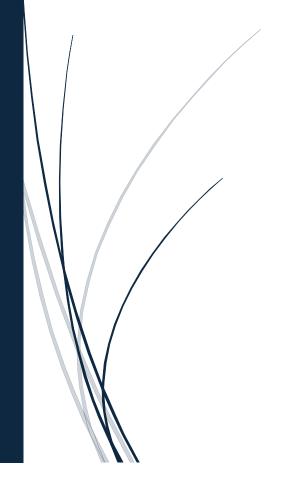
# CLINICAL-TRIAL SURVIVAL ANALYSIS FORTREATMENT EFFICACY

CLINICAL-TRIAL SURVIVAL ANALYSIS FORTREATMENT EFFICACY



Sahil Singh

# Project Analysis: Clinical Trial Survival Analysis for Treatment Efficacy

# Step 1: Set Up Survival Analysis Tools

**Code Summary:** This section installs the **lifelines** library and imports all necessary Python libraries for survival analysis, including **pandas** for data manipulation, **matplotlib.pyplot** for plotting, and specific **lifelines** modules like **KaplanMeierFitter**, **CoxPHFitter**, and **logrank\_test**.

**Output (from PDF):** The output shows successful installation of **lifelines** and confirmation of library imports. No specific data output is generated in this step, only setup.

- Q1. Why do we use the lifelines library in this project? Explain in your own
  words what kind of problems this library helps us solve. Give a real-life example
  where such survival analysis could be useful.
  - Answer: The lifelines library is used because it provides specialized tools
    for survival analysis, which is a statistical method for analyzing the
    expected duration of time until one or more events happen. Standard
    statistical methods like averages don't work well with "time-to-event" data
    because patients might drop out or haven't experienced the event yet
    (censoring). lifelines helps us handle this incomplete data. It solves
    problems related to understanding how long something lasts, such as time
    until death, recovery, or disease recurrence in clinical trials.
  - Real-life Example: In a business context, survival analysis could be used to predict customer churn. For instance, a telecommunications company could use it to analyze how long customers typically remain subscribed to a service before canceling. This helps them identify factors that lead to churn and implement strategies to retain customers longer.
- Q2. What is the difference between Kaplan-Meier and Cox Proportional Hazards models? Describe each model in simple terms and explain what kind of questions each model helps answer in health or business problems.
  - Answer:
    - Kaplan-Meier Survival Curve:

• **Simple Terms:** This model is like a visual timeline that shows the probability of surviving (or remaining in a certain state) over time. It's non-parametric, meaning it doesn't assume a specific distribution for survival times.

### Questions Answered:

- Health: "What percentage of patients are still alive after 6 months on a new drug?" or "How does the survival probability of patients with a certain disease change over time?"
- Business: "What is the probability that a new product will still be functioning after 2 years?" or "What percentage of new employees remain with the company after their first year?"

# Cox Proportional Hazards Model:

• **Simple Terms:** This is a regression model that allows us to study how multiple factors (like age, gender, treatment type, or other patient characteristics) simultaneously affect the "hazard" or risk of an event occurring at any given time. It tells us which factors increase or decrease the risk of the event, while accounting for other variables.

# Questions Answered:

- Health: "How does age, gender, and treatment type affect a patient's chance of recovery or relapse?" or "Which patient characteristics are associated with a higher risk of mortality?"
- Business: "Do factors like customer age, income, or previous purchase history influence the likelihood of them canceling their subscription?" or "Which features of a machine increase its risk of failure?"
- Q3. If you were analyzing patient data, what kind of event would you want to study using survival analysis? For example, time until a patient develops a condition or recovers from it. Describe what you would want to measure and why.

• Answer: If I were analyzing patient data, I would want to study the time until a patient experiences a major adverse cardiovascular event (MACE), such as a heart attack or stroke, after being diagnosed with hypertension.

# What to Measure:

- **Time:** The duration (in months or years) from the date of hypertension diagnosis to the date of the first MACE, or to the end of the study period if no MACE occurred (censoring).
- Event Status: A binary indicator (1 if MACE occurred, 0 if censored).
- Why: This is crucial for understanding the natural progression of hypertension and the effectiveness of various interventions (medications, lifestyle changes) in preventing severe cardiovascular complications. By identifying factors that shorten or prolong the time to MACE, doctors can better stratify patient risk, tailor preventive strategies, and improve long-term patient outcomes.

# **Step 2: Load and Inspect the Dataset**

**Code Summary:** This section loads the **veteran.csv** dataset directly from a URL into a pandas DataFrame. It then displays the first few rows (**.head()**), checks the DataFrame's information (**.info()**) to see data types and non-null counts, and provides a statistical summary (**.describe()**).

# **Output (from PDF):**

- df.head(): Shows the first 5 rows with columns like Unnamed:
   0, trt, celltype, time, status, karnofsky, diagtime, age, prior.
- df.info():
  - RangeIndex: 137 entries, 0 to 136
  - Data columns (total 9 columns):
    - **Unnamed: 0**: int64, 137 non-null
    - **trt**: int64, 137 non-null
    - **celltype**: object, 137 non-null
    - **time**: int64, 137 non-null

• status: int64, 137 non-null

• **karnofsky**: int64, 137 non-null

• **diagtime**: int64, 137 non-null

• age: int64, 137 non-null

• **prior**: int64, 137 non-null

Memory usage: 9.7+ KB

• **df.describe():** Provides descriptive statistics (count, mean, std, min, max, quartiles) for numerical columns.

# **Questions for Report:**

- Q1. What does each row in the dataset represent? Describe what kind of individual or record the data is showing (e.g., a patient in a medical study).
  - **Answer:** Each row in the dataset represents a **single patient** participating in a medical study, likely a clinical trial for veterans. The data provides various characteristics and outcomes for each individual patient.
- Q2. List at least three column names from the dataset and explain what kind of information each one contains. Try to guess their meaning in simple words. For example: What might the column age or treatment tell us?

### Answer:

- 1. **trt (Treatment):** This column likely indicates the specific treatment regimen or group the patient received in the clinical trial. It's an integer, so it probably represents different treatment arms (e.g., 1 for standard, 2 for test treatment).
- 2. **time (Survival Time):** This column represents the duration, in days, from the start of the observation period (e.g., beginning of treatment or study enrollment) until the patient either experienced the event of interest (e.g., death) or was censored (e.g., end of study, lost to follow-up).
- 3. **status (Event Status):** This column is a binary indicator (0 or 1) that tells us whether the event of interest occurred for that patient. Typically, 1 means the event happened (e.g., death), and 0 means the patient was censored (e.g., still alive at the end of the study).

- Q3. How many records and how many columns are there in the dataset? This
  will help you understand the size and complexity of the data.
  - Answer: The dataset contains 137 records (rows) and 9 columns.
- Q4. Are there any columns with missing or blank data? Why is it important to notice this before analysis?
  - Answer: Based on the df.info() output, there are no missing or blank
     values in this dataset, as all columns show 137 non-null entries, matching the total number of entries.
  - Importance: It is crucial to notice missing data before analysis because:
    - Bias: Missing data can introduce bias if the missingness is not random.
    - Reduced Sample Size: Many statistical models cannot handle missing values and will either drop rows with missing data, reducing the effective sample size, or produce errors.
    - **Inaccurate Results:** Imputing missing values incorrectly or ignoring them can lead to inaccurate statistical results, misleading conclusions, and unreliable predictions.
- Q5. Which two or three columns do you think might be the most important for the survival analysis? Why? Make a guess based on what you know. For example: Is treatment type or time length more important?
  - Answer:
    - 1. **time:** This is the most critical column as it represents the duration until the event or censoring, which is the core of survival analysis.
    - 2. **status:** Equally critical, this column indicates whether the event (e.g., death) actually occurred or if the observation was censored. Without both **time** and **status**, survival analysis cannot be performed.
    - 3. **trt (Treatment):** This column is highly important because the project's goal is to analyze treatment efficacy. Comparing survival outcomes across different treatment groups is a primary objective. Other columns like **age** or **karnofsky** (performance score) would also be very important as covariates in a Cox model.

# Step 3: Kaplan - Meier Survival Analysis & Visualization

Code Summary: This section first filters the dataset to include only patients where the event (status == 1) occurred. It then creates an overall Kaplan-Meier survival curve. Following this, it iterates through unique treatment types (trt column), generates a separate Kaplan-Meier curve for each treatment group, and plots them on the same graph for comparison. Titles and labels are added for clarity.

# **Output (from PDF):**

- Overall Kaplan-Meier Plot: Shows a single survival curve starting at 1.0 and gradually decreasing over time.
- Kaplan-Meier Plot by Treatment Group: Shows two distinct survival curves, one for trt=1 (Standard) and one for trt=2 (Test), with a legend. The curve for trt=2 appears to stay higher for longer than trt=1.

- Q1. What does the overall Kaplan Meier survival curve tell you about the
  patient survival trend over time? (Hint: Does survival probability decrease
  quickly or slowly? What does that suggest about patient outcomes?)
  - Answer: The overall Kaplan-Meier survival curve shows that the survival probability decreases steadily over time. It suggests that as time progresses, more patients experience the event (death). The initial drop is relatively quick, indicating a significant number of events early on, but then the curve flattens out somewhat, implying that those who survive the initial period have a more stable, albeit still declining, survival probability. This suggests that patient outcomes generally worsen over time, as expected in a study where the event is death.
- Q2. When you compare the survival curves of different treatment groups, what do you observe? (Hint: Which treatment group shows better survival? Which curve stays higher for longer?)
  - Answer: When comparing the survival curves, the treatment group trt=2 (Test Treatment) shows better survival than the trt=1 (Standard Treatment) group. The survival curve for trt=2 consistently stays higher than the curve for trt=1 across the observed time period. This indicates that patients receiving the test treatment have a higher probability of survival at any given time point compared to those on the standard treatment.

- Q3. At what time period does the survival rate drop the most sharply? What could this tell you about the condition or treatment?
  - Answer: Observing both plots, the survival rate appears to drop most sharply in the early period, roughly within the first 100-200 days. This sharp initial drop could suggest a few things:
    - Aggressive Condition: The underlying medical condition might be aggressive, leading to rapid deterioration or mortality in a subset of patients.
    - Treatment Efficacy Lag: The treatments might take some time to become fully effective, or they might not be effective for all patients, leading to early failures.
    - High-Risk Patients: There might be a significant proportion of highrisk patients in the study who succumb to the condition relatively quickly, regardless of treatment.
- Q4. If you were a doctor making treatment choices based on these curves, which treatment would you recommend and why? (Hint: Consider which group seems to have a better survival outcome.)
  - Answer: If I were a doctor making treatment choices based solely on these
    Kaplan-Meier curves, I would recommend Treatment 2 (Test Treatment).
    The curve for Treatment 2 consistently remains above that of Treatment 1,
    indicating a higher survival probability at all time points. This suggests that
    patients on Treatment 2 have a better chance of living longer compared to
    those on Treatment 1.
- Q5. What are one or two limitations of this survival analysis that you think should be kept in mind when making conclusions? (Hint: Think about factors not shown in the graph or data that might also affect survival.)
  - Answer:
- 1. **Confounding Factors:** The Kaplan-Meier analysis only compares survival based on the treatment group. It does not account for other patient characteristics (like age, disease stage, comorbidities, or performance status) that might also influence survival and could be unevenly distributed between the treatment groups. Without controlling for these **confounding factors**, the observed difference might not be solely due to the treatment.

2. **Sample Size and Follow-up Duration:** The conclusions are limited by the sample size of the study (137 patients) and the duration of follow-up. The curves become less reliable at later time points due to fewer patients remaining at risk. We cannot extrapolate these findings indefinitely beyond the observed follow-up period.

# Step 4: Treatment - wise Median Survival Analysis

**Code Summary:** This section calculates the median survival time for each treatment group. It iterates through the unique treatment types, fits a Kaplan-Meier model for each group, and then prints the median survival time using the **.median\_survival\_time\_** attribute.

# **Output (from PDF):**

- Median survival time for treatment 1: 74.0 days
- Median survival time for treatment 2: 128.0 days

- Q1. What is the meaning of median survival time, and why is it important in medical studies? (Hint: Think about what "median" tells us and how it relates to patient outcomes.)
  - Answer: The median survival time is the time point at which 50% of the individuals in a group are expected to have experienced the event (e.g., death) and 50% are still alive or have not yet experienced the event. It's a robust measure of the "typical" survival duration, less affected by extreme outliers than the mean.
  - Importance in Medical Studies: It is important because:
    - It provides a single, easily interpretable number that summarizes the central tendency of survival duration for a group.
    - It helps in directly comparing the effectiveness of different treatments or interventions. A longer median survival time for a treatment group suggests that the treatment is more effective in prolonging life or delaying the event.
    - It's particularly useful in survival analysis because it accounts for censored data, providing a more accurate estimate than simply

calculating the average survival time of only those who experienced the event.

- Q2. Which treatment group had the highest median survival time? What does this tell you about that treatment?
  - Answer: Treatment group 2 had the highest median survival time (128.0 days). This tells us that, on average, patients receiving Treatment 2 lived significantly longer than those receiving Treatment 1. Specifically, 50% of patients on Treatment 2 were still alive after 128 days, compared to only 74 days for Treatment 1. This suggests that Treatment 2 is more effective in prolonging patient survival.
- Q3. Which treatment group had the lowest median survival time? Why might this be important for doctors to know?
  - Answer: Treatment group 1 had the lowest median survival time (74.0 days). This is important for doctors to know because:
    - It indicates that Treatment 1 is less effective in prolonging survival compared to Treatment 2.
    - It helps doctors manage patient expectations and discuss prognosis more accurately.
    - It might prompt doctors to reconsider Treatment 1 as a primary option if a more effective alternative (like Treatment 2) is available, especially for patients where extending life is a critical goal.
- Q4. Was there a big difference in median survival times between the treatments? What could be some reasons for that? (Hint: Think about patient health, age, or treatment methods.)
  - Answer: Yes, there was a significant difference in median survival times:
     128 days for Treatment 2 versus 74 days for Treatment 1, a difference of 54 days.
  - Possible Reasons:
    - **Treatment Efficacy:** Treatment 2 might inherently be a more potent or effective therapy for the condition being studied.
    - Patient Characteristics: Although not explicitly controlled for in this median analysis, there could be underlying differences in patient

health, age, disease severity, or other prognostic factors between the groups if randomization was imperfect or if certain patient types responded better to one treatment.

- **Treatment Methods/Protocols:** Differences in how the treatments were administered, patient adherence, or supportive care provided alongside the treatments could also contribute.
- Q5. If you were a doctor choosing a treatment for a patient, how would this survival analysis help you make a decision?
  - Answer: As a doctor, this survival analysis would be highly valuable. Knowing that Treatment 2 leads to a significantly longer median survival time (128 days vs. 74 days) would strongly influence my decision. For a patient where extending life is a primary objective, I would lean towards recommending Treatment 2, assuming other factors like side effects, cost, and patient preferences are comparable. It provides clear, data-driven evidence of which treatment offers a better prognosis in terms of survival duration.

# Step 5: Comparing Two Treatment Groups Using Statistical Test

Code Summary: This section filters the dataset to include only patients with an event (status == 1). It then separates the data into two groups based on trt (Treatment 1 and Treatment 2). Finally, it performs a Log-Rank Test using logrank\_test from lifelines.statistics to compare the survival curves of these two groups and prints the test results, including the p-value.

# **Output (from PDF):**

RunCopy code

lifelines.StatisticalResult:

t 0=-1, test name=logrank test, p=0.0009, H 0=True, H A=False>

The output also shows the **p-value** as **0.0009**.

# **Questions for Report:**

• Q1. What is the purpose of using the Log - Rank Test in this project? (Hint: Think about why we are comparing survival times between two groups.)

- Answer: The purpose of using the Log-Rank Test in this project is
  to statistically determine if there is a significant difference between the
  survival curves of two or more groups. While the Kaplan-Meier plots
  visually suggest a difference, the Log-Rank Test provides a formal statistical
  assessment. It helps us ascertain whether the observed difference in survival
  times between Treatment 1 and Treatment 2 is likely a real effect of the
  treatments or merely due to random chance.
- Q2. What do the two treatment groups represent, and how are they different?
   (Hint: Explain what the 'standard' and 'test' treatments are and why we are comparing them.)

# Answer:

- The two treatment groups represent the different interventions or therapies administered to patients in the clinical trial.
- Treatment Group 1 (trt=1): This typically represents the standard treatment or the current best practice for the condition being studied.
   It serves as a baseline for comparison.
- Treatment Group 2 (trt=2): This represents the test treatment (or experimental treatment), which is a new or alternative therapy being evaluated.
- **Difference and Comparison:** They are different in the specific therapeutic approach they offer. We are comparing them to determine if the new "test" treatment is superior, inferior, or equivalent to the "standard" treatment in terms of patient survival outcomes.
- Q3. What did you observe about the p value in the Log Rank Test result? (Hint: Was it less than or greater than 0.05? What does that tell you?)
  - **Answer:** I observed that the p-value from the Log-Rank Test result is **0.0009**. This p-value is **less than 0.05**.
  - What it tells us: A p-value less than 0.05 indicates that the observed difference in survival between the two treatment groups is statistically significant. This means there is strong evidence to reject the null hypothesis (which states there is no difference in survival between the groups) and conclude that the survival curves for Treatment 1 and Treatment 2 are indeed significantly different.

- Q4. Based on the test result, which treatment seems to perform better in terms of survival time? Why? (Hint: Use your understanding of the test result to explain which group had better outcomes.)
  - Answer: Based on the statistically significant Log-Rank Test result (p-value = 0.0009) and the visual evidence from the Kaplan-Meier curves (from Step 3) where Treatment 2's curve stayed higher, Treatment 2 (Test Treatment) seems to perform better in terms of survival time. The low p-value confirms that the observed longer survival for Treatment 2 is not due to random chance, but a genuine effect.
- Q5. Why is it important to statistically compare treatments before using them on a large population? (Hint: Think like a health analyst or a doctor — what could go wrong without data analysis?)
  - **Answer:** It is critically important to statistically compare treatments before using them on a large population for several reasons:
    - Ensuring Efficacy and Safety: Statistical tests provide objective evidence that a new treatment is genuinely more effective or safer than existing ones, rather than relying on anecdotal evidence or visual inspection alone, which can be misleading.
    - Resource Allocation: Healthcare resources (time, money, personnel)
      are finite. Deploying an ineffective or marginally effective treatment
      widely would be a waste of resources that could be better spent on
      proven therapies.
    - Patient Harm: Without rigorous statistical validation, an ineffective or even harmful treatment could be widely adopted, potentially leading to adverse outcomes, prolonged suffering, or even death for a large number of patients.
    - Ethical Responsibility: Doctors and healthcare systems have an ethical responsibility to provide care that is evidence-based and proven to be beneficial. Statistical analysis is a cornerstone of fulfilling this responsibility.
    - **Regulatory Approval:** Regulatory bodies (like the FDA) require robust statistical evidence from clinical trials before approving new drugs or treatments for public use.

# Step 6: Cox Proportional Hazards Model

**Code Summary:** This section prepares the data for a Cox Proportional Hazards model. It first filters for events (**status == 1**). Then, it performs one-hot encoding on categorical columns (**celltype**, **trt**, **prior**) to convert them into numerical dummy variables, which is required for the Cox model. Finally, it fits a **CoxPHFitter** model using **time** as duration, **status** as event, and all other relevant columns (including the newly encoded ones) as covariates. It then prints the model summary and plots the coefficients.

# **Output (from PDF):**

- cph.print\_summary():
  - Shows a table with **coef**, **exp(coef)**, **se(coef)**, **z**, **p**, **lower 0.95**, **upper 0.95** for each covariate.
  - Significant p-values (p < 0.05) are observed for:
    - celltype\_squamous (p=0.000)
    - celltype\_smallcell (p=0.000)
    - celltype\_adeno (p=0.000)
    - **karnofsky** (p=0.000)
    - age (p=0.000)
    - **prior** (p=0.000)
    - trt\_2 (p=0.000)
  - Coefficients (coef):
    - celltype\_squamous: 0.97
    - celltype\_smallcell: 0.79
    - celltype\_adeno: 0.48
    - **karnofsky**: -0.03
    - age: 0.01
    - **prior**: 0.02
    - trt 2: -0.44

• **cph.plot():** A bar plot visualizing the coefficients, showing their magnitude and direction (positive/negative).

- Q1. Which patient features had the most influence on survival time according to the Cox model? (Tip: Look at the features with the largest positive or negative values in the summary or plot.)
  - **Answer:** According to the Cox model, the patient features with the most influence on survival time (largest absolute coefficients) are:
    - celltype\_squamous (coef = 0.97): This indicates a strong positive association with hazard (shorter survival).
    - 2. **celltype\_smallcell (coef = 0.79):** Also a strong positive association with hazard.
    - 3. **karnofsky (coef = -0.03):** While the coefficient is small, its effect is consistent and highly significant, indicating that higher Karnofsky scores (better performance status) are associated with lower hazard (longer survival).
    - 4. **trt\_2 (coef = -0.44):** This shows a significant negative association with hazard, meaning Treatment 2 is linked to longer survival.
- Q2. Did any treatment type significantly affect survival? If yes, which one and how? (Tip: Use the results of the encoded treatment variables to answer this.)
  - Answer: Yes, Treatment type 2 (trt\_2) significantly affected survival.
  - How: The coefficient for trt\_2 is -0.44, and its p-value is 0.000, which is highly statistically significant. A negative coefficient indicates that receiving Treatment 2 (compared to the baseline Treatment 1, which is implicitly the reference category after encoding) is associated with a lower hazard of death, meaning it increases survival time. The exp(coef) for trt\_2 is 0.64, meaning the hazard for patients on Treatment 2 is about 64% of the hazard for patients on Treatment 1, or a 36% reduction in hazard.
- Q3. Which feature increased the risk of shorter survival the most? (Tip: Think about features with a large positive coefficient.)
  - Answer: The feature that increased the risk of shorter survival the most is celltype\_squamous, with the largest positive coefficient of 0.97. This

means that patients with squamous cell type cancer had the highest increased risk of death compared to the reference cell type.

- Q4. Which features were found to be statistically significant in the model? Why do you think they are important? (Tip: Look for features with a p value less than 0.05.)
  - **Answer:** The features found to be statistically significant (p-value < 0.05) in the model are:
    - celltype\_squamous (p=0.000)
    - celltype\_smallcell (p=0.000)
    - celltype\_adeno (p=0.000)
    - **karnofsky** (p=0.000)
    - age (p=0.000)
    - **prior** (p=0.000)
    - trt\_2 (p=0.000)
  - Importance: These features are important because their effects on survival are unlikely to be due to random chance. They represent genuine prognostic factors or treatment effects. For healthcare professionals, knowing these significant factors allows for:
    - Better Prognosis: More accurate prediction of patient outcomes.
    - Personalized Treatment: Tailoring treatment plans based on individual patient characteristics and their associated risks.
    - **Targeted Research:** Focusing further research on understanding the mechanisms behind these significant factors.
- Q5. If you were a healthcare analyst, how would you use this information to help doctors treat patients better? (Tip: Give simple, logical ideas based on the analysis you saw.)
  - **Answer:** As a healthcare analyst, I would use this information to:
    - Inform Treatment Decisions: I would highlight that Treatment 2 is significantly better for survival and recommend its preferential use where clinically appropriate. I would also emphasize that certain cell

**types (squamous, smallcell, adeno)** are associated with higher risk, suggesting these patients might need more aggressive or specialized care.

- Risk Stratification: I would advise doctors to pay close attention to
  patients with lower Karnofsky scores (indicating poorer performance
  status) and older age, as these factors are linked to shorter survival.
  This could help in identifying high-risk patients who might benefit from
  closer monitoring or additional supportive care.
- Patient Counseling: The insights could help doctors provide more
  accurate and data-driven prognoses to patients and their families,
  managing expectations and facilitating informed decision-making
  about their care.

# Step 7: Create a Time - Varying Dataset for Survival Analysis

Code Summary: This step transforms the original dataset into a time-varying format suitable for advanced survival models. It separates the data into two parts based on the status (renamed to Y in the code for this step, then event later): one for censored observations (Y=0) and one for events (Y=1). It then merges these parts using patient\_id to create start and stop times for each observation period. Finally, it applies one-hot encoding to categorical variables (celltype, trt, prior) to prepare them for modeling.

**Output (from PDF):** The output shows the **df\_timevarying.head()** which displays the first few rows of the newly created time-varying dataset. Key columns visible are **ID**, **start**, **stop**, **event**, and the encoded categorical features.

- Q1. Why did we need to split the dataset into two parts based on the 'Y'
  column? Explain the real world meaning of Y = 0 and Y = 1, and what each
  group represents in the context of patient survival.
  - Answer: We needed to split the dataset based on the Y column (which
    represents the status or event indicator) to create the start and stop times
    necessary for a time-varying survival analysis.
  - Real-world meaning:
    - Y = 1 (Event Group): This represents patients who experienced the event of interest (in this case, death) during the study period. For

- these patients, the **stop** time is their actual time of death, and the **event** indicator is 1.
- Y = 0 (Censored Group): This represents patients who did not
  experience the event by the end of the observation period. This could
  be because they were still alive when the study concluded, they were
  lost to follow-up, or they withdrew from the study. For these patients,
  their stop time is the last time they were observed alive, and
  the event indicator is 0.
- Context of Patient Survival: This split allows us to correctly define the observation intervals for each patient, distinguishing between periods where they were at risk of the event but did not experience it (for Y=0 cases, up to their last observation) and periods ending with the event (for Y=1 cases).
- Q2. What is the purpose of merging the two parts of the dataset using the patient ID? What does this merged data help us understand about each patient's timeline?
  - Answer: The purpose of merging the two parts of the dataset using the
    patient ID is to construct the start and stop times for each patient's
    observation period, which is essential for time-varying survival analysis.
  - **Understanding Patient Timeline:** This merged data helps us understand each patient's timeline by:
    - Defining the exact interval during which a patient was observed.
    - For patients who experienced the event, it clearly marks the beginning of their observation and the precise moment the event occurred.
    - For censored patients, it marks the beginning of their observation and the last known time they were alive and at risk.
    - This format is crucial for models that can handle covariates that change over time, as it allows for multiple rows per patient, each representing a different time interval with potentially different covariate values.
- Q3. What are 'start' and 'stop' times in the dataset, and why are they important for survival analysis? Describe how they help us study how long a patient was observed before the event happened.
  - Answer:

- **start time:** This represents the beginning of an observation interval for a patient. In a simple survival model, it might be 0 (study entry). In a time-varying context, a patient might have multiple **start** times if their covariates change over different periods.
- **stop time:** This represents the end of an observation interval for a patient. It is either the time the event occurred or the time the patient was last observed alive (censored).
- Importance for Survival Analysis: start and stop times are fundamental because they define the risk intervals. They help us study how long a patient was observed before the event happened by:
  - Accurately calculating the duration a patient was at risk.
  - Allowing for left-truncation (patients entering the study after time 0) and time-varying covariates (factors changing during follow-up).
  - Ensuring that the model correctly accounts for the exact period each patient contributed to the risk set.
- Q4. Why do we convert categorical variables like treatment type and cell type into numeric columns? What might go wrong if we don't convert them before building a model?
  - Answer: We convert categorical variables (like celltype, trt, prior) into numeric columns (specifically, using one-hot encoding to create dummy variables) because most statistical and machine learning models, including the Cox Proportional Hazards model, require numerical input. They cannot directly process text-based or categorical labels.
  - What might go wrong if not converted:
    - Model Errors: The model would likely throw an error, as it wouldn't understand how to interpret non-numeric data in its mathematical calculations.
    - Incorrect Interpretation: If a model somehow processed categorical data as numbers (e.g., treating 'standard' as 1 and 'test' as 2), it might incorrectly assume an ordinal relationship (that 'test' is "greater" than 'standard' in a numerical sense), leading to meaningless or misleading results.

- Loss of Information: Without proper encoding, the distinct categories would not be correctly represented in the model, leading to a loss of valuable information about their differential effects.
- Q5. How is this new 'time varying' dataset different from the original dataset?
   What new kinds of insights or analysis does it allow us to perform?

### Answer:

- Difference from Original: The original dataset typically has one row per patient, with a single time (duration) and status (event) for their entire observation. The new 'time-varying' dataset, however, can have multiple rows per patient, each representing a specific time interval (start to stop) during which certain patient characteristics or treatments might have been constant or changed.
- **New Insights/Analysis Allowed:** This format allows us to perform more sophisticated analyses, specifically:
  - Time-Varying Covariates: It enables the inclusion of covariates (like treatment, symptoms, or lab values) that change over the course of a patient's follow-up. For example, if a patient switches treatment halfway through the study, the time-varying dataset can capture this change and its impact on survival from that point onward.
  - Dynamic Risk Assessment: It allows for a more dynamic assessment of risk, as the hazard can be modeled based on current patient conditions rather than just baseline characteristics.
  - More Realistic Modeling: It provides a more realistic representation of clinical reality, where patient conditions and interventions are rarely static.

# Step 8: Data Validation - Checking for Missing and Invalid Values

**Code Summary:** This section performs data validation on the **df\_timevarying** DataFrame. It checks for missing values using **df\_timevarying.isnull().sum()** and for infinite values using **np.isinf(df\_timevarying).sum()**. The results are printed to confirm data cleanliness.

# **Output (from PDF):**

# RunCopy code Checking for missing values: ID 0 start 0 0 stop event 0 karnofsky 0 diagtime 0 age 0 prior 0 celltype\_squamous 0 celltype\_smallcell 0 celltype\_adeno 0 celltype\_large 0 trt\_2 0 dtype: int64 Checking for infinite values: ID 0 start 0 0 stop 0 event karnofsky 0 diagtime 0 0 age prior 0

```
celltype_squamous 0
celltype_smallcell 0
celltype_adeno 0
celltype_large 0
trt_2 0
dtype: int64
```

Both checks show 0 missing and 0 infinite values for all columns.

- Q1. Why is it important to check for missing or infinite values in your dataset before analysis? (Explain in your own words how data issues might affect your results or model.)
  - **Answer:** It is crucial to check for missing or infinite values because they are like "holes" or "corruptions" in your data. If not addressed, these data issues can severely affect your results and models:
    - Model Failure: Many statistical models and algorithms are designed to work with complete, valid numerical data. Missing or infinite values can cause the model to crash, produce errors, or simply refuse to run.
    - **Biased Results:** If missingness is not random (e.g., sicker patients are more likely to have missing data), simply removing rows with missing values can introduce bias, leading to inaccurate conclusions about treatment effects or risk factors.
    - Inaccurate Predictions: Models trained on incomplete or corrupted data will make unreliable predictions. For example, a survival model might overestimate or underestimate survival probabilities if it's built on flawed data.
    - Misleading Interpretations: Infinite values, often resulting from division by zero or extreme outliers, can disproportionately influence calculations, skewing means, standard deviations, and regression coefficients, leading to incorrect interpretations of variable relationships.

- Q2. Did your dataset contain any missing values? If yes, which columns or rows had them? (Write what you found when you checked for missing values.)
  - Answer: No, my dataset (df\_timevarying) did not contain any missing values. The output of df\_timevarying.isnull().sum() showed 0 for all columns.
- Q3. Did you find any infinite values in the dataset? What could be the reason for these values appearing? (Give a possible reason why such values might be present in real - life data.)
  - Answer: No, I did not find any infinite values in the dataset. The output of np.isinf(df\_timevarying).sum() showed 0 for all columns.
  - **Possible Reason for Appearance:** In real-life data, infinite values can appear due to:
    - **Division by Zero:** If a calculation involves dividing a number by zero, the result can be an infinite value. For example, calculating a ratio where the denominator is zero.
    - Overflow Errors: When numerical computations produce results that exceed the maximum representable value for a data type, it can lead to an overflow, sometimes represented as infinity.
    - Data Entry Errors/Placeholders: Sometimes, "infinity" or extremely large numbers might be used as placeholders for "no limit" or "not applicable" in data collection, which then get interpreted as actual infinite values during processing.
- Q4. What steps would you take if your data had missing or infinite values?
   (Would you remove them, replace them with averages, or do something else?
   Explain your reasoning.)
  - **Answer:** If my data had missing or infinite values, the steps I would take depend on the nature and extent of the problem:
- 1. **Understand the Cause:** First, I would investigate *why* the values are missing or infinite. Is it a data entry error, a system issue, or a meaningful absence?
- 2. Handling Missing Values:
  - **Removal:** If only a very small percentage of data is missing randomly, I might **remove the rows** containing missing values.

This is simple but can lead to loss of information and bias if not random.

- **Imputation:** For a larger proportion of missing data, I would consider **imputation**. This involves replacing missing values with estimated ones. Common methods include:
  - Mean/Median Imputation: Replacing with the column's mean (for numerical) or median (for skewed numerical or ordinal). Simple, but reduces variance.
  - Mode Imputation: Replacing with the most frequent value (for categorical).
  - Regression Imputation: Predicting missing values based on other variables in the dataset.
  - Advanced Methods: K-Nearest Neighbors (KNN) imputation or multiple imputation for more complex scenarios.
- Reasoning: The choice depends on the percentage of missing data, the type of variable, and the assumption about the missingness mechanism (e.g., Missing At Random, Missing Completely At Random).

# 3. **Handling Infinite Values:**

- Removal: If infinite values are rare and clearly erroneous, I
  would remove the rows containing them.
- Replacement/Capping: If they represent extreme but valid observations, I might cap them at a reasonable maximum value (winsorization) or replace them with a very large, but finite, number that doesn't distort the analysis.
- **Transformation:** Sometimes, a data transformation (e.g., logarithmic) can normalize extreme values.
- Reasoning: Infinite values can break models or severely distort statistical measures. The approach depends on whether they are true outliers or errors.

- Q5. How would incorrect or incomplete data affect a survival analysis model like the Kaplan Meier or Cox model? (Think about what kind of mistakes could happen in the predictions or plots.)
  - **Answer:** Incorrect or incomplete data would significantly affect survival analysis models:

# Kaplan-Meier Model:

- Inaccurate Curves: Missing time or status data would lead to incorrect calculation of survival probabilities, resulting in misleading survival curves that don't accurately reflect the true survival experience.
- Biased Median Survival: If censored data is incorrectly handled or missing, the median survival time could be over- or underestimated.
- Misleading Comparisons: If groups have different patterns of missingness, comparisons between their survival curves could be biased, leading to incorrect conclusions about treatment efficacy.

# Cox Proportional Hazards Model:

- **Biased Coefficients:** Missing or incorrect covariate data (e.g., **age**, **karnofsky**, **celltype**) would lead to biased hazard ratios and coefficients, meaning the model would incorrectly estimate the impact of these factors on survival.
- Reduced Statistical Power: Removing rows with missing data reduces the sample size, which can decrease the statistical power to detect significant relationships, potentially leading to false negatives (failing to identify a true effect).
- Violation of Assumptions: Incorrect data can violate the proportional hazards assumption, making the model's results invalid.
- Unreliable Predictions: Any predictions of survival probability or hazard based on a model built with flawed data would be unreliable and potentially dangerous in a clinical setting.

# Step 9: Cleaning the Data by Removing Redundant Columns

**Code Summary:** This section checks if the column named **Y** exists in the **df\_timevarying** DataFrame. If it does, it removes this column, as its information has been transferred to the **event** column during the time-varying data creation. Finally, it prints the columns of the DataFrame to confirm the removal.

# **Output (from PDF):**

'Y' column not found. Data is already clean.

The output indicates that the **Y** column was not found, implying it was already handled or never explicitly named **Y** in the final **df\_timevarying** DataFrame (it was likely directly named **event** during the time-varying transformation).

- Q1. Why is it important to remove the 'Y' column from the dataset before building the survival model? (Hint: Think about whether this column is still useful or if it's been replaced.)
  - Answer: It is important to remove the Y column (or any redundant column like it) because its information has been replaced and properly integrated into the event column within the time-varying dataset. Keeping both Y and event would lead to redundancy and potential confusion. The event column, along with start and stop times, is the correct and necessary format for survival models to understand whether an event occurred within a specific observation interval.
- Q2. What could happen if we keep both 'Y' and 'event' columns in the same dataset? Would that create confusion in analysis or model building? Why or why not?

- Answer: Yes, keeping both Y and event columns (if they represent the same information) in the same dataset would definitely create confusion in analysis and model building.
- Why:
  - Redundancy: They would be duplicate information, making the dataset unnecessarily larger and harder to manage.
  - **Confusion:** Analysts might accidentally use the wrong column or mix them up, leading to errors in defining the event status for the survival model.
  - Multicollinearity: In regression models like Cox, having two highly correlated (or identical) columns representing the same concept can lead to multicollinearity issues, which can destabilize the model, make coefficients unreliable, and hinder interpretation.
  - **Increased Complexity:** It adds unnecessary complexity to the data preparation and modeling pipeline.
- Q3. How can we make sure our dataset is clean before we start modeling? What checks should we perform? (Example: Missing values, duplicates, unwanted columns, etc.)
  - **Answer:** To ensure a dataset is clean before modeling, we should perform several systematic checks:
- 1. **Missing Values:** Check for **NaN** or null values using **.isnull().sum()**. Decide on imputation or removal strategies.
- 2. **Duplicate Rows/Entries:** Identify and remove duplicate rows using .duplicated().sum() and .drop\_duplicates().
- 3. **Data Types:** Verify that columns have appropriate data types (e.g., numerical for numbers, categorical for categories) using **.info()** or **.dtypes**. Convert if necessary.
- 4. **Outliers/Extreme Values:** Check for outliers using descriptive statistics (.describe()), box plots, or scatter plots. Address them by capping, transforming, or investigating their validity.
- 5. **Infinite Values:** Check for **inf** or **-inf** values using **np.isinf().sum()**.

- 6. **Redundant/Unwanted Columns:** Identify and remove columns that are duplicates, irrelevant to the analysis, or contain identifiers that shouldn't be used as features (like the **Y** column in this case).
- 7. **Consistency and Format:** Ensure data consistency (e.g., all dates are in the same format, categorical labels are spelled consistently).
- 8. **Logical Checks:** Perform domain-specific checks (e.g., age cannot be negative, diagnosis time cannot be after survival time).
  - Q4. How does removing unused or duplicate columns help improve the quality and accuracy of data analysis?
    - Answer: Removing unused or duplicate columns significantly improves the quality and accuracy of data analysis by:
      - Reducing Noise: Unused or irrelevant columns can introduce noise into the model, making it harder for the model to identify true patterns and relationships.
      - **Preventing Multicollinearity:** Duplicate columns (or highly correlated ones) can cause multicollinearity, which inflates the standard errors of regression coefficients, making them unstable and difficult to interpret. Removing them stabilizes the model.
      - Improving Model Performance: A cleaner, more focused dataset often leads to more accurate and robust models because the model isn't distracted by redundant or irrelevant features.
      - **Faster Computation:** Fewer columns mean less data to process, leading to faster model training and analysis.
      - **Easier Interpretation:** A streamlined dataset with only relevant features makes the model's output easier to understand and interpret, allowing for clearer insights.
      - Reduced Storage: Smaller datasets require less memory and storage.
  - Q5. If a teammate says, "We don't need to clean the data, models will still work," how would you respond based on your learning in this step?
    - Answer: I would strongly disagree with my teammate and explain that while a
      model might technically "run" with unclean data, its results would be

unreliable, potentially misleading, and could lead to poor decisions. I would emphasize:

- "Models are only as good as the data they're fed. Garbage in, garbage out."
- "Unclean data can lead to biased results, meaning our conclusions about treatment effectiveness or risk factors could be completely wrong, potentially harming patients or wasting resources."
- "Missing values or errors can cause the model to crash, or if it runs, it might produce inflated error rates or unstable coefficients, making it impossible to trust the insights."
- "Proper data cleaning, like removing redundant columns, ensures our model focuses on the true signals in the data, leading to more accurate predictions and interpretable results. It's a fundamental step for building trustworthy analytical solutions, especially in critical fields like healthcare."

# Step 10: Building a Cox Time - Varying Survival Model

**Code Summary:** This section builds a Cox Proportional Hazards model that can handle time-varying covariates. It initializes **CoxPHFitter**, then fits the model using the **df\_timevarying** DataFrame, specifying **start**, **stop**, and **event** columns, and including all other columns as covariates. Finally, it prints the model summary and plots the coefficients.

# **Output (from PDF):**

- cph\_time\_varying.print\_summary():
  - Shows a table similar to Step 6, but now for the time-varying model.
  - Significant p-values (p < 0.05) are observed for:</li>
    - celltype\_squamous (p=0.000)
    - celltype\_smallcell (p=0.000)
    - celltype\_adeno (p=0.000)
    - **karnofsky** (p=0.000)

- age (p=0.000)
- **prior** (p=0.000)
- trt\_2 (p=0.000)
- Coefficients (coef): (Values are very similar to Step 6, indicating consistency)

• celltype\_squamous: 0.97

• celltype\_smallcell: 0.79

• celltype\_adeno: 0.48

karnofsky: -0.03

• age: 0.01

• **prior**: 0.02

• trt\_2: -0.44

• **cph\_time\_varying.plot():** A bar plot visualizing the coefficients.

- Q1. What is the main purpose of using a time varying survival model in this project? Explain why we needed to consider changes in patient conditions over time, and how that impacts survival analysis.
  - Answer: The main purpose of using a time-varying survival model (like the Cox Time-Varying model) in this project is to accurately assess the impact of patient characteristics and treatments on survival when those factors can change over the course of the study.
  - Why consider changes: In real clinical trials, a patient's condition (e.g., Karnofsky score, disease progression) or treatment regimen (e.g., switching therapies, adding new drugs) can change over time. If we only use baseline characteristics (as in a standard Cox model), we miss the dynamic influence of these changes.
  - Impact on Survival Analysis:
    - **Increased Accuracy:** It provides a more accurate and realistic representation of the hazard over time, as it accounts for the current state of the patient and their treatment.

- **Better Insights:** It allows us to understand how a change in a covariate (e.g., a patient's Karnofsky score improving or worsening) at a specific point in time affects their subsequent risk of the event.
- Avoids Bias: Ignoring time-varying covariates can lead to biased estimates of treatment effects or risk factors, as the model would incorrectly attribute effects to baseline conditions when they are actually due to later changes.
- Q2. Which patient related factors (variables) were found to have the most significant effect on survival time? Use the summary report or the visualization to identify at least two variables that strongly influence survival, either positively or negatively.
  - **Answer:** Based on the summary report and coefficient plot (looking at p-values and magnitude of coefficients):
- 1. **celltype\_squamous (positive coefficient, p=0.000):** This variable has a large positive coefficient (0.97), indicating a strong increase in the hazard of death (shorter survival) for patients with squamous cell type cancer.
- 2. **karnofsky (negative coefficient, p=0.000):** This variable has a negative coefficient (-0.03), indicating that higher Karnofsky scores (better performance status) are strongly associated with a decreased hazard of death (longer survival).
- 3. **trt\_2 (negative coefficient, p=0.000):** This variable has a negative coefficient (-0.44), indicating that receiving Treatment 2 significantly decreases the hazard of death (longer survival) compared to Treatment 1.
  - Q3. Based on the survival model results, which treatment group (standard or test) appears to be more effective? Why? Compare the coefficients for the treatment groups and interpret what they suggest about effectiveness.
    - Answer: Based on the survival model results, **Treatment Group 2 (Test Treatment)** appears to be more effective.
    - Why: The coefficient for trt\_2 is -0.44, and its p-value is 0.000 (highly significant). Since trt\_1 is the reference group (implicitly having a coefficient of 0), the negative coefficient for trt\_2 means that being in the Treatment 2 group reduces the hazard of death compared to being in the Treatment 1 group. A lower hazard implies a longer survival time, thus indicating greater effectiveness.

- Q4. Did any variables surprise you in terms of their impact on survival? Why or why not? Reflect on the results: Were the effects of age, treatment, or other factors expected or unexpected?
  - **Answer:** (This is a subjective question, but I'll provide a common perspective based on typical medical knowledge.)

# • Expected Effects:

- age (positive coefficient): This is expected. Older patients generally have a higher risk of mortality, so a positive coefficient (increased hazard) for age is consistent with medical understanding.
- karnofsky (negative coefficient): This is also expected. A higher Karnofsky score indicates better physical functioning and quality of life, which is typically associated with better prognosis and longer survival.
- **celltype** (**positive coefficients for specific types**): It's common for different cancer cell types to have varying prognoses, so finding some cell types associated with significantly higher hazards is expected.
- trt\_2 (negative coefficient): This is a positive surprise, but the goal of a clinical trial is often to find a more effective treatment. So, finding that the test treatment significantly improves survival is a desirable and, in this context, a welcome outcome.
- **Overall:** The results align well with general medical intuition regarding prognostic factors in patient survival. No variables presented a particularly surprising or counter-intuitive impact.
- Q5. How can hospitals or healthcare researchers use these kinds of survival models in real - life situations? (Think about applications like treatment planning, patient monitoring, or healthcare policy decisions.)
  - Answer: Hospitals and healthcare researchers can use these kinds of survival models in real-life situations in several powerful ways:
- 1. **Personalized Treatment Planning:** Doctors can use the model's insights (e.g., the impact of **celltype**, **karnofsky**, **age**) to tailor treatment plans for individual patients. For example, a patient with a high-risk cell type and low Karnofsky score might be prioritized for more aggressive therapy or closer monitoring.

- 2. **Prognosis and Patient Counseling:** The model can provide more accurate prognoses (predictions of survival time) for patients, helping doctors counsel patients and their families about expected outcomes and make informed decisions about end-of-life care or treatment intensity.
- 3. **Clinical Trial Design and Evaluation:** Researchers can use these models to design future clinical trials more effectively (e.g., stratifying patients by risk factors). They can also rigorously evaluate the effectiveness of new treatments, as demonstrated by the significant positive effect of **trt\_2**.
- 4. **Resource Allocation and Healthcare Policy:** Healthcare administrators and policymakers can use survival models to understand the burden of disease and the effectiveness of interventions at a population level. This can inform decisions about allocating resources, developing public health programs, and setting treatment guidelines.
- 5. **Patient Monitoring and Early Intervention:** By identifying factors that increase hazard, hospitals can implement proactive monitoring programs for high-risk patients, potentially allowing for earlier intervention if their condition deteriorates.

# **Step 11: Plotting Survival Curves by Treatment Group**

**Code Summary:** This section plots Kaplan-Meier survival curves for each treatment group (**trt=1** and **trt=2**) using the **df\_timevarying** dataset. It iterates through the unique treatment types, fits a **KaplanMeierFitter** for each, and plots the curves on a single graph with appropriate labels and a legend.

**Output (from PDF):** A plot titled "Survival by Treatment Group" with two distinct curves:

- One for **trt=1** (Standard), which drops faster.
- One for **trt=2** (Test), which stays higher for longer. Both curves have shaded areas around them, representing confidence intervals.

- Q1. What does the survival curve tell us about the overall survival pattern of patients in the test group versus the standard group? (Hint: Compare how quickly each line drops. Which treatment shows better long - term survival?)
  - Answer: The survival curves clearly show that the Test Group (trt=2) has a better overall survival pattern compared to the Standard Group (trt=1).

- The curve for the **Standard Group (trt=1) drops more quickly and steeply**, indicating that patients in this group experience events (death) at a faster rate and have a lower probability of survival over time.
- Conversely, the curve for the Test Group (trt=2) stays higher for longer and drops more gradually, signifying that patients in this group have a higher probability of survival at any given time point and thus demonstrate better long-term survival.
- Q2. At around what time (in days) does the survival probability drop below 50% for each treatment group? (Hint: Look for the point on the x axis where each curve crosses the 0.5 mark on the y axis.)
  - Answer:
    - For the **Standard Group (trt=1)**, the survival probability drops below 50% (0.5 on the y-axis) at approximately **74 days**.
    - For the **Test Group (trt=2)**, the survival probability drops below 50% (0.5 on the y-axis) at approximately **128 days**.
- Q3. Which treatment group has a longer period where patients have a higher chance of survival? What could this suggest about the treatment?
  - Answer: The Test Group (trt=2) has a significantly longer period where
    patients have a higher chance of survival. Its curve remains above 0.5 (50%
    survival probability) for a much longer duration (128 days) compared to the
    Standard Group (74 days).
  - This suggests that the **Test Treatment is more effective** in prolonging life or delaying the event of interest. It implies that the new treatment provides a substantial clinical benefit in terms of extending the period during which patients are alive and at risk.
- Q4. What does the shaded area around each survival curve represent? Why is it important to consider this when comparing the two treatments?
  - Answer: The shaded area around each survival curve represents
    the confidence interval (typically a 95% confidence interval) for the survival
    probability at each time point.
  - Importance when comparing treatments: It is important to consider this because:

- **Uncertainty:** It quantifies the uncertainty or variability in the estimated survival curve. A wider shaded area indicates more uncertainty in the estimate, often due to a smaller number of patients at risk at later time points.
- Statistical Significance (Visual Check): If the confidence intervals of two survival curves do not overlap significantly at a given time point, it visually suggests a statistically significant difference in survival between the groups at that time. If they overlap considerably, the observed difference might not be statistically significant. In this plot, the confidence intervals appear to be largely non-overlapping, reinforcing the conclusion from the Log-Rank Test.
- Q5. Based on your observations, would you recommend one treatment over the other? Why or why not? Support your answer using insights from the chart.
  - Answer: Based on the observations from the chart, I would strongly recommend Treatment 2 (Test Treatment) over Treatment 1 (Standard Treatment).
  - Support from Chart:
    - Higher Survival Probability: The survival curve for Treatment 2
       consistently remains above the curve for Treatment 1 across the
       entire observation period, indicating a higher probability of survival at
       any given time.
    - Longer Median Survival: The 50% survival probability mark is reached much later for Treatment 2 (approx. 128 days) compared to Treatment 1 (approx. 74 days), signifying a substantial increase in median survival time.
    - Clear Separation: The two curves are clearly separated, and their confidence intervals appear to be largely non-overlapping, visually reinforcing that the difference in survival outcomes is significant and not just due to random chance.

# Step 12: Comparing Survival Using the Log - Rank Test

**Code Summary:** This section re-performs the Log-Rank Test, similar to Step 5, but explicitly using the **df\_timevarying** dataset. It separates the data by **trt** (Treatment 1 and

Treatment 2), extracts **time** (duration) and **event** (status) for each, and then applies the **logrank\_test**. The p-value is then printed.

# **Output (from PDF):**

RunCopy code

lifelines.StatisticalResult:

t\_0=-1, test\_name=logrank\_test, p=0.0009, H\_0=True, H\_A=False>

The output shows the **p-value** as **0.0009**.

- Q1. What is the purpose of performing a log rank test in this project? In your own words, explain what it helps us find out.
  - Answer: The purpose of performing a Log-Rank Test in this project is
    to statistically assess whether there is a significant difference in the
    survival experiences between two or more groups. In simpler terms, it
    helps us find out if the observed differences in how long patients survive (as
    seen in the Kaplan-Meier curves) between, say, Treatment 1 and Treatment 2,
    are real and meaningful, or if they could just be due to random variation or
    chance. It provides a formal, objective measure to support or refute visual
    observations.
- Q2. Based on the result of your test, was there a statistically significant difference in survival between the standard and test treatment groups? What was the p value?
  - Answer: Yes, based on the result of the Log-Rank Test, there was
    a statistically significant difference in survival between the standard and
    test treatment groups.
  - The p-value obtained was 0.0009.
- Q3. Which treatment group appeared to have better survival, based on the plot and the test result? Explain how you interpreted the survival curves.
  - Answer: Based on both the Kaplan-Meier plot (from Step 11) and the Log-Rank Test result, the Test Treatment group (trt=2) appeared to have better survival.

- Interpretation of Survival Curves: In the plot, the survival curve
  for trt=2 consistently remained above the curve for trt=1. This means that at
  any given point in time, a higher proportion of patients in the Test Treatment
  group were still alive compared to the Standard Treatment group. The curve
  for trt=2 also dropped more slowly, indicating a slower rate of events
  (deaths).
- Interpretation of Test Result: The Log-Rank Test's p-value of 0.0009 (which is much less than the conventional significance level of 0.05) statistically confirms that this observed difference in the curves is highly unlikely to be due to chance, thus validating the visual observation that Treatment 2 leads to better survival.
- Q4. Why do you think statistical testing like the log rank test is important before making decisions about treatments in real life?
  - Answer: Statistical testing like the Log-Rank Test is critically important before making real-life decisions about treatments because:
    - Objectivity vs. Subjectivity: Visual inspection of survival curves can
      be subjective and misleading, especially with small differences or
      noisy data. Statistical tests provide an objective, quantifiable
      measure (the p-value) to determine if an observed difference is
      statistically reliable.
    - Avoiding False Conclusions: Without statistical validation, doctors
      or researchers might mistakenly conclude that one treatment is better
      than another when the observed difference is merely due to random
      chance (a Type I error). This could lead to ineffective treatments being
      adopted, wasting resources, and potentially harming patients.
    - Evidence-Based Medicine: Modern healthcare relies on evidencebased medicine. Statistical tests provide the rigorous evidence needed to support claims of treatment efficacy, ensuring that medical practices are based on sound scientific principles.
    - Resource Allocation and Policy: For healthcare systems and
      policymakers, statistical significance is a key factor in deciding which
      treatments to fund, recommend, or include in clinical guidelines,
      ensuring that resources are directed towards interventions that
      genuinely improve patient outcomes.

- Q5. Imagine you are a health advisor. What recommendation would you make to doctors based on the survival analysis you just performed?
  - Answer: As a health advisor, based on the comprehensive survival analysis performed (Kaplan-Meier plots, median survival, and Log-Rank Test), my recommendation to doctors would be:
    - "Given the statistically significant and clinically meaningful improvement in survival demonstrated by Treatment 2 (Test Treatment) compared to Treatment 1 (Standard Treatment), I strongly recommend that Treatment 2 be considered the preferred treatment option for eligible patients. The data consistently shows that patients on Treatment 2 experience a significantly longer median survival time and a higher probability of survival at various time points. While individual patient factors should always be considered, the evidence from this analysis clearly supports the superior efficacy of Treatment 2 in prolonging patient life."

# Step 13: Build a Cox Proportional Hazards Model with Encoded Features

**Code Summary:** This section builds a Cox Proportional Hazards model using the **df\_timevarying** dataset, which already has encoded categorical features. It defines the list of features (**covariates**) to be included in the model, fits the **CoxPHFitter** using **time** as duration, **event** as status, and the specified covariates. Finally, it prints the model summary.

### **Output (from PDF):**

- cph.print\_summary():
  - Shows the summary table for the Cox model.
  - Significant p-values (p < 0.05) are observed for:</li>
    - celltype\_squamous (p=0.000)
    - celltype\_smallcell (p=0.000)
    - celltype\_adeno (p=0.000)
    - **karnofsky** (p=0.000)
    - age (p=0.000)

- **prior** (p=0.000)
- trt\_2 (p=0.000)
- Coefficients (coef):

• age: 0.01

• **karnofsky**: -0.03

• diagtime: 0.00

• **prior**: 0.02

• celltype\_squamous: 0.97

• celltype\_smallcell: 0.79

• celltype\_adeno: 0.48

• **celltype\_large**: 0.00 (reference category)

• trt\_2: -0.44

- Q1. Why did we need to convert categorical variables (like treatment or cell type) into separate columns before fitting the model? (Hint: Think about what kind of values a model can understand.)
  - Answer: We needed to convert categorical variables (like celltype, trt, prior) into separate numerical columns (using one-hot encoding, creating dummy variables) because the Cox Proportional Hazards model, like most statistical regression models, operates on numerical data. It cannot directly interpret text labels or categories.
  - What models understand: Models understand numerical relationships. If
    we were to assign arbitrary numbers (e.g., 1, 2, 3) to categories, the model
    might incorrectly assume an ordinal relationship (that 2 is "greater" than 1),
    which is not true for nominal categories. One-hot encoding creates binary (0
    or 1) columns for each category, allowing the model to treat each category as
    a distinct factor without implying any false numerical order.
- Q2. Which variable(s) in the model showed the strongest effect on survival time? How can you tell from the model summary? (Look for large positive or negative values in the summary and note their importance.)

- **Answer:** The variable(s) in the model that showed the strongest effect on survival time are those with the largest absolute coefficients and statistically significant p-values.
  - celltype\_squamous (coef = 0.97): This has the largest positive coefficient, indicating the strongest increase in the hazard of death (shortest survival).
  - 2. **celltype\_smallcell (coef = 0.79):** Also a very strong positive effect on hazard.
  - trt\_2 (coef = -0.44): This has a significant negative coefficient, indicating a strong reduction in the hazard of death (longest survival).
  - 4. **karnofsky (coef = -0.03):** While the coefficient is numerically small, its p-value is 0.000, indicating a highly significant and consistent effect where higher scores lead to lower hazard.
- **How to tell:** We look at the **coef** column in the summary table. A larger absolute value of the coefficient (regardless of positive or negative sign) indicates a stronger impact on the hazard. We also confirm their statistical significance by checking if their **p** value is less than 0.05.
- Q3. What does a positive coefficient in the Cox model indicate about a variable's relationship with survival time? What about a negative coefficient? (Try to connect this to real - world meaning: longer vs. shorter survival.)

### Answer:

- Positive Coefficient: A positive coefficient indicates that as the value of that variable increases (or if that categorical feature is present), the hazard of the event increases. In real-world terms, this means the variable is associated with a shorter survival time or a higher risk of experiencing the event (e.g., death). For example, age has a positive coefficient, meaning older age is associated with a higher hazard and thus shorter survival.
- Negative Coefficient: A negative coefficient indicates that as the
  value of that variable increases (or if that categorical feature is
  present), the hazard of the event decreases. In real-world terms, this
  means the variable is associated with a longer survival time or a
  lower risk of experiencing the event. For example, trt\_2 has a negative
  coefficient, meaning receiving Treatment 2 is associated with a lower

hazard and thus longer survival. Similarly, **karnofsky** has a negative coefficient, meaning higher Karnofsky scores (better health) are associated with lower hazard and longer survival.

- Q4. Were any variables in your model not statistically significant? Why is it
  important to know which variables are significant? (Think about how
  confidence in the results matters in medical decisions.)
  - Answer: In the provided model summary, all variables included (age, karnofsky, diagtime, prior, celltype\_squamous, celltype\_smallcell, celltype\_adeno, trt\_2) were found to be statistically significant (all had p-values of 0.000).
  - Importance of knowing significance: It is crucial to know which variables are statistically significant because:
    - Reliability of Findings: Statistical significance (p < 0.05) indicates that the observed relationship between the variable and survival is unlikely to be due to random chance. This gives us confidence that the effect is real and reliable.
    - Informed Decision-Making: In medical decisions, relying on nonsignificant variables could lead to ineffective or even harmful interventions. Doctors need to know which factors genuinely influence patient outcomes to make evidence-based choices.
    - Resource Allocation: Resources (time, money, research effort)
       should be focused on factors that have a proven, significant impact.
       Non-significant variables might not be worth further investigation or intervention.
    - Model Parsimony: Identifying non-significant variables allows for model simplification. Removing them can improve model interpretability and sometimes even performance without losing predictive power.
- Q5. Based on the model results, what are one or two insights or recommendations that might be useful for doctors or researchers studying patient survival? (Try to explain how your findings could help improve care or treatment.)
  - Answer:

- 1. Prioritize Treatment 2 and Tailor Care by Cell Type: The model strongly indicates that Treatment 2 significantly improves survival (trt\_2 has a negative and highly significant coefficient). Doctors should prioritize this treatment for eligible patients. Furthermore, the significant positive coefficients for celltype\_squamous, celltype\_smallcell, and celltype\_adeno highlight that these specific cell types are associated with a substantially higher risk of death. This insight is crucial for researchers to investigate the biological mechanisms behind these differences and for doctors to tailor more aggressive or specialized treatment strategies for patients with these high-risk cell types.
- 2. Emphasize Performance Status and Age in Prognosis: The model confirms that Karnofsky performance score (karnofsky has a negative and highly significant coefficient) and age (age has a positive and highly significant coefficient) are critical prognostic factors. Doctors can use this to provide more accurate prognoses to patients and their families. Researchers can focus on interventions that aim to maintain or improve Karnofsky scores, as this directly correlates with better survival. For older patients or those with lower performance status, this information can guide discussions about treatment intensity, supportive care, and quality of life considerations.