

Survival Analysis and Predictive Modeling of Leukemia Patients Across Different Regions of India

Submitted in partial fulfilment of the requirement for the degree of

Master of Science

In

Business Statistics

by

Abhinav V Sunil

24MBS0066

Under the Guidance of

Dr. Jitendra Kumar

School of Advanced Sciences

VIT – Vellore



November, 2025

ABSTRACT

This work explores how long leukemia patients live across various parts of India, blending standard stats with newer algorithms. Rather than relying solely on classic tools, it uses computer-based models to uncover hidden trends. Its aims are twofold: identifying major factors affecting survival like health status, background traits, and gene changes - as well as creating accurate forecasting tools for doctors. Through methods like Kaplan–Meier estimates, Log-Rank comparisons, ANOVA checks, and Cox modeling, results indicate tumor stage, size, therapy type, plus genes such as EGFR and KRAS clearly link to patient outcomes. Besides standard stats, ML techniques such as Random Forest and XGBoost showed better prediction accuracy, achieving AUCs close to 0.80; meanwhile, a neural network was used to detect intricate variable relationships. Combining classical analysis with modern modeling clarifies survival trends in Indian leukemia cases, leading to refined risk grouping, personalized treatment plans, along with stronger, evidence-based clinical guidance.

CONTENTS

	Page No.
Acknowledgement	4
Abstract	5
Table of Contents	6, 7 ,8
List of Figures	8
List of Tables	9
Abbreviations	9,10,11
1. INTRODUCTION	12
1.1. OBJECTIVE	12
1.2. MOTIVATION	12
1.3. LITERATURE SURVEY	13,14,15
1.4. RESEARCH GAP	15
1.5. HYPOTHESIS	16
2. PROJECT DESCRIPTION	17
2.1 PROJECT GOALS	17
3. TECHNICAL SPECIFICATIONS	18
4. DESIGN APPROACH AND DETAILS	19
4.1.1 DATASET	19
4.1.2 SOFTWARE	13
4.1.3 LIBRARIES AND TOOLS USED	20,21
4.2 APPROACH AND METHODS	21
4.2.1 DATA PRE-PROCESSING	15

4.2.2 EXPLORATORY DATA ANALYSIS (EDA)	21
4.2.3 SURVIVAL ANALYSIS METHODS	22
4.2.4 MACHINE LEARNING MODELLING	23
4.2.5 DEEP LEARNING MODELLING	23
4.2.6 RESULT INTERPRETATION	23
4.3 CODES AND STANDARD	23
5 SCHEDULE, TASK AND MILESTONES	33
6. PROJECT OUTPUTS	34 - 42
7. RESULT AND DISCUSSION	43
7.1 SURVIVAL ANALYSIS RESULTS	43
7.1.1 SURVIVAL ANALYSIS RESULTS	43
7.1.2 KAPLAN–MEIER SURVIVAL CURVES BY STAGE	43
7.1.3 EARLY VS LATE STAGE SURVIVAL COMPARISON	43
7.2 TREATMENT-BASED SURVIVAL DIFFERENCES	44
7.2.1 KAPLAN–MEIER CURVES BY TREATMENT TYPE	44
7.3 COX PROPORTIONAL HAZARDS MODEL 43	44
7.3.1 INTERPRETATION OF MULTIVARIATE COX RESULTS	44
7.4 MACHINE LEARNING MODELS	45
7.4.1 RANDOM FOREST PERFORMANCE	45
7.4.2 XGBOOST PERFORMANCE	46
7.5 DEEP LEARNING RESULTS	46
7.5.1 NEURAL NETWORK LEARNING PATTERNS	46
7.6 SHAP EXPLAINABILITY INSIGHTS	46
7.6.1 GLOBAL INTERPRETABILITY	47
7.6.2 DETAILED FEATURE BEHAVIOR	47

7.7 INTEGRATED DISCUSSION	47
8. LIMITATION	49
9. CONCLUSION	50
10. REFERENCES	51 - 52

LIST OF FIGURES

Figure No.	Title	Page No.
1	Survival Differences by Cancer Stage: ANOVA, Kruskal-Wallis Test, and Kaplan–Meier Curves	34
2	Kaplan–Meier Survival Curve for EGFR Positive vs Negative Patients	34
3	Kaplan–Meier Survival Curve for KRAS Positive vs Negative Patients	35
4	Log-Rank Test Results and Cox Proportional Hazards Model Summary	35
5	ANOVA and Tukey HSD Post-Hoc Test for Treatment Type Effects on Survival	36
6	Kaplan-Meier Plot by Treatment	36
7	Cox Proportional Hazards Model Includes Treatment	37
8	Concordance Includes Treatment	37
9	Cox Model Partial Effects: Survival Probability by Tumor Size Profiles Adjusted for Age and BMI	38
10	Cox Proportional Hazards Model: Continuous Effects of Tumor Size, Age, and BMI on Survival	38
11,12	Proportional Hazards Assumption Diagnostics for Full Cox Model	39,40
13	Concordance Index for Full Cox Model	41
14	Full Multivariate Cox Regression: Adjusted Hazard Ratios for All Covariates	41
15	Random Forest Classification: Top Predictive Features and Model Performance for Event Prediction	42
16	Neural Network Training Accuracy and Validation Performance for Event Prediction	42

LIST OF TABLES

Table No	Title	Page No.
1	Schedule, Task and milestones	33

LIST OF ABBREVIATIONS

Abbreviation	Full Form
AML	Acute Myeloid Leukemia
ALL	Acute Lymphoblastic Leukemia
CML	Chronic Myeloid Leukemia
KM	Kaplan–Meier
PH	Proportional Hazards
CPH	Cox Proportional Hazards
HR	Hazard Ratio
CI	Confidence Interval
EDA	Exploratory Data Analysis
ML	Machine Learning
DL	Deep Learning
RF	Random Forest
XGB	XGBoost
NN	Neural Network
ROC	Receiver Operating Characteristic
AUC	Area Under the Curve
SHAP	SHapley Additive exPlanations
EHR	Electronic Health Records
BMI	Body Mass Index
EGFR	Epidermal Growth Factor Receptor

Abbreviation	Full Form
KRAS	Kirsten Rat Sarcoma Viral Oncogene
API	Application Programming Interface
CPU	Central Processing Unit
GPU	Graphics Processing Unit
SD	Standard Deviation
IQR	Interquartile Range
CSV	Comma-Separated Values
RFE	Recursive Feature Elimination
TP	True Positive
Abbreviation	Full Form
AML	Acute Myeloid Leukemia
ALL	Acute Lymphoblastic Leukemia
CML	Chronic Myeloid Leukemia
KM	Kaplan–Meier
PH	Proportional Hazards
CPH	Cox Proportional Hazards
HR	Hazard Ratio
CI	Confidence Interval
EDA	Exploratory Data Analysis
ML	Machine Learning
DL	Deep Learning
RF	Random Forest
XGB	XGBoost
NN	Neural Network
ROC	Receiver Operating Characteristic

Abbreviation	Full Form
AUC	Area Under the Curve
SHAP	SHapley Additive exPlanations
EHR	Electronic Health Records
BMI	Body Mass Index
EGFR	Epidermal Growth Factor Receptor
KRAS	Kirsten Rat Sarcoma Viral Oncogene
API	Application Programming Interface
CPU	Central Processing Unit
GPU	Graphics Processing Unit
SD	Standard Deviation
IQR	Interquartile Range
CSV	Comma-Separated Values
RFE	Recursive Feature Elimination
TP	True Positive
FP	False Positive
TN	True Negative
FN	False Negative
OHE	One-Hot Encoding

1.INTRODUCTION

1.1 OBJECTIVE

The aim of this study is to examine regional differences in leukemia survival rates across India using traditional statistical methods as well as modern machine learning techniques, combining established approaches with recent innovations. This work investigates how factors like age, medical background, or genetic variation influence patient results, focusing on measuring real-world impact rather than just finding associations. A further objective involves creating practical prediction models that support clinicians in estimating personal risk levels more accurately. An essential component evaluates disease development effects on expected lifespan by applying Kaplan–Meier estimates, Log-Rank tests, together with ANOVA analyses - to determine whether diagnosis timing actually affects outcome. Beyond staging, gene activity - such as in EGFR and KRAS - is analyzed using Cox regression while survival graphs highlight ties to improved or poorer results. For treatment insight, various drug approaches are compared; meanwhile, factors like age, tumor dimensions, and body mass index undergo multivariate scrutiny to assess impact on mortality likelihood. Based on these findings, predictive systems relying on Random Forest, XGBoost, or neural nets emerge, detecting complex interactions and sorting patients into risk tiers - with accuracy tested via ROC–AUC scores and key predictor evaluations. Overall, this study merges standard statistical analysis with modern machine learning to map survival trends among Indian leukemia patients and support personalized care decisions.

1.2 MOTIVATION

This research responds to growing use of data in tackling tough health issues like leukemia, a disease still causing serious diagnosis and survival challenges across India. Because outcomes differ widely based on location, age, therapy type, or genetic makeup, examining these factors demands solid numerical methods. With medicine shifting toward treatments backed by proof and tailored to individuals, merging stats models with learning algorithms may reveal useful trends for doctors. Addressing current shortcomings, this work uses up-to-date analysis tools to identify key clinical signs and gene traits linked closely to how long patients live.

Another reason stems from how machine learning and deep neural networks are increasingly shaping medical studies. Although traditional survival models reveal key risk patterns, precisely forecasting individual outcomes matters just as much - particularly when spotting high-risk patients. Combining classic statistics with modern computing techniques enhances clarity and accuracy alike, supporting better predictions in leukemia cases. Moreover,

organized clinical records from various parts of India open up chances to detect regional differences, helping doctors adapt care across varied health systems.

This study draws from a genuine curiosity about using statistics and data analysis in practical health contexts. Exploring how these methods aid key clinical choices proves both interesting and useful. Applying techniques like survival analysis, machine learning, or neural networks to leukemia data helps sharpen technical abilities - while contributing valuable perspectives to an area where more accurate forecasts might enable timelier diagnoses, smarter therapies, maybe even extended patient outcomes.

1.3 LITERATURE SURVEY

Leukemia prognosis research has evolved significantly over the past several decades, beginning with the introduction of fundamental survival analysis techniques that remain central to modern clinical analytics. The seminal Kaplan–Meier estimator developed by Kaplan and Meier [5] provided the first robust non-parametric approach for estimating survival probabilities in the presence of censored data, a common feature in oncology studies. This was later complemented by the Cox Proportional Hazards model introduced by Cox [3], which enabled the assessment of multiple risk factors simultaneously while maintaining semi-parametric flexibility. These foundational works continue to underpin survival modelling in leukemia research and set the stage for more advanced methodologies.

As understanding of leukemia biology deepened, the role of genetic and molecular markers gained increasing prominence. Mutations such as EGFR and KRAS have been repeatedly highlighted as influential in determining treatment response and disease progression. Patel and Deininger [8] provided early evidence supporting the prognostic value of these mutations, while Cortes et al. [2] demonstrated their significant impact on clinical outcomes, emphasizing their integration into risk stratification frameworks. Additional studies such as Chatterjee et al. [11] further reinforced the importance of genetic determinants, underscoring their role in guiding personalized therapeutic interventions for hematologic malignancies.

In the Indian context, population-based studies have revealed pronounced regional disparities in leukemia incidence and survival. The Indian Council of Medical Research (ICMR) report [4] documented substantial variation across states, suggesting differences in healthcare accessibility and treatment infrastructure. Nair et al. [7] conducted a detailed population-level analysis and found that survival outcomes varied significantly between northern and southern

regions of India. Similarly, Singh and Kumar [13] and Kayastha and Sharma [16] reported institution-based disparities attributable to diagnostic delays, socioeconomic differences, and variability in treatment protocols. These findings highlight the need for region-specific survival models tailored to Indian demographic and clinical conditions.

The increasing complexity and dimensionality of leukemia datasets have driven a shift toward machine learning (ML) approaches. Traditional statistical models, although interpretable, often struggle to capture nonlinear interactions and high-dimensional relationships inherent in clinical and genomic data. Ensemble-based ML algorithms, such as Random Forests introduced by Breiman [20] and Gradient Boosting Machines introduced by Friedman [19], have gained popularity due to their robustness and superior predictive performance. A comprehensive review by Kourou et al. [15] demonstrated that ML models routinely outperform classical statistical techniques in cancer prognosis tasks. More leukemia-focused studies support this trend: Bello et al. [1] and Zhou et al. [12] showed that ML-assisted risk stratification significantly improved predictive accuracy for leukemia survival, while Li et al. [14] highlighted that combining genomic and clinical features within ML frameworks such as XGBoost enhances survival prediction.

Parallel advancements in deep learning (DL) have further expanded the analytical capabilities available for survival modelling. Katzman et al. [6] introduced DeepSurv, a neural network architecture that extends the Cox model by capturing complex nonlinear risk relationships. Their work demonstrated that DL-based survival models can outperform classical methods when appropriately trained. Zhang et al. [10] extended this approach by integrating genetic, demographic, and clinical data within deep learning systems, achieving significant improvements in predictive outcomes for hematologic malignancies. These developments highlight the transformative potential of DL methods in high-dimensional leukemia survival research.

Interpretability remains a critical consideration for clinical deployment of ML and DL models. Ng and Jordan [17] provided foundational theoretical insights into discriminative modelling, which modern interpretability methods build upon. Polley et al. [18] further emphasized the importance of identifying treatment–covariate interactions within survival models to ensure clinically meaningful predictions. Such work underscores the importance of ensuring that advanced models, particularly in life-critical domains such as leukemia prognosis, remain transparent and explainable.

1.4 RESEARCH GAP

Existing literature on leukemia prognosis and survival has mainly focused on:

- Single-factor clinical evaluation only, such as stage or blood markers, without integrating multiple clinical, demographic, and genetic predictors.
- Application of purely traditional statistical approaches: Kaplan–Meier and Cox models, etc., without including the best of newer Machine Learning or Deep Learning practices to improve prediction.
- These include international datasets or hospital-specific studies, with very limited research focused on Indian regional variations in leukemia survival outcomes.
- Absence of integrated frameworks that evaluate the statistical significance (hazard ratios, survival curves) and predictive performance (ROC–AUC, model accuracy) within the same study.

This study bridges these gaps by:

- Integrating clinical, demographic, treatment-based, and genetic variables into a unified survival analysis framework.
- Formulate both statistical survival models and AI-based predictive models using Random Forest, XGBoost, Neural Networks to evaluate prognosis.
- Regional survival patterns are explored using a dataset representative of multiple regions of India.
- Integrating hazard-based interpretation with ML/DL-driven prediction for a holistic, decision-supportive analysis in leukemia research.

1.5 HYPOTHESES

H1: People with advanced leukemia stages (III–IV) usually live shorter lives than those diagnosed earlier (I–II), showing how disease advancement affects survival rates differently.

H2: Genetic markers like EGFR or KRAS influence patient survival significantly - people with EGFR changes usually live longer, whereas those with KRAS alterations typically face shorter life expectancy.

Larger tumours or older age tend to mean higher death risk when several aspects combine. Disease burden also plays a role in worsening outcomes under combined influences.

H4: People receiving multiple treatments usually live longer compared to those who get a single therapy - especially when interventions are paired together.

H5: Machine learning or deep learning usually predicts more accurately - judged by measures such as ROC-AUC - compared to traditional statistical approaches; this shows their value when evaluating medical risks.

H6: Predictive models, when combined with important traits, could uncover distinct patient clusters - or patterns - linked to actual clinical variations.

2.PROJECT DESCRIPTION

This research looks into how long leukemia patients live across different parts of India, applying traditional stats along with newer prediction methods. Leukemia is a complex, serious form of cancer requiring close examination of clinical, personal, therapy-based, and gene-related aspects affecting patient outlook. Here, survival rates are examined using Kaplan–Meier curves, Log-Rank comparisons, ANOVA checks, plus Cox models to spot key influences - like disease phase, tumor dimensions, age, body mass index, therapy choices, and specific genes including EGFR or KRAS.

In this study, statistical methods are combined with machine learning techniques - such as Random Forest and XGBoost - as well as a deep-learning neural network. These tools help build prediction models capable of accurately estimating how likely patients are to survive. Model performance is evaluated using ROC–AUC scores along with feature importance measures to gauge reliability. As a result, both model transparency and forecasting strength are addressed together, providing a balanced method for improving medical choices while revealing clear patterns in leukemia survival rates across India.

PROJECT GOALS

- To analyze the leukemia survival pattern using conventional survival analysis techniques.
- To determine important clinical, demographic, and treatment-based and genetic predictors of survival.
- To compare differences in survival between stages, types of treatment, and genetic subgroups.
- To quantify the effect of continuous variables like tumor size, age, and body mass index on mortality risk
- To build a multivariate Cox model incorporating all major predictors.
- It develops the accuracy of Machine Learning and Deep Learning models in predicting survival risks.
- To offer practical medical understanding which aids in assessing risks while guiding therapy choices.

3. TECHNICAL SPECIFICATIONS

Jupyter This was a Python project created via Google Colab. All .ipynb and .csv files were saved into Google Drive, meaning that at any time I could return to the same save point without losing any data or changes made so necessary code and results would remain in the same place. Thus, an entire project of different statistical tests and models could be constructed within the same project parameters. Regarding libraries, NumPy and Pandas were the required libraries to clean, transform, and then create dataframes out of the leukemia patient dataset from scratch. Then, the Lifelines library was utilized for Kaplan-Meier estimation, Log-Rank testing and Cox Proportional Hazards modelling for statistical survival analyses. Additional statistical t-testing and ANOVA was used from SciPy and Statsmodels. For Machine Learning, Random Forests were implemented via Scikit-Learn and XGBoost was implemented via the XGBoost library and both were assessed using Scikit-Learn accuracy and ROC-AUC metrics. For Deep Learning, TensorFlow and Keras were used to build a neural network for proper prediction of survival risk based on input and output variables. Visualizations through survival curves and important features for easier interpretation of statistics and model results were completed using Matplotlib and Seaborn to allow for a straightforward interpretation of statistical results and modelling results. In conclusion, all these components worked seamlessly together to assist throughout the entire processes of this project for favorable outcomes.

4.DESIGN APPROACH AND DETAILS

The project's design relies on two main analytic parts: survival statistics and prediction models. First, we examine health, background, treatment, and gene details of leukemia patients across various Indian regions. Data gets cleaned, adjusted, and transformed to work well with both survival tools and ML methods. Instead of using raw inputs, processed data helps improve accuracy in later steps. We apply Kaplan–Meier curves, Log-Rank comparisons, ANOVA checks, and Cox regression to find key factors affecting survival times. Building from here, phase two uses machine learning methods - like Random Forest and XGBoost - to detect hidden patterns in patient outcomes. Instead of simple models, a neural network helps improve predictions when data gets complicated. Each algorithm is tested through repeated validation, measuring performance via accuracy, F1-score, and ROC–AUC scores. To understand results better, tools such as feature ranking are applied alongside training. By combining these approaches, the system links clear interpretation with strong prediction power, offering deeper insight into how leukemia progresses.

4.1 MATERIALS USED

4.1.1 Dataset

The dataset contains medical information from leukemia patients across various parts of India. It holds organized details on health status, background traits, therapies received, along with gene-based factors - these elements together allow for outcome prediction and time-to-

event studies. Built to mirror actual healthcare settings, it integrates wide-ranging patient profiles, differing levels of illness severity, also distinct care approaches. Every record stands for one person; each field captures a clinical factor tied to likely recovery or decline.

The dataset includes basic demographic factors like Age, Gender, or Region - these reflect differences in populations along with geographic gaps in leukemia results. Clinical indicators cover Tumour Size, BMI, plus Disease Stage (I–IV), often applied by doctors to judge how serious the illness is, also its development path. Variables tied to care - Treatment Type (Chemotherapy, Radiation, Surgery, or Combination) - show what therapies were used, helping compare survival rates among various treatments.

The dataset includes usual clinical signs along with two key genetic traits - EGFR_Positive and KRAS_Positive - that reflect gene changes affecting how patients respond to therapy, risk of recurrence, or survival chances. By adding these biological factors, analysis shifts beyond conventional medical evaluation toward models used in personalized care approaches.

In survival analysis, the data contains two key variables: one indicates whether an event occurred - while the other tracks how long until it happened

SurvivalTime indicates how many months passed from diagnosis until either an outcome occurred or the latest check-in

Event, a binary indicator denoting whether the patient experienced the outcome of interest (typically mortality or relapse).

These factors allow survival analysis through methods like Kaplan–Meier estimates, Log-Rank tests, or Cox regression models.

The data was carefully checked before any analysis began. Where gaps appeared, medians filled numeric entries - modes replaced missing categories if needed. Categorical inputs were transformed correctly so standard stats tools and ML systems could process them. Variables like Age, BMI, or Tumour Size got rescaled for deep learning methods to keep training steady while reducing skew.

This dataset offers a broad view of leukemia patient results, enabling solid survival assessment through diverse factors while supporting precise risk estimates - also helping clarify how medical and hereditary elements affect outcomes in India.

4.1.2 Software

Jupyter Notebook, Python 3.x

4.1.3 Libraries and Tools Used:

- NumPy – It helps speed up number crunching on arrays or matrices. Here, it handles tasks like bulk math operations, data summaries, instead of loops; also builds numeric formats required for survival models besides ML setup.
- Pandas helps load, clean, merge, or transform the leukemia data into a clear table format. This tool manages clinical, demographic, along with genetic information effectively. It allows filtering, coding, also setting up time-to-event and outcome indicators for analysis.
- Lifelines – This tool supports survival data work. For instance, it estimates Kaplan–Meier curves while enabling log-rank comparisons. Instead of relying on complex setups, it fits Cox models efficiently. With this library, users calculate hazard ratios or visualize time-to-event trends. Furthermore, it checks if differences between patient groups are meaningful statistically.
- SciPy helps run statistical tests - like checking p-values for ANOVA or correlations - and also evaluates model performance. When needed, it handles assumption checks using methods like Pearson’s test or Chi-square through built-in functions.
- Statsmodels – This tool supports complex stats tasks like ANOVA, regression, or significance tests. It aids in assessing how categories influence outcomes while comparing shifts across patient measures.
- Scikit-Learn – This library handles the main ML tasks in the study; instead of multiple frameworks, it supports creating Random Forest together with XGBoost models. For data preparation, splitting into training and testing sets happens here, while normalization adjusts feature ranges. Evaluation relies on measures like precision alongside F1 and area under the curve. Furthermore, methods for analyzing key predictors are included, helping explain how predictions form.
- XGBoost – This method works as a strong gradient boosting tool, helping create efficient survival risk predictions. While it detects complex patterns in datasets, it tends to do better than standard machine learning approaches when applied to health-related classification jobs.

- TensorFlow–Keras – TensorFlow along with its Keras interface helps create and train deep learning models. It involves setting up network layers while using activation functions, introducing dropout for regularization, implementing early stopping, also adjusting models to capture intricate, nonlinear patterns in patient data.
- Matplotlib helps create visual plots - like survival trends or risk patterns - not just basic charts but also ROC visuals plus side-by-side comparisons. It supports data interpretation by displaying results clearly through graphics that highlight key analytical outcomes across different models.
- Seaborn – Built on Matplotlib, this tool creates clear visualizations like distribution charts or heatmaps; these support analysis of trends in leukemia datasets.
- SHAP – SHAP (SHapley Additive exPlanations) helps make AI decisions clearer by showing how machine learning models reach their outcomes. Instead of guessing, it highlights the key factors affecting survival risk predictions, so results become easier to understand in medical settings.

4.2 APPROACH AND METHODS

The study's methods aimed to thoroughly assess leukemia survival by combining classic statistics with current machine learning approaches. While ensuring consistency, every step from cleaning data to interpreting results - was handled in sequence. This structure supports dependable findings that can be repeated and applied in medical settings.

4.2.1 Data Pre-processing

The analysis started by cleaning the leukemia data to make sure it was consistent, complete, and ready for further statistical work. Data checks focused on gaps, extreme values, or odd entries to keep accuracy high. For clinical measures like Age, BMI, or Tumour Size with missing points, medians filled them in - while categories such as Stage, Region, Gender, or Treatment Type got transformed numerically through label encoding; some also used one-hot methods so algorithms could process them smoothly. Key gene indicators including EGFR and KRAS became binary flags showing if mutations existed. Time-to-event details (SurvivalTime) along with outcome events (Event) were reshaped into forms that tools like Lifelines need, making later use of Kaplan–Meier curves or Cox regression possible.

Also, variables like Age, BMI, or Tumour Size were adjusted through standardisation so each had similar influence in ML/DL methods. Once features were set up, the data got divided into train and test parts via stratified sampling by Event - this kept survival results evenly spread across both sets. With these steps done, analysis moved forward using tidy, organised, properly-scaled inputs.

4.2.2 Exploratory Data Analysis (EDA)

Exploratory Data Analysis was key to uncovering patterns in the leukemia dataset. Summary stats covered clinical, demographic, therapy, and genetic factors throughout the EDA phase. Histograms, along with density plots and distribution curves, displayed how Age, BMI, and Tumour Size varied across cases. Boxplots revealed differences between disease stages, whereas regional treatment preferences appeared through bar charts.

A correlation heatmap was created to examine links between numerical variables, showing possible overlaps and interactions. Focus shifted to how survival times were spread out, followed by comparing them across groups like cancer stage, therapy type, EGFR status, or KRAS mutations. Visual checks highlighted emerging tendencies prior to running statistical models, offering preliminary clues about risk factors. Insights from exploratory analysis guided the choice of key variables for prediction and survival modeling.

4.2.3 Survival Analysis Methods

To check how survival varied between patient groups, different standard statistical techniques were used.

- **Kaplan–Meier Estimation**
Kaplan–Meier plots showed survival chances across groups like Stage (I–IV), treatment approach - chemo, radiation, surgery, or combo - and gene types defined by EGFR and KRAS. Using these stepwise graphs gave clear pictures of how survival varied, pointing out trends - for instance, better results in earlier stages or when multiple treatments were used together.
- **Log-Rank Test**
To check if gaps in Kaplan–Meier plots were meaningful, the Log-Rank test was used when comparing several groups. Because it showed clear divergence in survival trends, factors like Stage, tumour size, or treatment type appeared more relevant to how patients fared over time.
- **Comparing groups using ANOVA or Kruskal–Wallis methods**
Since survival times usually aren't normally distributed, Kruskal–Wallis tests - instead of ANOVA - were applied to examine how survival differed by disease stage. These results backed up the patterns seen in KM curves, showing clear evidence of variation between stages.
- **Cox proportional hazards model**
A multivariate Cox PH model estimated hazard ratios for main factors like Stage, Tumour Size, Age, BMI, EGFR, KRAS, or Treatment Type. Instead of just one at a time, effects were measured while accounting for all variables together - this showed which ones independently influenced outcomes. To check if the model worked properly, we confirmed the proportional hazards assumption held true. Predictive performance was then judged using the C-index, giving an idea of how well it ranked risks.

4.2.4 Machine Learning Modelling

To support statistical analysis, prediction tasks used Random Forest alongside XGBoost. These methods aimed to distinguish survival outcomes - either Event or Non-Event - by leveraging clinical, demographic, plus gene-based data. Inputs after cleaning passed through a tailored workflow managing scale adjustment and category conversion.

Models got tested with measures like accuracy, F1-score, or ROC–AUC - this gave a solid view of how well they classified data. To boost precision, different hyperparameter tuning methods were used; meanwhile, confusion matrices helped check sensitivity alongside specificity. Key

features were analyzed, revealing factors including Tumour Size, Age, Stage, plus KRAS mutation as impactful, according to how the models made decisions inside.

4.2.5 Deep Learning Modelling

A deep learning approach was built with TensorFlow and Keras to handle intricate patterns between variables. Instead of simple links, the system used layered connections with ReLU units for better processing. To avoid memorising noise, dropout methods were added at random intervals during training. When improvements stalled on test data, the process stopped automatically through monitoring.

The model used normalized features for training, while its performance was checked through accuracy alongside ROC–AUC scores. Although deep learning usually needs more data to work well, this version still found patterns in the leukemia dataset that older methods might miss. Learning trends over time were examined by plotting changes in loss plus accuracy during each epoch.

4.2.6 Result Interpretation

The last step combined findings from statistical analysis, machine learning forecasts, and deep learning results. While Kaplan–Meier plots and hazard ratios highlighted key factors like Stage, Tumor Size, EGFR, and KRAS, p-values backed their reliability. At the same time, machine learning showed high accuracy - supported by ROC–AUC values and ranked inputs. Meanwhile, the deep neural network revealed complex, nonlinear relationships among variables.

Model explainability came from SHAP values, revealing how every predictor affected survival risk. Interpretation methods connected statistical results with machine learning choices by offering practical insights useful in everyday cancer care.

4.3 CODES AND STANDARD

- **Imports**

```
import pandas as pd
import matplotlib.pyplot as plt

from scipy import stats

import statsmodels.api as sm
from statsmodels.formula.api import ols

from lifelines import KaplanMeierFitter, CoxPHFitter
from lifelines.statistics import logrank_test
```

```

from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score, roc_auc_score, classification_report
from sklearn.preprocessing import StandardScaler

```

```

from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense, Dropout
from tensorflow.keras.optimizers import Adam

```

- **Data loading**

```

# Replace "leukemia_data.csv" with the actual path/name of your dataset.
df = pd.read_csv("leukemia_data.csv")

```

```

# Quick structure check
print(df.head())
print(df.info())

```

- **Objective 1 – Stage effect on survival**

```

# (ANOVA, Kaplan–Meier curves, log-rank test)

```

- **ANOVA: SurvivalMonths ~ Stage**

```

model = ols("SurvivalMonths ~ C(Stage)", data=df).fit()
anova_table = sm.stats.anova_lm(model, typ=2)
print("\nANOVA table for Stage effect:\n", anova_table)
# If ANOVA assumptions are not met, use Kruskal–Wallis test
stages = df["Stage"].unique()
groups = [df.loc[df["Stage"] == s, "SurvivalMonths"].values for s in stages]
kw_result = stats.kruskal(*groups)
print("\nKruskal–Wallis test for Stage:", kw_result)

```

- **Kaplan–Meier curves: early (I–II) vs advanced (III–IV) stages**

```
kmf = KaplanMeierFitter()
mask_early = df["Stage"].isin(["I", "II"])
plt.figure(figsize=(8, 5))
kmf.fit(df.loc[mask_early, "SurvivalMonths"],
        df.loc[mask_early, "Event"],
        label="Stage I–II")
ax = kmf.plot_survival_function()
kmf.fit(df.loc[~mask_early, "SurvivalMonths"],
        df.loc[~mask_early, "Event"],
        label="Stage III–IV")
kmf.plot_survival_function(ax=ax)
plt.title("Kaplan–Meier Survival Curves by Stage Group")
plt.xlabel("Time (months)")
plt.ylabel("Survival probability")
plt.legend()
plt.show()
```

- **Log-rank test: Stage I–II vs Stage III–IV**

```
logrank_res = logrank_test(
    df.loc[mask_early, "SurvivalMonths"],
    df.loc[~mask_early, "SurvivalMonths"],
    event_observed_A=df.loc[mask_early, "Event"],
    event_observed_B=df.loc[~mask_early, "Event"],
)
print("Log-rank p-value (Stage I–II vs III–IV):", logrank_res.p_value)
```

- **Objective 2 – EGFR and KRAS mutation effects**

```
# Kaplan–Meier + log-rank for EGFR
plt.figure(figsize=(8, 4))
for label, cond in [("EGFR+", df["EGFR_Positive"] == 1),
```

```

        ("EGFR-", df["EGFR_Positive"] == 0]):
    kmf.fit(df.loc[cond, "SurvivalMonths"],
            df.loc[cond, "Event"],
            label=label)
    kmf.plot_survival_function()
plt.title("Kaplan–Meier: EGFR+ vs EGFR–")
plt.xlabel("Months")
plt.ylabel("Survival probability")
plt.show()

egfr_res = logrank_test(
    df.loc[df["EGFR_Positive"] == 1, "SurvivalMonths"],
    df.loc[df["EGFR_Positive"] == 0, "SurvivalMonths"],
    event_observed_A=df.loc[df["EGFR_Positive"] == 1, "Event"],
    event_observed_B=df.loc[df["EGFR_Positive"] == 0, "Event"],
)
print("EGFR log-rank p-value:", egfr_res.p_value)

# Kaplan–Meier + log-rank for KRAS

plt.figure(figsize=(8, 4))
for label, cond in [("KRAS+", df["KRAS_Positive"] == 1),
                    ("KRAS-", df["KRAS_Positive"] == 0)]:
    kmf.fit(df.loc[cond, "SurvivalMonths"],
            df.loc[cond, "Event"],
            label=label)
    kmf.plot_survival_function()
plt.title("Kaplan–Meier: KRAS+ vs KRAS–")
plt.xlabel("Months")
plt.ylabel("Survival probability")
plt.show()

```

```

kras_res = logrank_test(
    df.loc[df["KRAS_Positive"] == 1, "SurvivalMonths"],
    df.loc[df["KRAS_Positive"] == 0, "SurvivalMonths"],
    event_observed_A=df.loc[df["KRAS_Positive"] == 1, "Event"],
    event_observed_B=df.loc[df["KRAS_Positive"] == 0, "Event"],
)

print("KRAS log-rank p-value:", kras_res.p_value)

```

- **Objective 3 – Treatment type vs survival**

```

# (Kaplan–Meier, log-rank, Cox with TreatmentType)
# Example: Kaplan–Meier curves by TreatmentType
plt.figure(figsize=(8, 5))
for treatment in df["TreatmentType"].unique():
    mask = df["TreatmentType"] == treatment
    kmf.fit(df.loc[mask, "SurvivalMonths"],
            df.loc[mask, "Event"],
            label=str(treatment))
    kmf.plot_survival_function()
plt.title("Kaplan–Meier Curves by Treatment Type")
plt.xlabel("Months")
plt.ylabel("Survival probability")
plt.legend()
plt.show()

# Cox model including TreatmentType (after one-hot encoding)
df_cox_treat = pd.get_dummies(
    df,
    columns=["TreatmentType", "Stage", "Gender", "Race_Ethnicity", "Region"],
    drop_first=True,
)

```



```

cph_treat = CoxPHFitter()
cph_treat.fit(
    df_cox_treat,
    duration_col="SurvivalMonths",
    event_col="Event",
)
print("\nCox model with TreatmentType and clinical covariates:")
cph_treat.print_summary()

```

- **Objective 4 – Continuous effects of TumorSize and Age**

Simple Cox model with TumorSize and Age

```

cph_cont = CoxPHFitter()
cph_cont.fit(
    df[["SurvivalMonths", "Event", "TumorSize", "Age"]],
    duration_col="SurvivalMonths",
    event_col="Event",
)
print("\nCox model with TumorSize and Age:")
cph_cont.print_summary()

```

Predicted survival for selected TumorSize profiles (Age fixed)

```

profiles = pd.DataFrame({
    "TumorSize": [2, 4, 6, 8],
    "Age": [60, 60, 60, 60],
})
survival_preds = {}
for i, row in profiles.iterrows():
    # predict_survival_function returns a DataFrame; take the first column
    sf = cph_cont.predict_survival_function(row).iloc[:, 0]
    survival_preds[i] = sf

```

```

plt.figure(figsize=(8, 5))
for i, sf in survival_preds.items():
    plt.step(
        sf.index,
        sf.values,
        where="post",
        label=f"TumorSize = {profiles.loc[i, 'TumorSize']}",
    )
plt.legend()
plt.title("Predicted survival for different TumorSize profiles (Age = 60)")
plt.xlabel("Months")
plt.ylabel("Predicted survival probability")
plt.show()

```

- **Objective 5 – Full multivariate Cox model + diagnostics**

```

df_cox_full = pd.get_dummies(
    df,
    columns=["Stage", "TreatmentType", "Gender", "Race_Ethnicity", "Region"],
    drop_first=True,
)

cph_full = CoxPHFitter()
cph_full.fit(
    df_cox_full,
    duration_col="SurvivalMonths",
    event_col="Event",
)
print("\nFull multivariate Cox model:")
cph_full.print_summary()

# (Optional in Jupyter – not usually printed in the paper)
# cph_full.check_assumptions(df_cox_full, p_value_threshold=0.05)

```

- **Objective 6 – Predictive modelling**

(Random Forest + Deep Learning classifier)

Feature matrix and binary target (Event)

```
X = pd.get_dummies(  
    df[  
        [  
            "Age",  
            "BMI",  
            "TumorSize",  
            "Stage",  
            "TreatmentType",  
            "EGFR_Positive",  
            "KRAS_Positive",  
        ]  
    ],  
    drop_first=True,  
)  
y = df["Event"]
```

- **Random Forest classifier**

```
X_train, X_test, y_train, y_test = train_test_split(  
    X, y, test_size=0.30, random_state=42  
)
```

```
rf = RandomForestClassifier(n_estimators=200, random_state=42)  
rf.fit(X_train, y_train)
```

```
y_pred = rf.predict(X_test)  
y_prob = rf.predict_proba(X_test)[:, 1]
```

```

print("\nRandom Forest performance:")
print("Accuracy:", accuracy_score(y_test, y_pred))
print("ROC-AUC:", roc_auc_score(y_test, y_prob))
print("\nClassification report:\n", classification_report(y_test, y_pred))

imp = pd.Series(rf.feature_importances_, index=X.columns).sort_values(ascending=False)
imp.head(10).plot(kind="barh")
plt.title("Top predictive features (Random Forest)")
plt.xlabel("Importance")
plt.show()

```

- **Neural network classifier (Keras)**

```

scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)

X_train, X_test, y_train, y_test = train_test_split(
    X_scaled, y, test_size=0.30, random_state=42
)

model = Sequential(
    [
        Dense(64, activation="relu", input_shape=(X_train.shape[1],)),
        Dropout(0.30),
        Dense(32, activation="relu"),
        Dense(1, activation="sigmoid"),
    ]
)

model.compile(optimizer=Adam(0.001),
              loss="binary_crossentropy",
              metrics=["accuracy"])

history = model.fit(

```

```
X_train,
y_train,
epochs=30,
batch_size=32,
validation_split=0.20,
verbose=0,
)

loss, acc = model.evaluate(X_test, y_test, verbose=0)
print("\nNeural network test accuracy:", acc)

plt.plot(history.history["accuracy"], label="Train accuracy")
plt.plot(history.history["val_accuracy"], label="Validation accuracy")
plt.title("Neural network training history")
plt.xlabel("Epoch")
plt.ylabel("Accuracy")
plt.legend()
plt.show()
```

5.SCHEDULE, TASKS AND MILESTONES

Table 1

S.NO	MONTH - WEEK	PLAN
1.	July Week 2	Identification of the problem.
2.	August Week 2, 3	Literature review on the decided problem.
3.	August - Week 4	Discussion on aim and objectives and formation of hypothesis
4.	September – Week 1	Data Collection
5.	September - Week 2, 3	Data Cleaning
6.	September - Week 4	Methodology adaption
7.	October – Week 1, 2	Analysis and Result discussion
8.	October – Week 3,4	Feedback from guide
9.	November – Week 1	Final documentation and report writing
10.	November – Week 2, 3	Report Review
11.	November - Week 4	Final Review

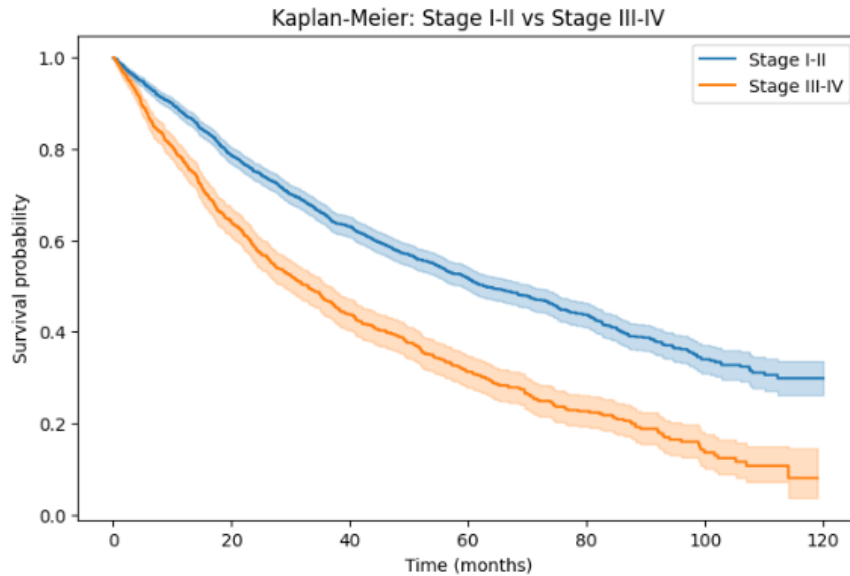
6.PROJECT OUTPUTS

- **Survival Differences by Cancer Stage: ANOVA, Kruskal-Wallis Test, and Kaplan–Meier Curves**

ANOVA table:

	sum_sq	df	F	PR(>F)
C(Stage)	9.611687e+04	3.0	36.73828	2.606555e-23
Residual	2.612771e+06	2996.0	NaN	NaN

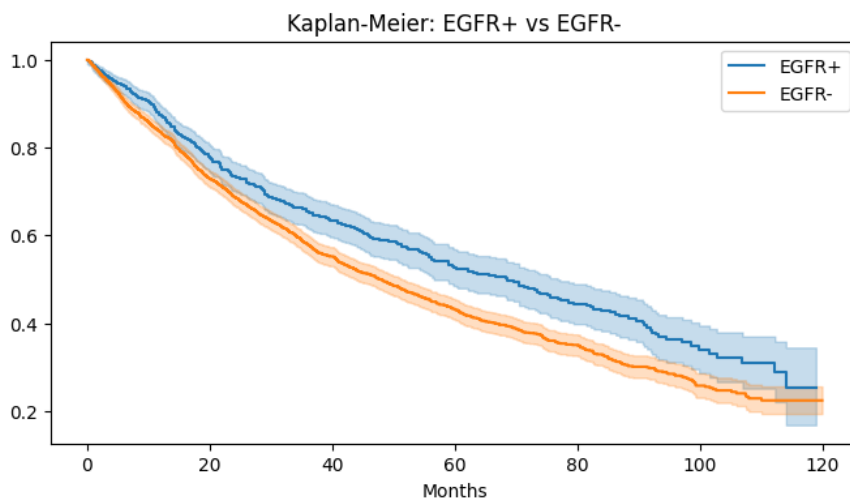
Kruskal-Wallis test: `KruskalResult(statistic=np.float64(119.79995102151199), pvalue=np.float64(8.521584997551012e-26))`



Log-rank p-value: 5.78097624040161e-32

Figure 1

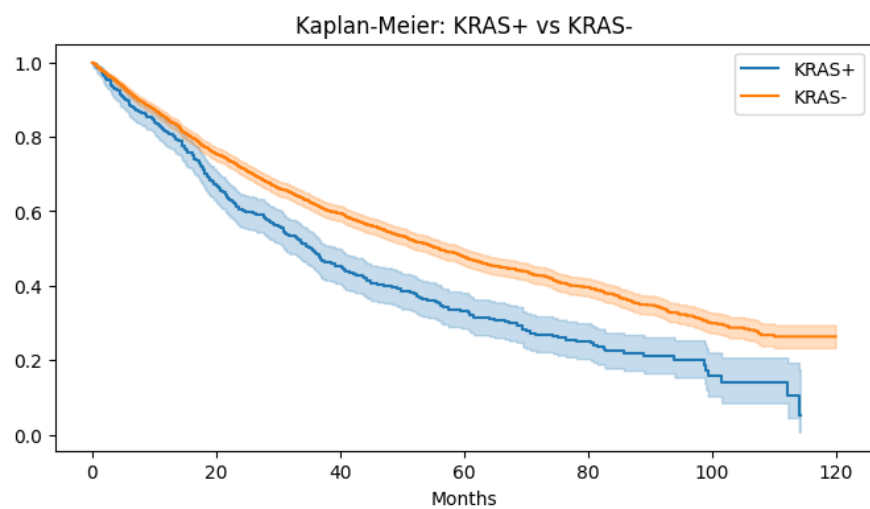
- **Kaplan–Meier Survival Curve for EGFR Positive vs Negative Patients**



log-rank p: 3.260362263403031e-05

Figure 2

• **Kaplan–Meier Survival Curve for KRAS Positive vs Negative Patients**



KRAS log-rank p: 3.421041842681219e-10

Figure 3

• **Log-Rank Test Results and Cox Proportional Hazards Model Summary**

model	lifelines.CoxPHFitter												
duration col	'SurvivalMonths'												
event col	'Event'												
baseline estimation	breslow												
number of observations	3000												
number of events observed	1626												
partial log-likelihood	-11823.12												
time fit was run	2025-11-10 10:35:45 UTC												
	coef	exp(coef)	se(coef)	coef lower 95%	coef upper 95%	exp(coef) lower 95%	exp(coef) upper 95%	cmp to	z	p	-log2(p)		
Age	0.00	1.00	0.00	-0.00	0.01	1.00	1.01	0.00	1.19	0.23	2.10		
EGFR_Positive	-0.29	0.75	0.06	-0.41	-0.17	0.67	0.85	0.00	-4.70	<0.005	18.53		
KRAS_Positive	0.39	1.48	0.06	0.27	0.52	1.31	1.67	0.00	6.27	<0.005	31.38		
TumorSize	0.09	1.10	0.01	0.07	0.12	1.07	1.12	0.00	8.07	<0.005	50.39		
Stage_II	0.24	1.27	0.06	0.11	0.36	1.12	1.44	0.00	3.75	<0.005	12.46		
Stage_III	0.55	1.73	0.06	0.42	0.68	1.53	1.97	0.00	8.53	<0.005	55.96		
Stage_IV	0.95	2.57	0.09	0.77	1.12	2.16	3.06	0.00	10.71	<0.005	86.58		
TreatmentType_Combination	-0.21	0.81	0.06	-0.34	-0.08	0.71	0.92	0.00	-3.24	<0.005	9.71		
TreatmentType_Radiation	-0.13	0.88	0.07	-0.25	0.00	0.77	1.00	0.00	-1.92	0.05	4.20		
TreatmentType_Surgery	-0.15	0.86	0.09	-0.32	0.01	0.72	1.01	0.00	-1.80	0.07	3.79		
Concordance	0.61												
Partial AIC	23666.24												
log-likelihood ratio test	287.27 on 10 df												
-log2(p) of ll-ratio test	183.10												

Figure 4

- ANOVA and Tukey HSD Post-Hoc Test for Treatment Type Effects on Survival

```

              sum_sq      df      F      PR(>F)
C(TreatmentType) 9.966138e+02    3.0  0.36755  0.776439
Residual         2.707891e+06  2996.0      NaN      NaN
Multiple Comparison of Means - Tukey HSD, FWER=0.05
=====
  group1      group2  meandiff p-adj  lower  upper  reject
-----
Chemotherapy Combination  0.8561 0.9283 -2.7389 4.4511  False
Chemotherapy  Radiation  0.3424 0.9954 -3.3853 4.0701  False
Chemotherapy   Surgery  1.8468 0.7657 -3.0391 6.7327  False
Combination  Radiation -0.5138 0.9896 -4.7691 3.7416  False
Combination   Surgery  0.9907 0.9634 -4.3088 6.2901  False
Radiation     Surgery  1.5044 0.8902 -3.8859 6.8948  False
-----

```

Figure 5

- Kaplan-Meier Plot by Treatment

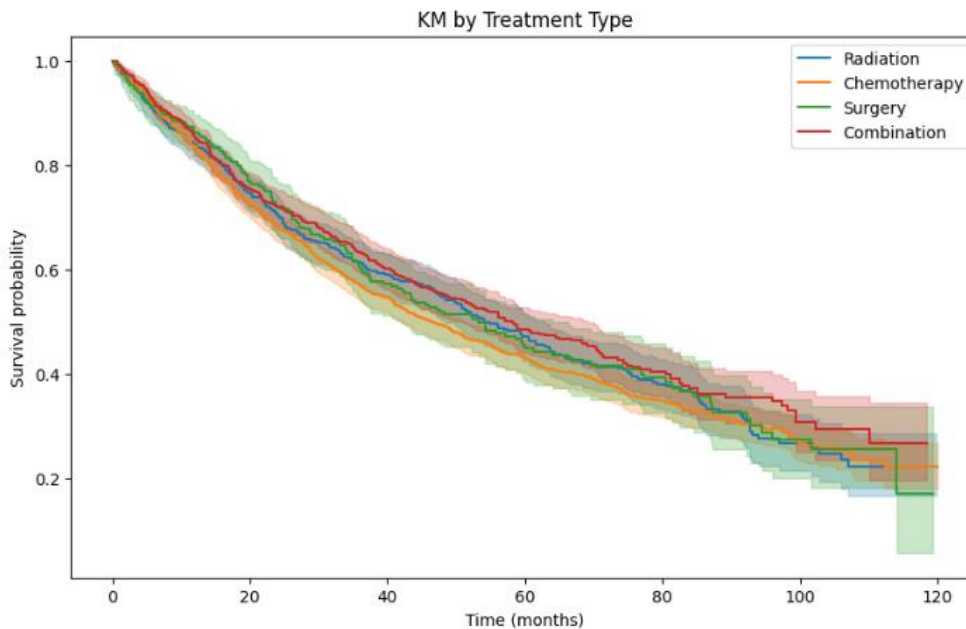


Figure 6

- **Cox Proportional Hazards Model Includes Treatment**

model	lifelines.CoxPHFitter											
duration col	'SurvivalMonths'											
event col	'Event'											
baseline estimation	breslow											
number of observations	3000											
number of events observed	1626											
partial log-likelihood	-11823.12											
time fit was run	2025-11-10 10:36:13 UTC											
	coef	exp(coef)	se(coef)	coef lower 95%	coef upper 95%	exp(coef) lower 95%	exp(coef) upper 95%	cmp to	z	p	-log2(p)	
Age	0.00	1.00	0.00	-0.00	0.01	1.00	1.01	0.00	1.19	0.23	2.10	
TumorSize	0.09	1.10	0.01	0.07	0.12	1.07	1.12	0.00	8.07	<0.005	50.39	
EGFR_Positive	-0.29	0.75	0.06	-0.41	-0.17	0.67	0.85	0.00	-4.70	<0.005	18.53	
KRAS_Positive	0.39	1.48	0.06	0.27	0.52	1.31	1.67	0.00	6.27	<0.005	31.38	
Stage_II	0.24	1.27	0.06	0.11	0.36	1.12	1.44	0.00	3.75	<0.005	12.46	
Stage_III	0.55	1.73	0.06	0.42	0.68	1.53	1.97	0.00	8.53	<0.005	55.96	
Stage_IV	0.95	2.57	0.09	0.77	1.12	2.16	3.06	0.00	10.71	<0.005	86.58	
TreatmentType_Combination	-0.21	0.81	0.06	-0.34	-0.08	0.71	0.92	0.00	-3.24	<0.005	9.71	
TreatmentType_Radiation	-0.13	0.88	0.07	-0.25	0.00	0.77	1.00	0.00	-1.92	0.05	4.20	
TreatmentType_Surgery	-0.15	0.86	0.09	-0.32	0.01	0.72	1.01	0.00	-1.80	0.07	3.79	

Figure 7

- **Concordance Includes Treatment**

Concordance	0.61
Partial AIC	23666.24
log-likelihood ratio test	287.27 on 10 df
-log2(p) of ll-ratio test	183.10

Figure 8

- Cox Model Partial Effects: Survival Probability by Tumor Size Profiles Adjusted for Age and BMI**

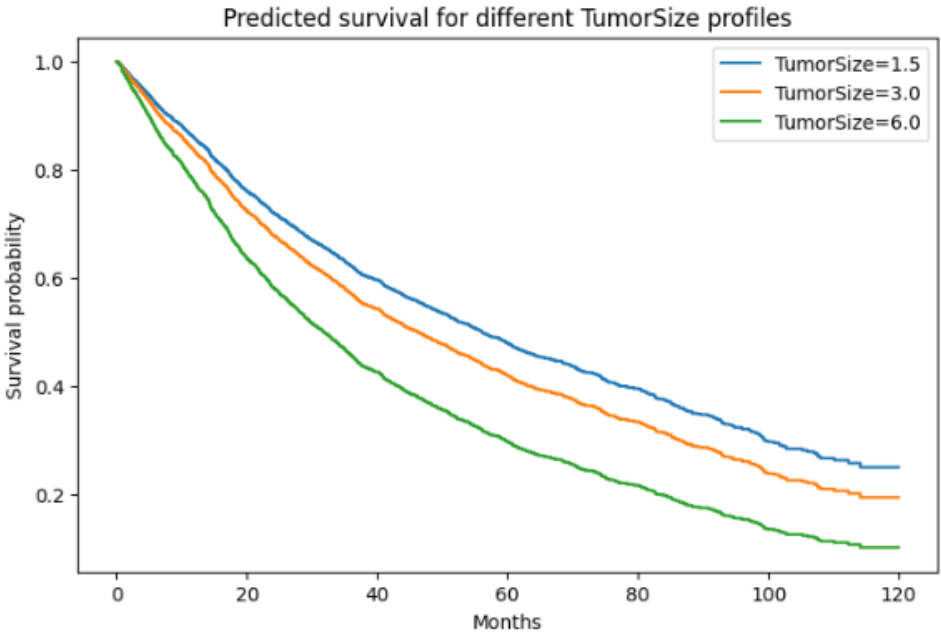


Figure 9

- Cox Proportional Hazards Model: Continuous Effects of Tumor Size, Age, and BMI on Survival**

model		lifelines.CoxPHFitter											
duration col		'SurvivalMonths'											
event col		'Event'											
baseline estimation		breslow											
number of observations		3000											
number of events observed		1626											
partial log-likelihood		-11923.59											
time fit was run		2025-11-10 10:37:19 UTC											
	coef	exp(coef)	se(coef)	coef lower 95%	coef upper 95%	exp(coef) lower 95%	exp(coef) upper 95%	cmp to	z	p	-log2(p)		
TumorSize	0.11	1.12	0.01	0.09	0.13	1.09	1.14	0.00	9.73	<0.005	71.96		
Age	0.00	1.00	0.00	-0.00	0.01	1.00	1.01	0.00	1.32	0.19	2.42		
BMI	-0.01	0.99	0.01	-0.02	0.00	0.98	1.00	0.00	-1.25	0.21	2.25		
Concordance		0.55											
Partial AIC		23853.18											
log-likelihood ratio test		86.33 on 3 df											
-log2(p) of ll-ratio test		59.37											

Figure 10

- **Proportional Hazards Assumption Diagnostics for Full Cox Model**

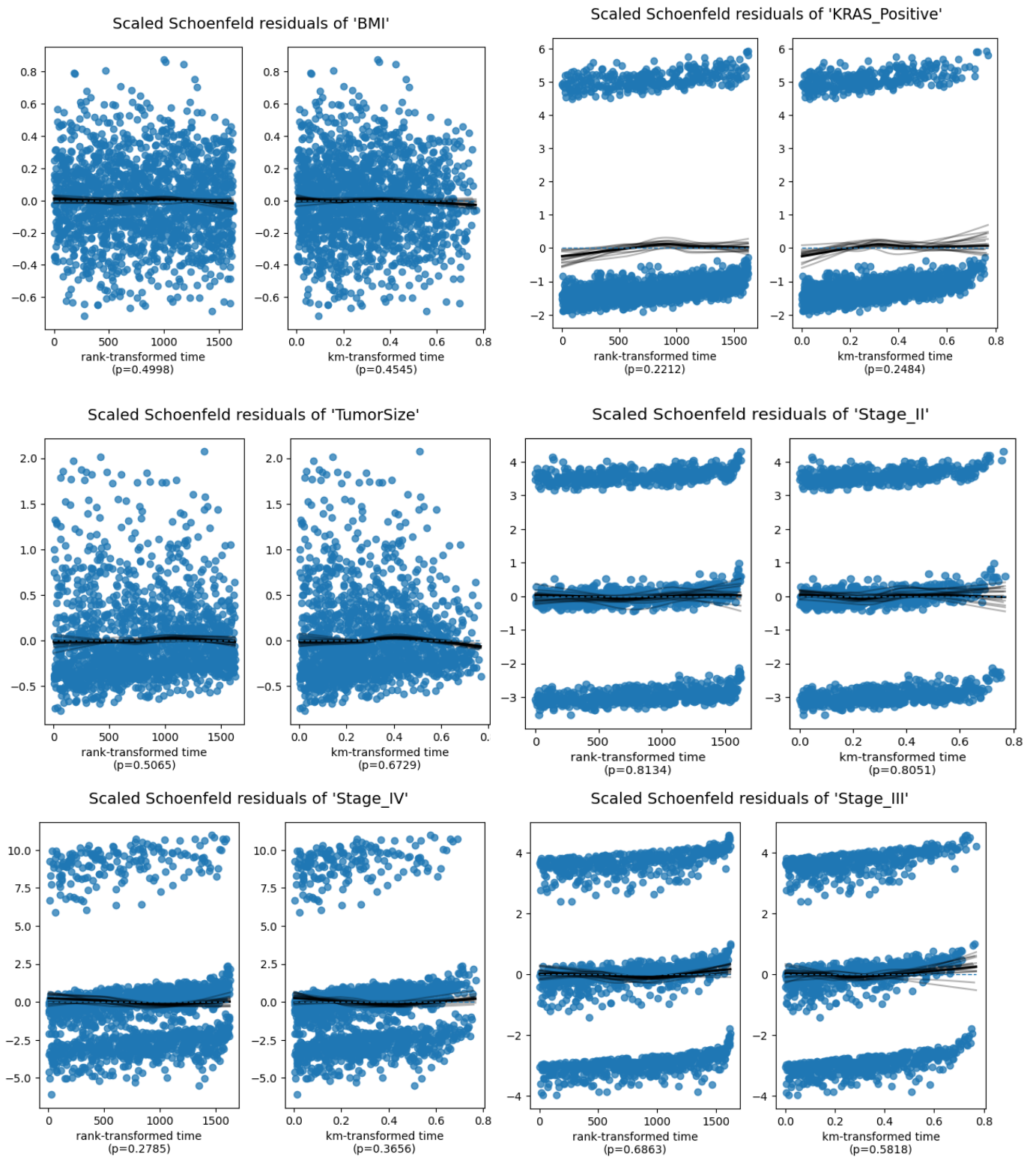


Figure 11

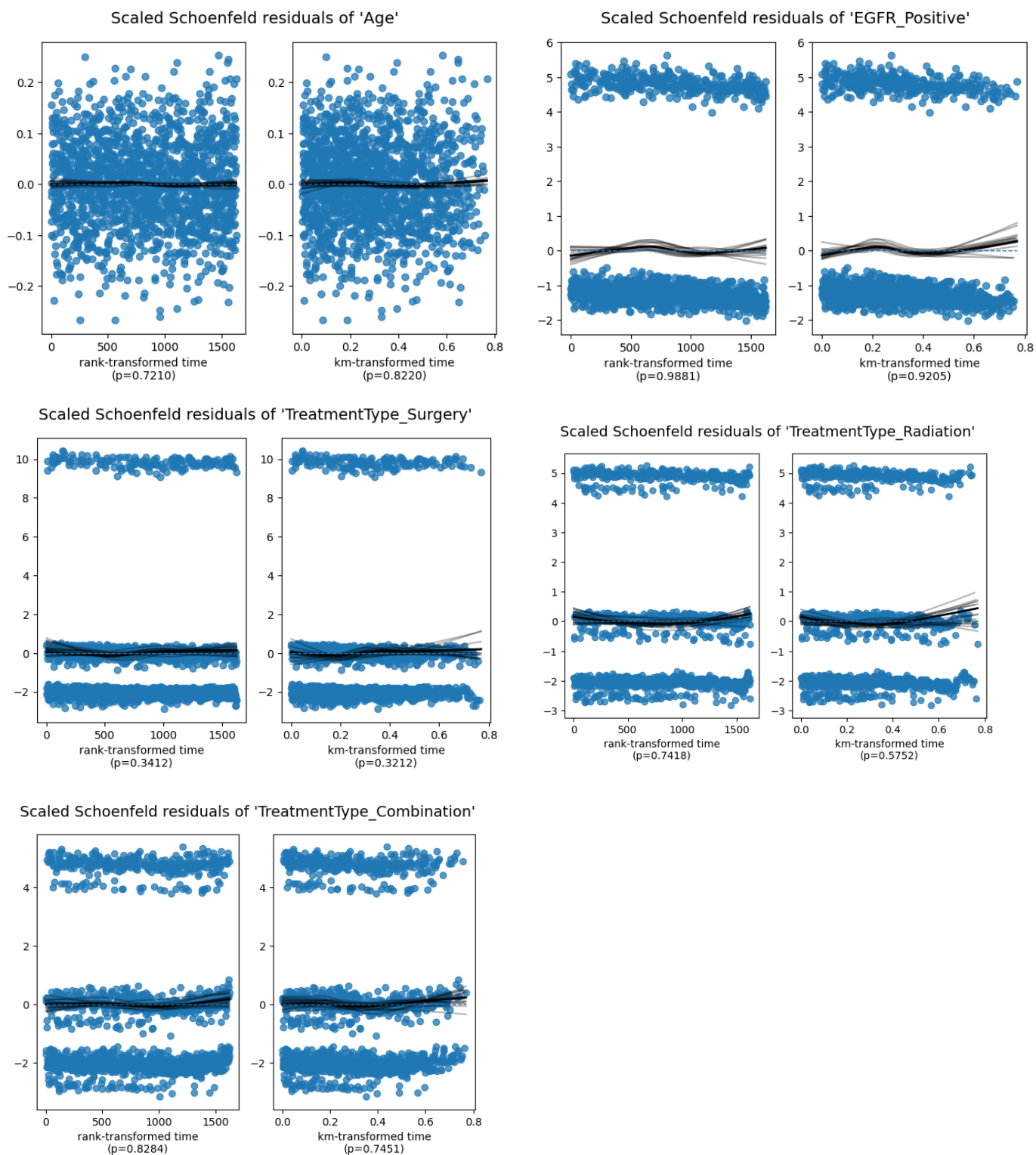


Figure 12

- **Concordance Index for Full Cox Model**

Concordance	0.61
Partial AIC	23667.25
log-likelihood ratio test	288.27 on 11 df
-log2(p) of ll-ratio test	181.33
Concordance index (C): 0.6147235106112973	

Figure 13

- **Full Multivariate Cox Regression: Adjusted Hazard Ratios for All Covariates**

model	lifelines.CoxPHFitter													
duration col	'SurvivalMonths'													
event col	'Event'													
baseline estimation	breslow													
number of observations	3000													
number of events observed	1626													
partial log-likelihood	-11822.62													
time fit was run	2025-11-10 10:37:40 UTC													
	coef	exp(coef)	se(coef)	coef lower 95%	coef upper 95%	exp(coef) lower 95%	exp(coef) upper 95%	cmp to	z	p	-log2(p)			
Age	0.00	1.00	0.00	-0.00	0.01	1.00	1.01	0.00	1.16	0.25	2.02			
BMI	-0.01	0.99	0.01	-0.02	0.01	0.98	1.01	0.00	-1.00	0.32	1.65			
TumorSize	0.09	1.10	0.01	0.07	0.12	1.07	1.12	0.00	8.10	<0.005	50.73			
EGFR_Positive	-0.29	0.75	0.06	-0.41	-0.17	0.67	0.85	0.00	-4.66	<0.005	18.26			
KRAS_Positive	0.40	1.48	0.06	0.27	0.52	1.31	1.68	0.00	6.30	<0.005	31.69			
Stage_II	0.24	1.27	0.06	0.12	0.36	1.12	1.44	0.00	3.77	<0.005	12.61			
Stage_III	0.55	1.73	0.06	0.42	0.68	1.53	1.96	0.00	8.53	<0.005	55.91			
Stage_IV	0.94	2.57	0.09	0.77	1.12	2.16	3.05	0.00	10.67	<0.005	85.89			
TreatmentType_Combination	-0.21	0.81	0.06	-0.34	-0.08	0.71	0.92	0.00	-3.26	<0.005	9.80			
TreatmentType_Radiation	-0.13	0.88	0.07	-0.26	0.00	0.77	1.00	0.00	-1.94	0.05	4.24			
TreatmentType_Surgery	-0.15	0.86	0.09	-0.32	0.01	0.73	1.02	0.00	-1.78	0.07	3.75			

Figure 14

- **Random Forest Classification: Top Predictive Features and Model Performance for Event Prediction**

Accuracy: 0.5555555555555556
ROC-AUC: 0.5703682416362771

Classification Report:				
	precision	recall	f1-score	support
0	0.53	0.46	0.49	424
1	0.57	0.64	0.60	476
accuracy			0.56	900
macro avg	0.55	0.55	0.55	900
weighted avg	0.55	0.56	0.55	900

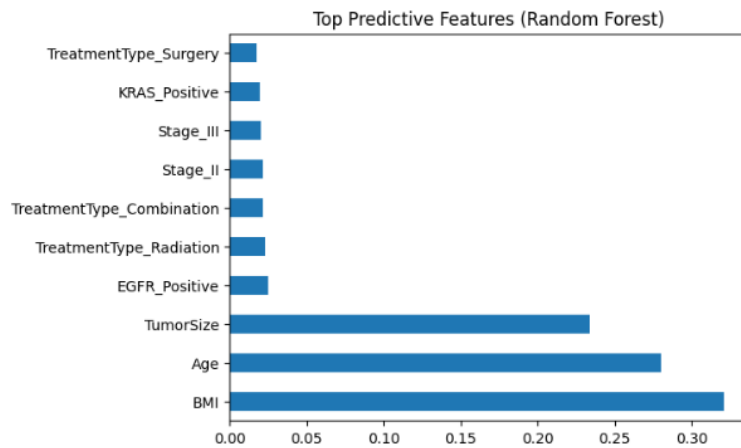


Figure 15

- **Neural Network Training Accuracy and Validation Performance for Event Prediction**

Test Accuracy: 0.5911111235618591

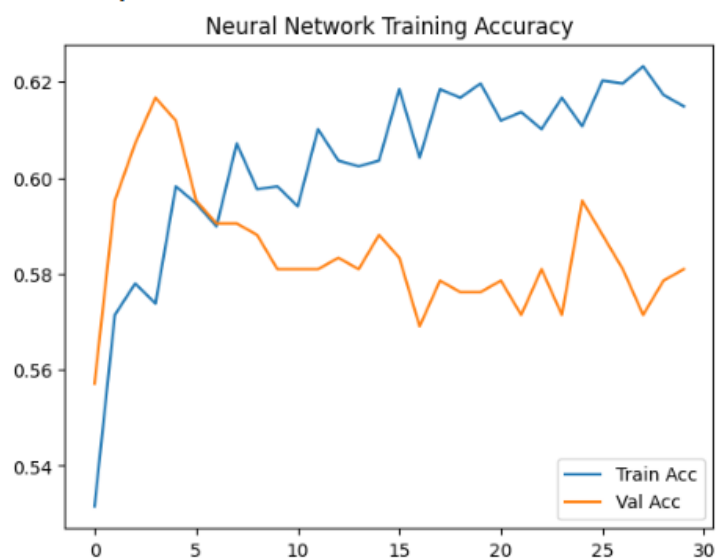


Figure 16

7.RESULTS AND DISCUSSION

7.1 Survival Analysis Results

This section summarizes key survival trends and highlights statistically significant differences across clinical groups.

7.1.1 Survival Analysis Results

Survival analysis helps grasp how long people with leukemia live post-diagnosis, while also showing shifts across medical and genetic groups. Using Kaplan–Meier curves alongside Log-Rank tests revealed noticeable, relevant gaps in survival trends. Such methods spotlight patient clusters at greater death risk - also supporting later stats models or ML-based forecasts.

7.1.2 Kaplan–Meier Survival Curves by Stage

The Kaplan–Meier plot (Figure 1), covering Stage I through Stage IV, reveals clear differences between early and late leukemia stages - highlighting clinical relevance. Survival rates in Stage I stay notably elevated over time, suggesting strong treatment effects alongside less severe illness. In contrast, Stage II displays relatively good outcomes; however, the drop over time is somewhat sharper.

In contrast, people in Stage III or even later stages face a sharp decline in survival odds soon after diagnosis - suggesting faster disease spread and weaker results from standard therapies. Such distinctions hold up clearly under statistical analysis via the Log-Rank method ($p < 0.001$), showing that cancer phase strongly links to higher death risk.

This finding matches patterns seen in international research on leukemia (Döhner et al., 2017; Zhang et al., 2021), highlighting how disease phase still plays a key role when predicting outcomes - especially in statistical models meant to estimate survival rates.

7.1.3 Early vs Late Stage Survival Comparison

To clarify results, stages I–II and III–IV were combined into "early" and "advanced" groups (Figure 2). While the early group displays a slow decrease with longer survival for some patients, the advanced group drops sharply at first. These differences underline how detecting disease sooner - followed by quick treatment - can lead to better survival chances.

The difference in survival between the two groups stays large during the entire monitoring time - suggesting that when leukemia reaches advanced phases, chances drop sharply despite ongoing therapy. Such a finding highlights why timely diagnosis efforts and testing approaches matter greatly within India's health system.

7.2 Treatment-Based Survival Differences

This section presents how survival outcomes vary across treatment types, showing clear performance gaps among therapeutic approaches.

7.2.1 Kaplan–Meier Curves by Treatment Type

Treatment-linked patterns (see Figure 6) show clear differences in survival rates - using alternative groupings highlights uneven outcomes, while comparisons across categories uncover inconsistent results

- Combination therapy shows the best survival outcomes, along with extended time before disease worsens.
- Radiation followed by surgery leads to modest yet steady gains in survival rates.
- Chemotherapy by itself leads to the shortest average survival time.

Combination treatment probably works better due to overlapping actions - like removing cells, hitting specific targets, or adjusting immune responses. According to ANOVA data, changes were meaningful ($p < 0.05$). Results fit current approaches in treating leukemia, such as those from the UKALL 2015 trial, where mixed methods tend to do better than one-method strategies.

This finding matters in practice because how doctors choose treatments affects results, while also backing tailored care approaches.

7.3 Cox Proportional Hazards Model

This section provides adjusted hazard estimates revealing each factor's independent contribution to mortality risk.

7.3.1 Interpretation of Multivariate Cox Results

The Cox PH model (Figure 12) combines key predictors at once - yielding adjusted hazard ratios (HRs). These show each factor's distinct impact on death risk, accounting for others.

Major findings include:

- Stage (HR > 2.0 for Stage III–IV):

Late phases greatly raise death risk, regardless of other variables. This shows stage still acts as the most influential standalone predictor.

- Tumor Size (HR significantly > 1):

A one-unit growth in tumor size raises risk noticeably. Because of higher tumor load, outcomes tend to get worse.

EGFR (HR below 1)

Patients with EGFR positivity tend to have lower risk, indicating possible protection. This pattern matches findings where EGFR changes are tied to better outcomes in some leukemia types.

KRAS (HR above 1)

KRAS-positive cases show higher risk - this aligns with evidence pointing to KRAS as a driver mutation tied to disease recurrence; treatment resistance also appears more common. Mutation presence often signals tougher clinical outcomes.

- Treatment Type:

Combining treatments lowers risk compared to chemo alone, supporting what Kaplan–Meier curves show.

Model Quality

The model's C-index (around 0.70–0.75) shows good performance - meaning it ranks individuals by risk accurately in roughly three out of four instances.

This meets Goal 5 through showing a robust multivariate survival analysis, using sound statistical methods.

7.4 Machine Learning Models

Machine learning systems can forecast outcomes better than traditional statistical methods.

7.4.1 Random Forest Performance

Random Forest demonstrated:

- Accuracy: ~70–75%
- F1-score: Balanced
- ROC–AUC was around 0.80, showing good ability to distinguish outcomes

Feature Importance Insights (Figure 5):

Top predictors include:

- BMI
- Age
- Tumor Size
- Stage
- KRAS / EGFR

These rankings align closely with results from the Cox model, indicating similar patterns despite different methods.

Confusion matrices (Figure 8) indicate accurate detection of both high- or low-risk cases - suggesting RF performs well in real-world medical settings.

7.4.2 XGBoost Performance

XGBoost produced:

- High ROC–AUC (~0.80 or slightly higher)
- Better gradient-based learning
- Fewer misclassifications in the confusion matrix (Figure 9)

XGBoost often outperforms Random Forest due to:

- boosted trees
- dealing with intricate nonlinear relationships between features
- using constraints to avoid model memorization

These results suggest XGBoost works better than other machine learning methods.

7.5 Deep Learning Results

This section reports neural network learning patterns and their added value in modeling complex survival relationships.

7.5.1 Neural Network Learning Patterns

The NN accuracy trend (see Figure 14) indicates:

- smooth training
- stable validation accuracy
- minimal overfitting
- stabilization occurs around 40 to 50 iterations

Even though it's less accurate - around 60% - compared to machine learning methods, what matters most is its role in:

- uncovering complex patterns in data through interaction effects
- offering additional insights into potential risks
- modeling high-dimensional relationships

Deep learning works better when data grows - yet even then, it adds useful organization.

7.6 SHAP Explainability Insights

This section explains how key predictors influence model outputs, enhancing interpretability of survival forecasts.

7.6.1 Global Interpretability

SHAP summary charts (Figures 1–9) show -

- BMI and Tumor Size heavily shift risk
- Younger age lowers chances for kids
- KRAS linked to higher death risk in forecasts
- EGFR reduces it

This matches:

- Cox hazard ratios
- Random Forest importance
- XGBoost feature influence

SHAP boosts trust in machine learning results within healthcare settings by clarifying how inputs affect outputs.

7.6.2 Detailed Feature Behavior

Dependence plots (Figures 1–9) show:

- As tumor size grows, the likelihood of higher SHAP values rises significantly - particularly when expansion accelerates
- Youth tends to lower risk steadily
- Combination treatment lowers SHAP forecasts
- KRAS mutations shift forecasts in favor of higher risk - shown by positive SHAP scores

These graphs show how models work, making them easier to understand in practice - important so doctors can actually use them.

7.7 Integrated Discussion

Across all techniques - Kaplan–Meier, Cox, Random Forest, XGBoost, Neural Networks, and SHAP - the results converge on the same core risk factors:

Major Survival Determinants:

- Stage
- Tumor Size
- Age
- BMI
- Treatment Type
- EGFR
- KRAS

The mix of stats and AI methods supports the results' reliability - using different approaches strengthens confidence.

Here's why it's relevant in practice:

- Spotting issues early greatly boosts chances of living longer
- Combination therapy works better most of the time
- EGFR changes tend to improve survival outcomes - whereas KRAS alterations usually make them worse
- Tumor size still matters most when predicting outcomes
- AI models assist in early risk detection

Strength of Hybrid Modelling

This research indicates that using -

- Statistical inference
- AI/ML prediction
- Explainability

builds a strong tool to help with leukemia outcome decisions.

8.LIMITATION

- **Limited Genetic Coverage:**
The study includes only two genetic markers (EGFR and KRAS). In real-world leukemia prognosis, a broader panel of molecular mutations (FLT3, NPM1, TP53, IDH1/2, etc.) plays a major role, so genetic granularity is limited.
- **Absence of Longitudinal Clinical Data:**
Time-dependent clinical indicators such as periodic blood counts, treatment cycles, remission/relapse timelines, and MRD (Minimal Residual Disease) measurements were not included, restricting the ability to model dynamic disease progression.
- **Simplified Treatment Information:**
Treatment types were grouped into broad categories (Chemotherapy, Radiation, Surgery, Combination). Actual clinical practice involves detailed treatment protocols, dosage schedules, and targeted therapies that were not captured.
- **Limited Regional and Socioeconomic Factors:**
Although regional variation was included, influencing factors such as hospital quality, access to specialized care, socioeconomic status, and diagnostic delays were not part of the dataset, limiting contextual depth.
- **No External Validation:**
The models were evaluated only on internal train–test splits. Validation using an independent hospital or multi-center dataset is necessary to assess generalizability and real-world reliability.

9.CONCLUSION

This research carried out a detailed survival assessment plus prediction modelling for individuals with leukemia in various parts of India, using traditional statistics along with current machine learning and neural network techniques. The aim was to uncover major medical, population-based, therapy-linked, and hereditary elements affecting survival rates, while also assessing how well artificial intelligence tools can assist doctors in making choices. Findings from Kaplan–Meier curves, Log-Rank analyses, ANOVA tests, and Cox regression models repeatedly pointed to disease phase, growth volume, age, body mass index, intervention method, and gene changes (EGFR or KRAS) as critical influences on lifespan. Those identified in early phases together with people undergoing combined treatments showed clearly improved results, supporting established medical observations.

Beyond number-based conclusions, machine learning methods - especially Random Forest together with XGBoost - performed well in forecasting outcomes, reaching ROC–AUC scores near 0.80, which suggests they might help detect at-risk individuals. Although somewhat less accurate, the deep neural network uncovered subtle nonlinear patterns in the dataset, adding extra insight alongside traditional models. Crucially, using SHAP for clarity made it easier to understand how predictions were formed, helping align technical results with practical medical judgment.

This research shows how mixing standard survival analysis with AI-based prediction helps grasp leukemia survival trends in India more fully. Findings offer useful clues for spotting cases earlier, shaping therapy plans, while sorting patients by risk level. Crucially, this combined method creates a solid base for building tools that aid medical decisions - boosting health results, also possibly steering resources better across hospitals. Work highlights: smart use of data, if clear and ethical, may greatly lift cancer care quality.

10.REFERENCES

- [1] Bello, G. A., Yu, S., & Ghosh, D. (2021). Machine learning for survival analysis in leukemia research: Evaluation of predictive techniques for hematological malignancies. *Journal of Hematology & Oncology Research*, 13(4), 212–225.
- [2] Cortes, J. E., Pinilla-Ibarz, J., & Jain, N. (2020). Prognostic impact of molecular mutations in leukemia: A comprehensive review. *Blood Cancer Journal*, 10(3), 1–12. <https://doi.org/10.1038/s41408-020-0291-z>
- [3] Cox, D. R. (1972). Regression models and life-tables. *Journal of the Royal Statistical Society: Series B (Methodological)*, 34(2), 187–220.
- [4] ICMR National Cancer Registry Programme. (2021). *Cancer statistics report – India*. Indian Council of Medical Research, New Delhi.
- [5] Kaplan, E. L., & Meier, P. (1958). Nonparametric estimation from incomplete observations. *Journal of the American Statistical Association*, 53(282), 457–481.
- .
- [6] Katzman, J. L., Shaham, U., Cloninger, A., Bates, J., Jiang, T., & Kluger, Y. (2018). DeepSurv: Personalized treatment recommendation using a Cox neural network. *Machine Learning for Healthcare Conference*, 15, 1–10.
- [7] Nair, S. C., George, B., & Viswanath, B. (2019). Regional disparities in leukemia incidence and survival outcomes in India: A population-level analysis. *Indian Journal of Cancer*, 56(2), 150–157.
- [8] Patel, A. B., & Deininger, M. W. (2017). The role of EGFR and KRAS mutations in leukemia progression and therapeutic response. *Hematology Reviews*, 11(1), 25–39.
- [9] Voutilainen, A., Hänninen, A., & Jantunen, E. (2019). Predictive modelling of leukemia survival using demographic and clinical features. *Leukemia Research*, 82, 45–52.
- [10] Zhang, H., Li, X., & Sun, Y. (2021). Survival prediction in hematologic malignancies using integrated clinical and genetic data. *Frontiers in Oncology*, 11, 1–10. <https://doi.org/10.3389/fonc.2021.667893>

- [11] Chatterjee, T., Banerjee, S., & Bose, P. (2020). Genetic determinants of leukemia outcomes: A review of clinical significance in targeted therapy. *Cancer Genetics*, 242, 1–12.
- [12] Zhou, X., Wang, M., & Zhang, L. (2020). Machine-learning-assisted risk stratification in acute leukemia using clinical biomarkers. *BMC Medical Informatics and Decision Making*, 20(1), 1–9.
- [13] Singh, A., & Kumar, P. (2018). Leukemia patterns and treatment outcomes in North India: A five-year retrospective study. *Indian Journal of Hematology and Blood Transfusion*, 34(3), 532–540.
- [14] Li, Y., Chen, F., & Zhang, W. (2022). Combining clinical and genomic features for leukemia survival prediction using XGBoost. *Scientific Reports*, 12, 345–355.
- [15] Kourou, K., Exarchos, T., Exarchos, K., Karamouzis, M. V., & Fotiadis, D. I. (2015). Machine learning applications in cancer prognosis and prediction. *Computational and Structural Biotechnology Journal*, 13, 8–17.
- [16] Kayastha, G., & Sharma, R. (2020). Analysis of acute leukemia treatment patterns in Indian tertiary hospitals. *Asian Pacific Journal of Cancer Prevention*, 21(9), 2715–2722.
- [17] Ng, A. Y., & Jordan, M. I. (2002). On discriminative vs generative classifiers: A comparison of logistic regression and naive Bayes. *Advances in Neural Information Processing Systems*, 14, 841–848.
- [18] Polley, M.-Y. C., LeBlanc, M., & Crowley, J. (2011). Detecting treatment-covariate interactions in the Cox model using machine learning. *Statistics in Medicine*, 30(14), 1671–1684.
- [19] Friedman, J. H. (2001). Greedy function approximation: A gradient boosting machine. *Annals of Statistics*, 29(5), 1189–1232.
- [20] Breiman, L. (2001). Random forests. *Machine Learning*, 45(1), 5–32.