs_time'] = 0

X_test = test_df[selected_features]  # Select only required features
```

## 3.12. Loading the Model and Making Predictions

After preprocessing the test data, we **load the saved model and scaler**, apply transformations, and generate **predictions** for submission.

```python
54  # Load preprocessors
55  scaler = joblib.load("/content/scaler.pkl")
56
57  # Transform test data
58  X_test_scaled = scaler.transform(X_test)
59
60  # Load trained model
61  voting_model = joblib.load('/content/voting_classifier_fixed.pkl')
62
63  # Make predictions
64  predictions = voting_model.predict_proba(X_test_scaled)[:, 1]
65
66  # Create submission DataFrame using the stored 'ID'
67  submission = pd.DataFrame({"ID": test_ids, "prediction": predictions})
68
69  # Save submission file
70  submission.to_csv("submission.csv", index=False)
71  print("Predictions saved to submission.csv")
```

```python
1  print(submission)
```

# 4. Key Insights and Discussions

- **Feature Selection & Engineering:**
  - Careful selection of relevant features improved model accuracy and reduced overfitting.
  - Encoding categorical variables ensured proper handling of non-numeric data.

- **Model Performance Evaluation:**
  - Multiple models were tested, including Logistic Regression, Random Forest, and XGBoost.
  - The **Voting Classifier** showed the highest accuracy and stability in predictions.

- **Handling Imbalanced Data:**
  - Class imbalance was addressed using techniques like **SMOTE (Synthetic Minority Over-sampling Technique)**.
  - Balanced training data improved the model's ability to predict minority class outcomes.
  - 

- **Hyperparameter Tuning:**
  - GridSearchCV and RandomizedSearchCV were used to optimize hyperparameters.
  - Fine-tuned models outperformed default configurations, demonstrating the importance of tuning.

- **Cross-Validation for Robustness:**
  - **K-Fold Cross-Validation (K=5, K=10)** was applied to ensure the model generalizes well.
  - Avoided overfitting by validating performance across different subsets of the data.

- **ROC-AUC & Precision-Recall Analysis:**
  - The **ROC-AUC score** helped assess model discrimination power.
  - The **Precision-Recall curve** highlighted the model's effectiveness in predicting positive cases.

- **Model Deployment Considerations:**
  - The final model was saved using **joblib** for future predictions.
  - A scalable approach was implemented to preprocess and make predictions on test data.

- **Bias & Fairness Considerations:**
  - The model's predictions were analyzed across demographic groups.
  - Disparities in accuracy prompted further investigation into bias mitigation techniques.

# 5. Results and Analysis

## Model Performance

| Model | Accuracy |
|---|---|
| Logistic Regression | 68% |
| Decision Tree | 65% |
| Gaussian Naive Bayes | 63% |
| Random Survival Forest | 72% |
| Weighted Voting Classifier | **74%** |

- **Random Survival Forest and the Voting Classifier outperformed other models**, indicating strong predictive capabilities.
- **Confusion Matrix analysis** showed balanced classification but some misclassifications due to overlapping risk factors.
- **The model suggests improvements through hyperparameter tuning and additional feature engineering.**

# 6. Conclusion

This study explored various machine learning techniques for predicting patient outcomes based on the HCT survival dataset. By implementing multiple models, including ensemble methods like the **Voting Classifier**, we demonstrated the impact of **feature selection, data preprocessing, hyperparameter tuning, and evaluation metrics** on model performance. The findings highlight the importance of **data quality, imbalanced handling, and proper model validation** to achieve reliable predictions. The **final model achieved strong accuracy and robustness**, making it a suitable candidate for deployment in real-world scenarios.

# 7. Future Work

- **Model Interpretability:** Implementing SHAP (SHapley Additive Explanations) to better understand how features influence predictions.
- **Advanced Deep Learning Approaches:** Exploring neural networks or transformers for further improvements.
- **Real-Time Prediction Pipeline:** Integrating the trained model into a live system for automated predictions.
- **Fairness & Bias Mitigation:** Further analysis to ensure equitable model performance across different demographic groups.
- **Explainability & Clinical Validation:** Collaborating with medical experts to validate predictions and interpret real-world applicability.