

Survival Analysis: Kaplan-Meier Method

Time-to-Event Analysis in Statistics

Joerg R. Osterrieder
www.joergosterrieder.com

December 20, 2025

Main Focus:

- Kaplan-Meier estimator
- Survival curves
- Log-rank test
- Practical applications

Learning Objectives:

- Understand censoring
- Calculate KM estimates
- Interpret survival curves
- Compare groups

Course Structure:

- 1 Introduction (Why survival?)
- 2 Kaplan-Meier method
- 3 Applications
- 4 Advanced topics (Appendix)

Prerequisites:

- Basic probability
- Hypothesis testing
- Data visualization

Foundation for understanding the methods covered in this section

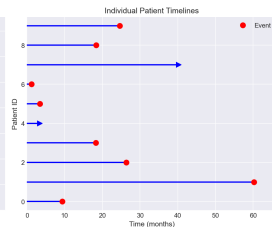
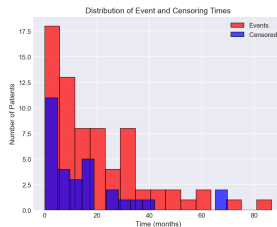
What is Survival Analysis?

Definition:

Statistical methods for analyzing time until an event occurs

The Event:

- Death
- Disease recurrence
- Machine failure
- Customer churn
- Graduation
- Employment



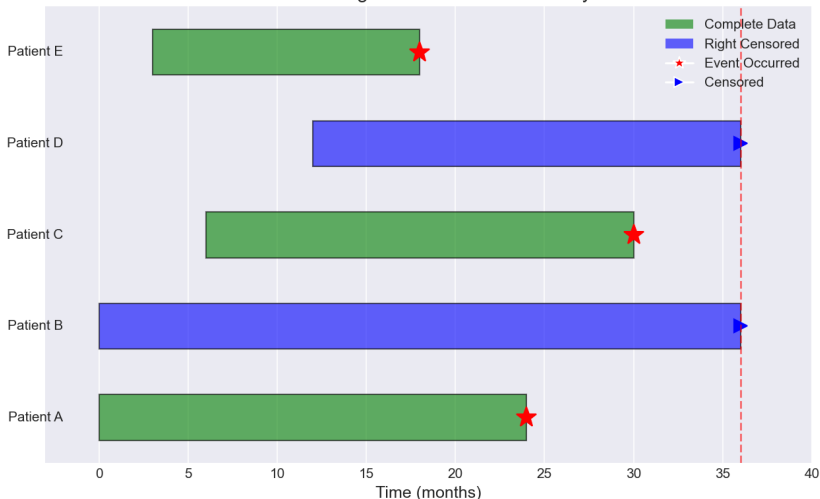
Key Feature: Not everyone experiences the event during study period!

Central Question: How long until the event occurs?

Time-to-event data requires specialized methods

Why Special Methods? The Censoring Problem

The Censoring Problem in Survival Analysis



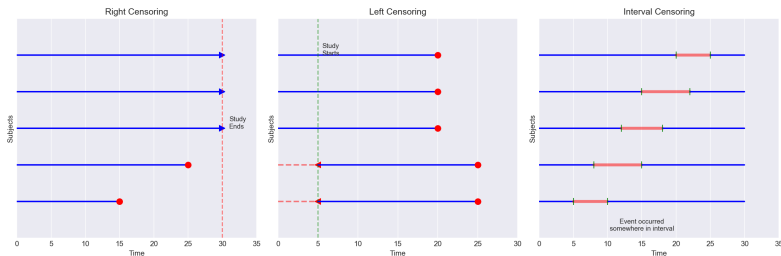
Complete Data:

- Patient A: Event at 24 months

Censored Data:

- Patient B: Still alive at 36 months
- We know: survival > 36 months

Types of Censoring



- **Right Censoring:** Event hasn't occurred by end of study (most common)
- **Left Censoring:** Event occurred before observation began
- **Interval Censoring:** Event occurred between two observations

Note: Kaplan-Meier handles right censoring

Censoring must be independent of the outcome

Time Origin (t_0):

- Diagnosis
- Treatment start
- Birth
- Purchase date

Time Scale:

- Days, months, years
- Cycles, uses
- Must be clearly defined

Event Definition:

- Must be unambiguous
- Binary (occurred/not occurred)
- Observable

At Risk:

- Subjects who could experience event
- Not yet had event
- Not yet censored

Time-to-event data requires specialized methods

Survival and Hazard Functions

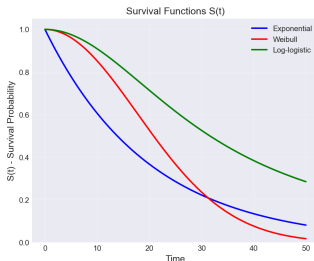
Survival Function $S(t)$:

$$S(t) = P(T > t)$$

Probability of surviving beyond time t

Properties:

- $S(0) = 1$ (all alive at start)
- $S(\infty) = 0$ (eventually all die)
- Non-increasing
- Right-continuous



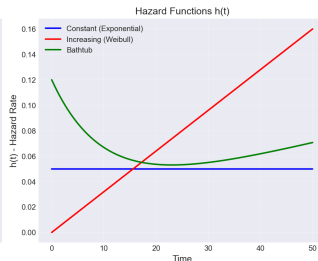
Hazard Function $h(t)$:

$$h(t) = \lim_{\Delta t \rightarrow 0} \frac{P(t \leq T < t + \Delta t | T \geq t)}{\Delta t}$$

Instantaneous risk of event at time t

Interpretation:

- Rate of event occurrence
- Given survival to time t
- Can increase, decrease, or stay constant



Time-to-event data requires specialized methods

Required Variables:

- 1 **Time:** Duration to event or censoring
- 2 **Status:** Event occurred (1) or censored (0)

Optional:

- Group/treatment
- Covariates
- Entry time

Example Dataset:

Patient	Time	Status	Treatment
1	6	1	A
2	12	0	A
3	21	1	B
4	9	1	A
5	30	0	B
6	15	1	B

Status: 1 = Event, 0 = Censored

Key Point: Both time and status are essential!

Time-to-event data requires specialized methods

Survival Analysis Methods

Non-parametric

Kaplan-Meier (FOCUS)

Life Tables

Nelson-Aalen

Semi-parametric

Cox Regression

Stratified Cox

Time-varying Cox

Parametric

Exponential

Weibull

Log-normal

Gamma

Non-parametric:

• Kaplan-Meier

MSc Statistical Methods - Survival Analysis

Semi-parametric:

• Cox regression

Survival Analysis: Kaplan-Meier Method

Parametric:

• Exponential

December 20, 2025

9 / 41

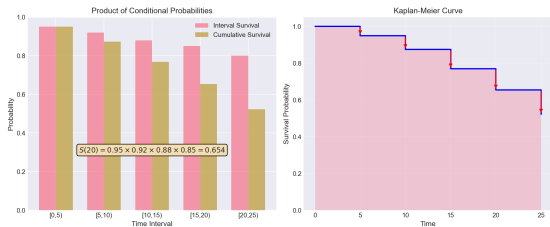
The Kaplan-Meier Approach - Intuition

Core Idea: Estimate survival probability at each event time

Key Innovation: Use conditional probabilities

Formula Intuition:

$$S(t) = P(\text{survive all intervals up to } t)$$



Product of survival probabilities for each interval

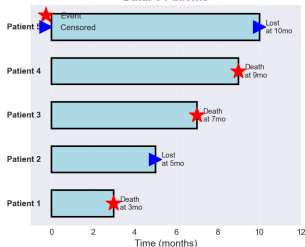
Why “Product-Limit”? $S(t)$ is product of interval survival probabilities

Kaplan-Meier handles censored observations appropriately

Simple Example: Complete KM Calculation with 5 Patients

The Kaplan-Meier Method: Complete Example

Data: 5 Patients



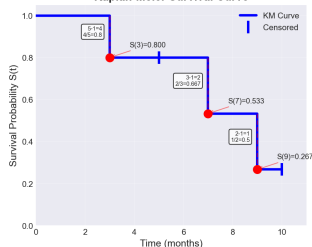
Calculations Step-by-Step

Time	n (at risk)	d (deaths)	(n-d)/n	S(t)
t=0	5	-	-	1.000
t=3	5	1	$(5-1)/5 = 0.8$	$1.0 \times 0.8 = 0.800$
t=5	4	0 (cens)	no change	0.800
t=7	3	1	$(3-1)/3 = 0.667$	$0.800 \times 0.667 = 0.533$
t=9	2	1	$(2-1)/2 = 0.5$	$0.533 \times 0.5 = 0.267$
t=10	1	0 (cens)	no change	0.267

Formula: $S(t) = S(t-1) \times (n-d)/n$

n = patients at risk, d = deaths at time t

Kaplan-Meier Survival Curve



Critical Formula Components:

- **n** = number of patients still “at risk” (being followed) at time t
- **d** = number of deaths (events) occurring exactly at time t
- Formula at each event: $S(t) = S(t_{prev}) \times \frac{n-d}{n}$
- Notice: We multiply all survival fractions together!

This pattern applies to similar real-world scenarios

Step-by-Step KM Formula Explanation

The Kaplan-Meier Formula: Detailed Explanation

Step 1: Understanding the Data & Notation

Time (t)	Event?	n (at risk)	d (deaths)	Survival Factor
t = 2	Death	5	1	$(5-1)/5 = 4/5 = 0.80$
t = 4	Death	4	1	$(4-1)/4 = 3/4 = 0.75$
t = 6	Censored	3	0	No change (censored)
t = 8	Death	2	1	$(2-1)/2 = 1/2 = 0.50$
t = 10	Censored	1	0	No change (censored)

KEY: n = number of patients still at risk before time t
d = number of deaths at time t

Step 2: Chain Multiplication Calculation

Initial: $S(0) = 1.000$

Event at t=2: $n=5, d=1$
 $S(2) = 1.000 \times (5-1)/5$
 $= 1.000 \times 0.800 = 0.800$

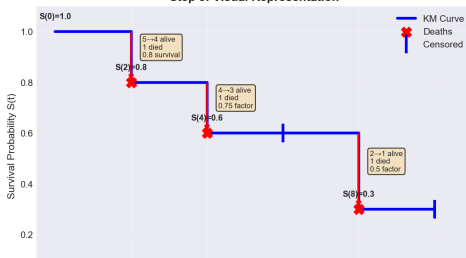
Event at t=4: $n=4, d=1$
 $S(4) = 0.800 \times (4-1)/4$
 $= 0.800 \times 0.750 = 0.600$

Censored at t=6: $n=3, d=0 \rightarrow$ NO CHANGE
 $S(6) = 0.600$

Event at t=8: $n=2, d=1$
 $S(8) = 0.600 \times (2-1)/2$
 $= 0.600 \times 0.500 = 0.300$

FORMULA: $S(t) = S(\text{previous}) \times (n-d)/n$

Step 3: Visual Representation



Step 4: Key Concepts & Interpretation

REMEMBER:

- n = patients still being followed
- d = deaths at that exact time
- Factor = $(n-d)/n$ = survival rate

KEY INSIGHTS:

1. Multiply all factors together
2. **Only update at death times**
3. **Censoring does NOT change $S(t)$**
4. Each factor is ≤ 1.0
5. $S(t)$ always decreases or stays flat

INTERPRETATION EXAMPLE:
 "At 4 months, 60% of patients are still alive"
 "The probability of surviving past 8 months is 30%"

Critical Question: When Do Censored Patients Die?

Censoring in Kaplan-Meier: What Happens to Lost Patients?

WRONG Interpretation

What we DONT do:

- Assume censored patients die immediately
- Assume they die at study end
- Assume they never die
- Make ANY assumption about their death time

These would all BIAS our estimates!

Common Misconception!

CORRECT Interpretation

What we ACTUALLY do:

- Remove censored patients from "at risk" pool
- They contribute to n UNTIL censoring time
- After censoring:
 - Not in numerator (didn't die)
 - Not in denominator (not at risk)
- Make NO assumption about their eventual outcome

KEY: Censoring is "non-informative" (independent of death risk)

Example: Patient Censored at t=6



CRITICAL INSIGHT: We never assume when censored patients die. They simply exit the risk set at censoring time. The KM method handles this uncertainty by only using available information.

The Answer: We DON'T assume anything!

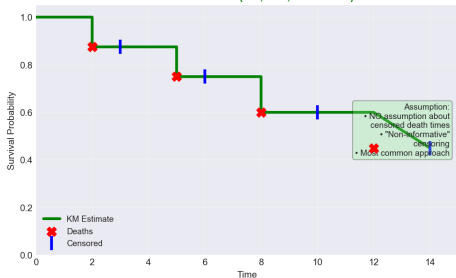
- Censored patients are **removed from risk set** at censoring time
- They contribute to denominator (n) **only until** censoring
- **No assumptions** made about their eventual outcome
- This is the "non-informative censoring" assumption

Key insight from Critical Question: When Do Censored Patients Die?

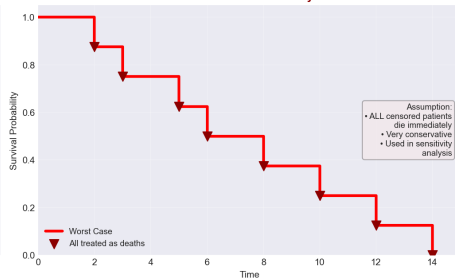
Do Other Methods Assume When Censored Patients Die?

How Different Methods Handle Censored Patients

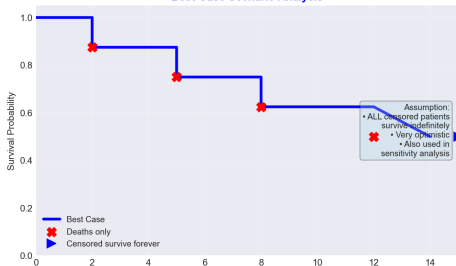
Standard Methods (KM, Cox, Parametric)



Worst-Case Scenario Analysis



Best-Case Scenario Analysis



Comparison of Approaches

Method	Censoring Assumption	When to Use
KM/Cox/Parametric	None (non-informative)	DEFAULT - primary analysis
Worst-case	All die immediately	Sensitivity - lower bound
Best-case	All survive forever	Sensitivity - upper bound
Imputation	Model-based prediction	Research/special cases

KEY INSIGHT:
 Competing **only** sensitivity analyses **make assumptions about censored patients**.
 Standard survival methods (KM, Cox) **make NO such assumptions!**

Breaking Time into Intervals:

At each event time t_i :

- n_i = number at risk
- d_i = number of events
- c_i = number censored

Interval Survival:

$$\hat{p}_i = \frac{n_i - d_i}{n_i}$$

Probability of surviving interval i

Overall Survival:

$$\hat{S}(t) = \prod_{t_i \leq t} \hat{p}_i = \prod_{t_i \leq t} \frac{n_i - d_i}{n_i}$$

Example:

- Interval 1: $\hat{p}_1 = 9/10 = 0.9$
- Interval 2: $\hat{p}_2 = 7/8 = 0.875$
- $\hat{S}(t_2) = 0.9 \times 0.875 = 0.788$

Key insight from Product-Limit Estimator Concept

The Kaplan-Meier Estimator:

$$\hat{S}(t) = \prod_{t_i \leq t} \left(1 - \frac{d_i}{n_i}\right)$$

where:

- $t_1 < t_2 < \dots < t_k$ are the distinct event times
- d_i = number of events at time t_i
- n_i = number at risk just before time t_i

Properties:

- Step function that decreases at event times
- Constant between events
- Handles ties (multiple events at same time)
- Maximum likelihood estimator
- Non-parametric (no distribution assumed)

Key insight from Mathematical Formulation

Step-by-Step Calculation Example

Data: 10 patients

Time	Status
2	1
3	1
5+	0
7	1
8+	0
11	1
12+	0
15	1
17+	0
20	1

+ indicates censored

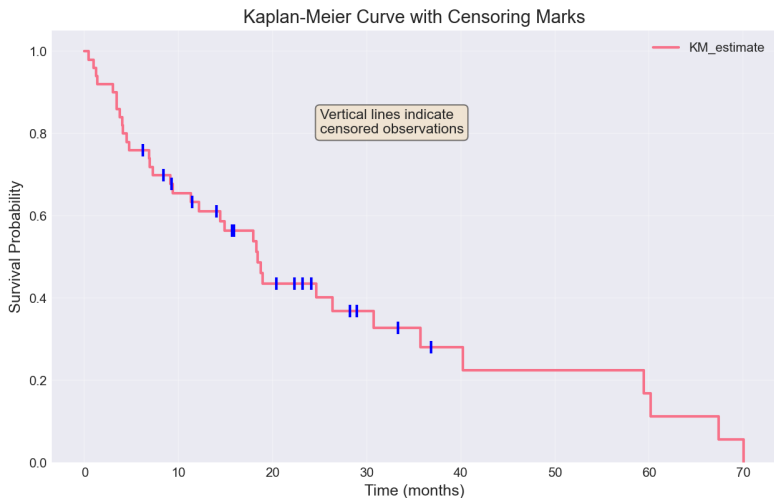
Calculation Table:

t_i	n_i	d_i	$\frac{n_i - d_i}{n_i}$	$\hat{S}(t_i)$
0	10	0	1.00	1.00
2	10	1	0.90	0.90
3	9	1	0.89	0.80
7	7	1	0.86	0.69
11	5	1	0.80	0.55
15	3	1	0.67	0.37
20	1	1	0.00	0.00

Note: Censored observations reduce n_i but not counted in d_i

This pattern applies to similar real-world scenarios

Handling Censored Observations

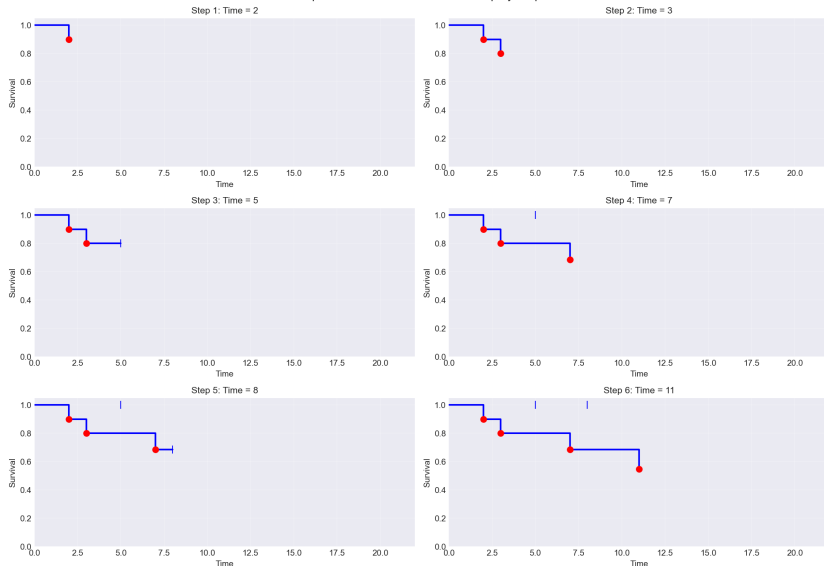


Key Rules:

- Censored subjects contribute to risk set until censored
- After censoring, removed from risk set
- Do not cause step down in survival curve

Building the Kaplan-Meier Curve

Kaplan-Meier Curve Construction Step-by-Step



Step-by-step construction:

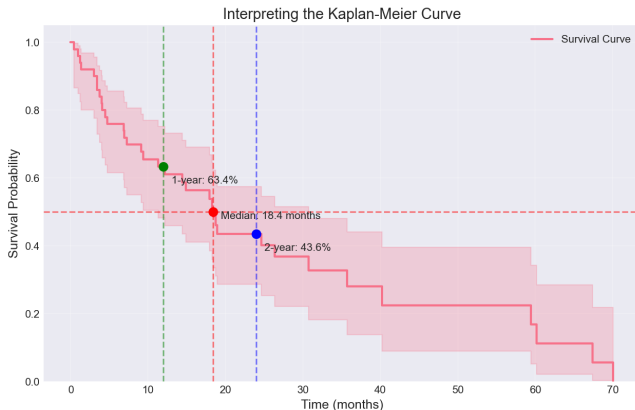
Interpretation of Kaplan-Meier Curves

Reading the Curve:

- Y-axis: Survival probability
- X-axis: Time
- Steps: Events
- Ticks: Censoring

Key Values:

- Median survival
- 1-year survival
- 5-year survival



Example: At 12 months, 65% probability of survival

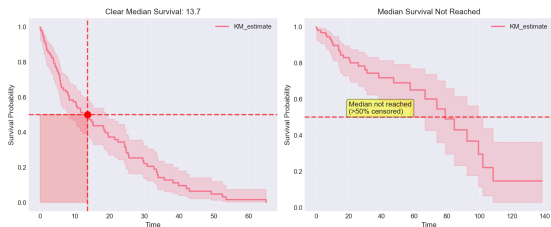
Kaplan-Meier handles censored observations appropriately

Definition: Time at which $S(t) = 0.5$

Advantages:

- Robust to outliers
- Meaningful summary
- Easy to interpret

Issue: May not be reached if $>50\%$ censored



Reading Median:

- Draw horizontal line at 0.5
- Find intersection with curve
- Read time on x-axis

Time-to-event data requires specialized methods

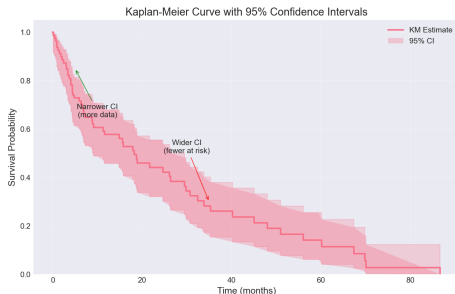
Greenwood's Formula:

$$\text{Var}[\hat{S}(t)] = \hat{S}(t)^2 \sum_{t_i \leq t} \frac{d_i}{n_i(n_i - d_i)}$$

95% CI:

$$\hat{S}(t) \pm 1.96 \times \text{SE}[\hat{S}(t)]$$

Alternative: Log-log transformation for better properties near 0 and 1



Interpretation:

- Wider at tail (fewer observations)
- Narrower early (more data)
- Accounts for censoring

Confidence intervals provide more information than p-values alone

Standard Error Calculation:

$$SE[\hat{S}(t)] = \hat{S}(t) \sqrt{\sum_{t_j \leq t} \frac{d_j}{n_j(n_j - d_j)}}$$

Example Calculation:

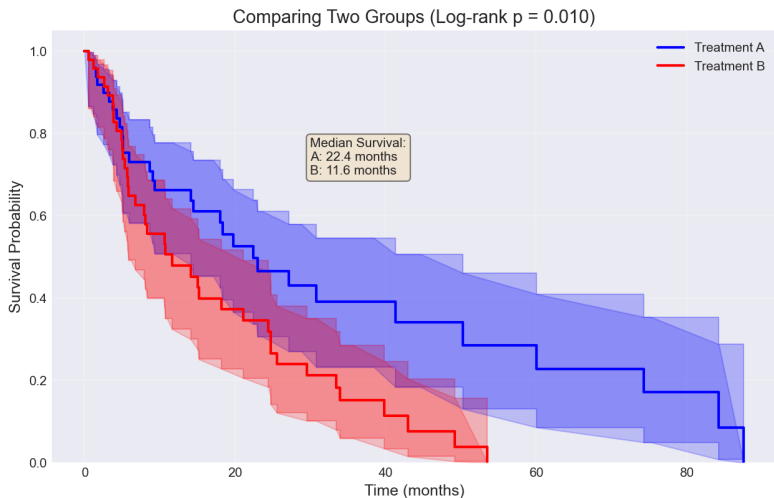
t_j	n_j	d_j	$\frac{d_j}{n_j(n_j - d_j)}$	Cumulative	SE
2	10	1	0.0111	0.0111	0.095
3	9	1	0.0139	0.0250	0.126
7	7	1	0.0238	0.0488	0.152

Properties:

- Accounts for sample size at each time
- Incorporates censoring information
- Asymptotically normal

Key insight from Greenwood's Formula - Details

Comparing Two Groups Visually



Visual Assessment:

- Separation of curves suggests difference
- Overlapping confidence bands suggest similarity
- Early vs late differences matter

Purpose: Test if survival curves are statistically different

Null Hypothesis:

$$H_0 : S_1(t) = S_2(t) \text{ for all } t$$

Alternative:

$$H_a : S_1(t) \neq S_2(t) \text{ for some } t$$

Key Features:

- Non-parametric
- Most powerful for proportional hazards
- Uses all time points
- Handles censoring

Test Statistic:

$$\chi^2 = \frac{(O - E)^2}{V}$$

where O = observed, E = expected

Key concepts that will be built upon throughout this lesson

At each event time t_i :

	Group 1	Group 2	Total
Events	d_{1i}	d_{2i}	d_i
At risk	n_{1i}	n_{2i}	n_i

Expected events in Group 1:

$$E_{1i} = n_{1i} \times \frac{d_i}{n_i}$$

Test Statistic:

$$Z = \frac{\sum_i (d_{1i} - E_{1i})}{\sqrt{\sum_i V_i}}$$

where $V_i = \frac{n_{1i}n_{2i}d_i(n_i - d_i)}{n_i^2(n_i - 1)}$

Decision: $Z^2 \sim \chi_1^2$ under H_0

Choose the appropriate test based on your data characteristics

Log-Rank Test Example

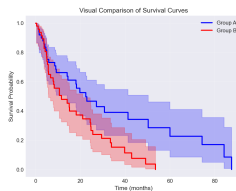
Data:

Time	n_1	d_1	n_2	d_2	E_1	V
2	10	1	8	0	0.56	0.24
4	9	0	8	1	0.53	0.25
6	9	1	7	1	1.13	0.49
8	8	1	6	0	0.57	0.24
10	7	0	6	1	0.54	0.25
Total		3		3	3.33	1.47

Calculation:

$$Z = \frac{3 - 3.33}{\sqrt{1.47}} = -0.27$$

$$\chi^2 = 0.073, \text{ p-value} = 0.79$$



Log-Rank Test Results

Test Statistic: 0.571
p-value: 0.0504
Degrees of Freedom: 1

Conclusion: Significant difference at $\alpha = 0.05$ level

Group Statistics:
Group A: 58 patients, 31 events
Group B: 58 patients, 39 events

Conclusion: No significant difference between groups ($p = 0.79$)

Choose the appropriate test based on your data characteristics

Significant Result ($p < 0.05$):

- Survival curves differ
- Groups have different prognosis
- Treatment may be effective

Non-significant ($p \geq 0.05$):

- Cannot reject equality
- May lack power
- Check sample size

Limitations:

- Tests overall difference
- Not specific about when/how
- Assumes proportional hazards
- May miss crossing curves

Report:

- Test statistic
- p-value
- Median survival by group
- Hazard ratio (if appropriate)

Log-rank test compares survival curves between groups

Assumptions:

- Independent observations
- Non-informative censoring
- Event times measured accurately
- Well-defined start time
- Clear event definition

When Violated:

- Biased estimates
- Invalid confidence intervals
- Misleading comparisons

Limitations:

- No covariate adjustment
- Assumes homogeneous groups
- Can't handle time-varying effects
- Reduced precision with heavy censoring
- May be unstable in tails

Solutions:

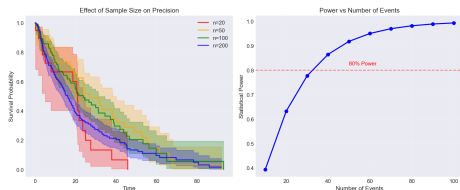
- Cox regression for covariates
- Stratified analysis
- Sensitivity analyses

Always verify assumptions before interpreting results

Factors Affecting Power:

- Number of events (not just n)
- Censoring proportion
- Effect size (hazard ratio)
- Follow-up duration
- Accrual pattern

Rule of Thumb: Need ≈ 20 -30 events per group for reasonable power



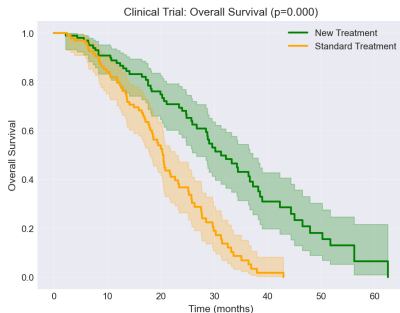
Planning Studies:

- Focus on events, not patients
- Account for loss to follow-up
- Consider interim analyses

Key insight from Sample Size Considerations

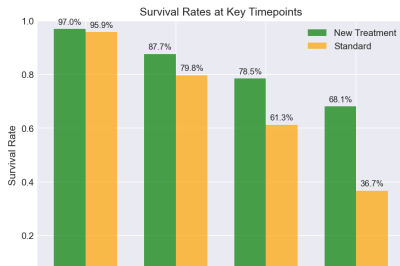
Complete Example - Clinical Trial

Clinical Trial Example: Cancer Treatment



Number at Risk:

Time (months):	0	6	12	18	24	30
New Treatment:	100	97	82	67	50	34
Standard:	100	92	72	52	26	12



Treatment Comparison Summary

Median Survival:
New Treatment: 31.3 months
Standard: 20.5 months
Difference: 10.8 months

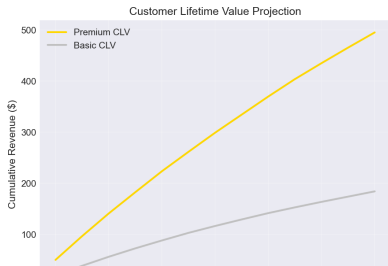
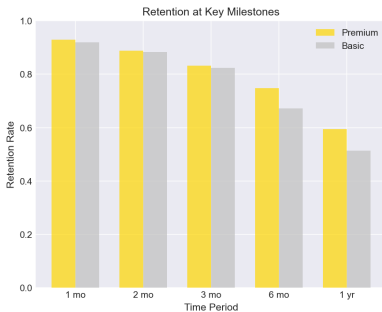
1-Year Survival:
New Treatment: 87.7%
Standard: 79.8%

Log-rank Test:
Chi-square: 31.26
p-value: 0.0000

Conclusion: Superior efficacy

Complete Example - Customer Churn

Business Application: Customer Churn Analysis



Customer Retention Insights

Median Retention Time:

Premium: 435 days
Basic: 381 days

Critical Period: First 90 days

Premium: 83.3% retained
Basic: 82.4% retained

Annual Retention:

Premium: 59.4%
Basic: 51.4%

Recommendations:

- Focus on first 30 days onboarding
- Premium has 2x better retention
- Target interventions at day 60-90

R Code:

```
surv_obj <- Surv(time, status) km_fit <- survfit(surv_obj ~ group)
ggsurvplot(km_fit, conf.int = TRUE, risk.table = TRUE, pval = TRUE)
survdif(surv_obj ~ group)
```

Python Code:

```
kmf = KaplanMeierFitter() kmf.fit(time, event_observed = status)
kmf.plot_survival_function()
results = logrank_test(time_A, time_B, event_A, event_B) print("p - value : ", results.p_value)
```

Key insight from Software Implementation

Common Mistakes:

- Ignoring censoring
- Wrong time origin
- Unclear event definition
- Informative censoring
- Overinterpreting tail
- Missing competing risks

Quality Checks:

- Verify data structure
- Check censoring patterns
- Examine assumptions
- Sensitivity analyses

Best Practices:

- Plot KM curves with CI
- Show number at risk
- Mark censoring times
- Report median survival
- Include log-rank test
- Consider stratification

Reporting:

- Follow CONSORT/STROBE
- Show survival table
- Describe censoring
- State assumptions

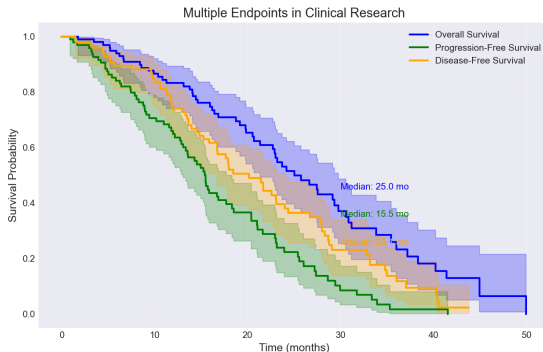
Avoid these common mistakes in your analysis

Common Uses:

- Overall survival
- Disease-free survival
- Progression-free survival
- Time to recurrence

Considerations:

- Competing risks (other causes of death)
- Protocol violations
- Loss to follow-up



Example: Comparing chemotherapy regimens

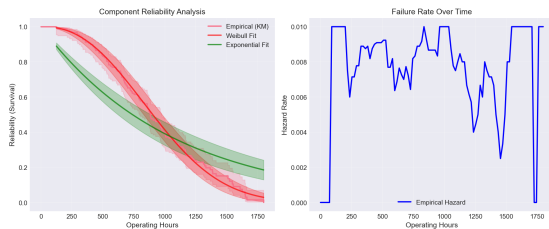
Apply these concepts to your domain-specific problems

Reliability Testing:

- Component lifetime
- System failures
- Warranty analysis
- Maintenance planning

Special Features:

- Often heavy censoring
- Accelerated testing
- Multiple failure modes

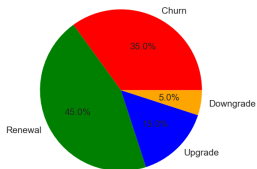


Example: LED bulb lifetime testing

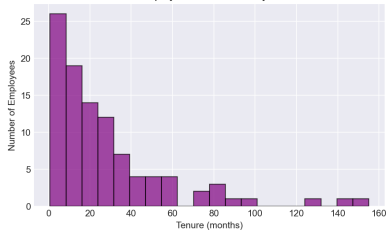
Apply these concepts to your domain-specific problems

Business Applications of Survival Analysis

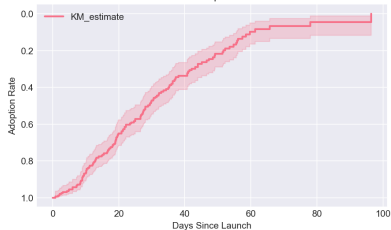
Customer Analytics Applications



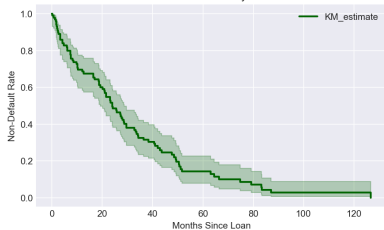
Employee Retention Analysis



Product Adoption Curve



Credit Risk Analysis



Business Applications:

- Customer lifetime value

- Contract renewal

Educational Research:

- Time to graduation
- Dropout analysis
- Program completion
- Skill acquisition

Economics:

- Unemployment duration
- Time to first job
- Business survival
- Poverty spells

Sociology:

- Marriage duration
- Time to first child
- Recidivism studies
- Social program effectiveness

Key Advantage: Handles incomplete observations common in longitudinal studies

Apply these concepts to your domain-specific problems

Essential Elements:

- Study design and duration
- Event definition
- Censoring reasons
- Number at risk over time
- Median survival (95% CI)
- Survival rates at key times
- Log-rank test results

Graphical Display:

- KM curves with CI
- Number at risk table
- Censoring marks
- Clear axis labels

Example Report: “The median overall survival was 18.2 months (95% CI: 15.3-21.5) in the treatment group versus 12.1 months (95% CI: 10.2-14.8) in the control group. The 1-year survival rates were 68% and 51%, respectively (log-rank $p = 0.012$). During the 36-month follow-up period, 45 patients (38%) were censored due to loss to follow-up ($n=12$) or administrative censoring ($n=33$).”

Key insight from Reporting Guidelines

Study Design:

- Define clear endpoints
- Plan follow-up duration
- Minimize loss to follow-up
- Consider interim analyses
- Account for competing events

Data Quality:

- Verify event dates
- Check for impossible values
- Document censoring reasons
- Handle missing data appropriately

Analysis Decisions:

- Choice of time origin
- Handling of ties
- Truncation vs censoring
- Sensitivity analyses

Interpretation:

- Clinical vs statistical significance
- Generalizability
- Selection bias
- Follow-up adequacy

Key insight from Practical Considerations

Kaplan-Meier Method:

- Handles censored data
- Non-parametric approach
- Product-limit estimator
- Visual and intuitive
- Foundation for survival analysis

Key Concepts:

- Censoring is informative
- Focus on events, not just time
- Risk set changes over time
- Confidence intervals important

Next Steps: Cox regression for covariate adjustment (see appendix)

Applications:

- Medical research
- Reliability engineering
- Business analytics
- Social sciences

Remember:

- Check assumptions
- Report comprehensively
- Consider alternatives when needed
- Interpret in context

Review these key points before moving to the next section

Advanced Topics in Survival Analysis:

- 1 Cox Proportional Hazards Model
- 2 Parametric Survival Models
- 3 Competing Risks
- 4 Time-Varying Covariates
- 5 Stratified Analysis
- 6 Sample Size Calculations
- 7 Multiple Comparisons
- 8 Interval Censoring
- 9 Cure Models
- 10 Advanced Diagnostics

Note: These topics extend beyond basic Kaplan-Meier but are important for comprehensive survival analysis

Foundation for understanding the methods covered in this section

Cox Proportional Hazards Model

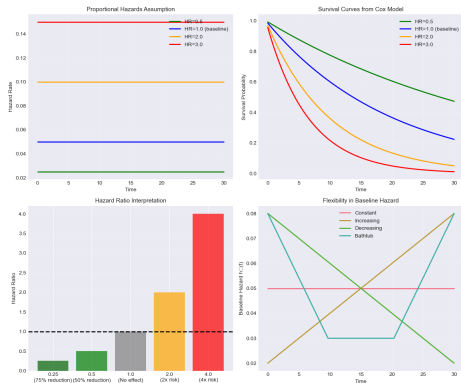
Model Form:

$$h(t|x) = h_0(t) \exp(\beta_1 x_1 + \dots + \beta_p x_p)$$

Key Features:

- Semi-parametric
- Adjusts for covariates
- Estimates hazard ratios
- No distribution assumed

Interpretation: $\exp(\beta_i)$ = hazard ratio for unit increase in x_i

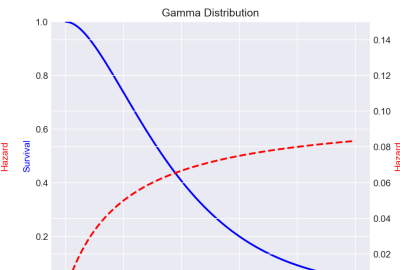
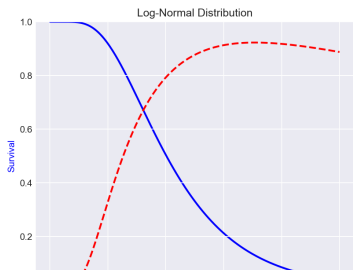
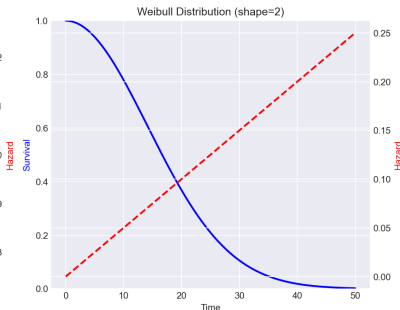
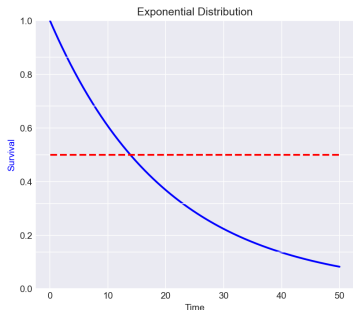


Advantages over KM:

- Multiple variables
- Continuous covariates
- More efficient
- Prediction capability

Hazard ratios quantify relative risk over time

Parametric Survival Models



Competing Risks

Problem: Multiple possible events

Example:

- Death from cancer
- Death from heart disease
- Death from other causes

Approach:

- Cumulative incidence function
- Cause-specific hazards
- Fine-Gray model



Key Point: Standard KM may overestimate event probability

Key insight from Competing Risks

Types:

- External: Predictable (e.g., age, season)
- Internal: Related to individual (e.g., biomarker levels)

Data Structure:

ID	Start	Stop	Event	Covariate
1	0	6	0	Low
1	6	12	0	Medium
1	12	18	1	High
2	0	8	0	Low
2	8	15	0	Low

Analysis: Extended Cox model with time-dependent terms

Key insight from Time-Varying Covariates

When to Use:

- Confounding variable
- Different baseline hazards
- Maintain non-parametric approach

Example: Compare treatments stratified by:

- Center in multi-center trial
- Disease stage
- Age group

Method:

- Calculate log-rank within strata
- Combine across strata
- Weighted average

Formula:

$$Z = \frac{\sum_s (O_{1s} - E_{1s})}{\sqrt{\sum_s V_s}}$$

where s indexes strata

Choose the appropriate test based on your data characteristics

Different Weights for Different Questions:

Standard Log-rank:

- Weight = 1
- Equal emphasis all times
- Best for proportional hazards

Gehan-Wilcoxon:

- Weight = n_i
- Emphasizes early differences
- Robust to outliers

Tarone-Ware:

- Weight = $\sqrt{n_i}$
- Compromise
- Moderate early emphasis

Peto-Prentice:

- Weight = $\hat{S}(t_i)$
- Similar to Gehan
- Handles censoring better

Choose the appropriate test based on your data characteristics

Freedman's Method:

Required number of events:

$$d = \frac{(z_{\alpha/2} + z_{\beta})^2(\theta + 1)^2}{\theta(\log HR)^2}$$

where:

- θ = allocation ratio
- HR = hazard ratio to detect
- α = Type I error
- β = Type II error

Total Sample Size:

$$n = \frac{d}{\text{Pr(event)}}$$

Considerations:

- Accrual period
- Follow-up duration
- Loss to follow-up rate

Key insight from Sample Size Calculation Details

Problem: Testing multiple groups inflates Type I error

Solutions:

- Bonferroni correction
- Holm-Bonferroni
- False Discovery Rate
- Pre-planned comparisons

Example: 3 groups = 3 pairwise comparisons

Bonferroni: Use $\alpha/3 = 0.017$ for each test

Better Approach:

- Overall test first
- If significant, then pairwise
- Adjust for multiple testing

Key insight from Multiple Comparisons

Problem: Event time known only within interval

Examples:

- Disease detected at periodic checkup
- Equipment inspected monthly
- Survey conducted annually

Data Structure:

ID	Left	Right	Status
1	0	6	Interval
2	6	12	Interval
3	12	∞	Right
4	18	18	Exact

Methods:

- Turnbull estimator
- Parametric models
- Multiple imputation

Censoring must be independent of the outcome

Concept: Some subjects will never experience event

Model:

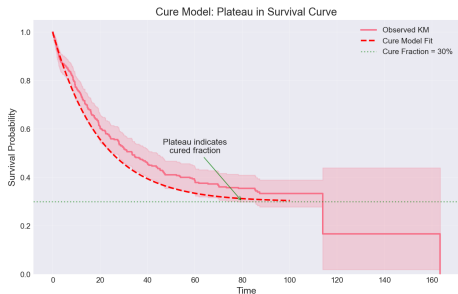
$$S(t) = \pi + (1 - \pi)S^*(t)$$

where:

- π = cure fraction
- $S^*(t)$ = survival for uncured

Applications:

- Some cancers
- Customer churn
- Credit default



Identifying Cure:

- Plateau in survival curve
- Long follow-up needed
- Biological plausibility

Key insight from Cure Models

Random Effects in Survival:

Purpose: Account for unobserved heterogeneity

Model:

$$h(t|Z) = Z \cdot h_0(t) \exp(\beta X)$$

where Z is random frailty term

Applications:

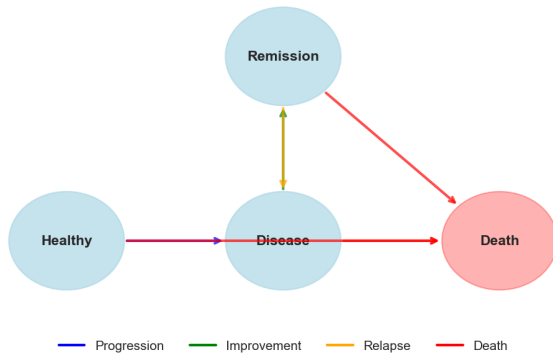
- Clustered data (patients within hospitals)
- Recurrent events
- Family studies
- Multi-center trials

Distributions:

- Gamma (most common)
- Log-normal
- Positive stable

Key insight from Frailty Models

Multistate Model: Disease Progression



Complex Event Processes:

- Multiple states (healthy \rightarrow ill \rightarrow dead)
- Transitions between states
- Reversible transitions possible

Applications:

Cox Model Diagnostics:

- Schoenfeld residuals (PH assumption)
- Martingale residuals (functional form)
- Deviance residuals (outliers)
- Score residuals (influence)

Goodness of Fit:

- Concordance index
- AIC/BIC
- Likelihood ratio test
- Calibration plots



Validation:

- Cross-validation
- Bootstrap
- External validation
- Time-dependent ROC

Always check diagnostics before trusting model results

When Standard Methods Fail:

Bootstrap Procedure:

- 1 Resample data with replacement
- 2 Calculate KM estimate
- 3 Repeat B times (e.g., 1000)
- 4 Use percentiles for CI

Advantages:

- No distributional assumptions
- Works for complex statistics
- Better coverage in small samples
- Can handle median survival

Types:

- Percentile method
- BCa (bias-corrected accelerated)
- Bootstrap-t

Confidence intervals provide more information than p-values alone

Foundational Texts:

- Kaplan & Meier (1958) - Original paper
- Klein & Moeschberger - Survival Analysis Techniques
- Collett - Modelling Survival Data in Medical Research
- Hosmer, Lemeshow & May - Applied Survival Analysis

Advanced Topics:

- Therneau & Grambsch - Modeling Survival Data
- Hougaard - Analysis of Multivariate Survival Data
- Ibrahim et al. - Bayesian Survival Analysis

Software Resources:

- R: survival, survminer, flexsurv packages
- Python: lifelines, scikit-survival
- SAS: PROC LIFETEST, PROC PHREG
- Stata: sts, stcox commands

Consult primary sources for detailed methodology