

Survival analysis for healthcare data

M2 Données massives en santé

Juliette Murriss

Januray 2025

Outline

Competing risks

Context

Non-parametric estimation

Semi-parametric estimation

Recurrent events

Introduction

Non-parametric approach

Semi-parametric approaches

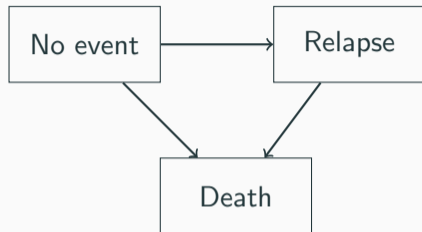
Parametric approaches

Further survival models for survival problems

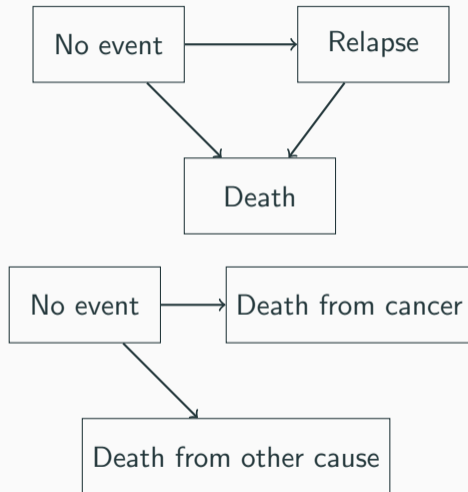
Conclusion

Competing risks

Context



Context



Context

- Subjects may be at risk of experiencing multiple different events
- Examples:
 - Non-independent events: relapse and death
 - Mutually exclusive events: different causes of death

Definition

A competing risk is another event that either:

- Prevents the occurrence of the event of interest, or
- Fundamentally alters its probability of occurrence

Cumulative incidence

Main objective: Estimate the cumulative incidence function

$$F_1(t) = Pr(T \leq t, \epsilon = 1)$$

- $\epsilon = 1$: Represents the **event of interest** (e.g., death due to a specific cause or disease recurrence).
- $\epsilon \neq 1$: Represents **competing events** (e.g., death due to other causes or alternative types of failure).
- $\epsilon = 0$: Sometimes used to indicate **censored observations** (i.e., no event occurred before the end of follow-up).

This is the proportion of patients at time t who experienced the event of interest, accounting for competing events

Two non-parametric approaches:

- **1-KM** (Kaplan-Meier):
 - Treats competing events as censored
 - Often overestimates true incidence
- **Kalbfleisch-Prentice** (1980):
 - Incorporates competing events information
 - More accurate when events aren't independent

Event Coding:

- For 1-KM method:
 - 1: Event of interest
 - 0: Censored or competing event
- For Kalbfleisch-Prentice method:
 - 1: Event of interest
 - 0: Censored
 - 2: Competing event

Important Note

Choice of coding affects interpretation and results!

In practice

```
1 bladder_v2 <- bladder1[!duplicated(bladder1$id,  
2     fromLast = F),  
3     c("id", "treatment", "stop", "status", "number")]  
4 bladder_v2$recidive <- ifelse(bladder_v2$status == 1, 1,  
5     ifelse(bladder_v2$status > 1, 2, 0))  
6 bladder_v2 <- subset(bladder_v2, treatment != "pyridoxine")  
7 bladder_v2$treatment <- factor(bladder_v2$treatment)
```

In practice

```
1 par(mfrow = c(1,3))
2
3 # KM
4 KM <- plot(survfit(Surv(stop, recidive) ~ treatment, data =
   bladder_v1), col = c("blue","red"), xlab = "Temps", ylab = "
   Survie", main = "KM", bty = "l")
5 legend("topright", c("Placebo", "Thiotepa"),
6       col = c("blue", "red"), lty = c(1,1), bty = "n")
7
8 # 1 - KM
9 KM1 <- plot(cuminc(bladder_v1$stop, bladder_v1$recidive,
   bladder_v1$treatment), col = c("blue", "red"), lty = c(1,1),
   xlab = "Temps", ylab = "Incidence cumulee", main = "1 - KM"
  )
```

```
1 # Kalbfleisch et Prentice
2 KM2 <- plot(cuminc(bladder_v2$stop, bladder_v2$recidive,
  bladder_v2$treatment), col = c("blue", "red", "blue", "red")
  , lty = c(1,1,2,2), xlab = "Temps", ylab = "Incidence
  cumulee", main = "Kalbfleish et Prentice")
```

Comparison of cumulative incidence curves

Cumulative incidence curves are compared using Gray's test

Unlike the Log-rank test, the curves are allowed to cross

```
1 cuminc(bladder_v2$stop, bladder_v2$recidive, bladder_v2$  
  treatment)
```

P-value > 0.05 → no significant difference between the cumulative incidence curves across treatment groups.

Fine and Gray Model (FG)

The Fine and Gray model focuses on modeling the **subdistribution hazard function**, which corresponds to the instantaneous risk at time t of experiencing the event of interest, given that it has not occurred before:

$$\alpha_1(t|Z) = \lim_{\Delta t \rightarrow \infty} \frac{P(t \leq T < t + \Delta t, \epsilon = 1 | T \geq t \cup (T \leq t \cap \epsilon \neq 1))}{\Delta t} \quad (1)$$

- **Censoring:** Includes only "true" censoring (patients lost to follow-up or still event-free at the end of the study).
- **Event:** Only the event of interest is considered.

The cumulative incidence associated with this model is calculated using the method of Kalbfleisch and Prentice.

Fine and Gray Model (FG)

The FG model assumes **proportional hazards** for the subdistribution hazard:

$$\alpha_1(t|Z) = \alpha_{01}(t) \exp(\beta Z) \quad (2)$$

Here, $\alpha_{01}(t)$ is the baseline subdistribution hazard, and βZ represents the covariate effect

Cause-Specific (CS) Method

Competitive risk analysis using the traditional Cox model: **Cause-specific hazard analysis.**

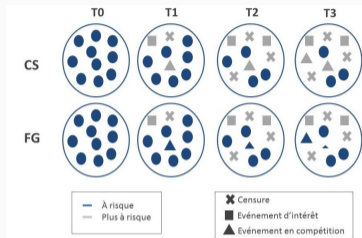
- **Censoring:** Includes both patients experiencing competing events and "true" censoring.
- **Event:** Focuses solely on the event of interest.

This method effectively analyzes the **net survival**, which assumes a hypothetical world where only the event of interest can occur

→ Violates the assumption of **non-informative censoring** since competing risks carry relevant information.

→ Requires specialized methods to account for this violation. The cumulative incidence for the cause-specific method is computed using **1-Kaplan-Meier (1-KM)**.

Difference Between FG and CS Methods: Definition of At-Risk Population



The primary distinction lies in how the **at-risk population** is defined:

- **CS**: Patients who are neither censored nor have experienced the event of interest or competing events.
- **FG**: Patients who are not censored, including those who have already experienced the competing event.
→ **The contribution of patients with competing events to the risk set decreases over time.**

FG Model: In practice

To perform competing risks analysis using the Fine and Gray model, two variables are required:

- A **categorical variable** to indicate the type of event (e.g., event of interest or competing event).
- A **time-to-event variable** to represent survival duration.

Example in R:

```
1 cov1 <- model.matrix(~ factor(treatment) + number,  
2   data = bladder_v2)[, -1]  
3 summary(crr(bladder_v2$stop, bladder_v2$recidive, cov1 = cov1))
```

This code fits a Fine and Gray proportional hazards model using the 'crr' function from the 'cmprsk' package. The categorical variable 'treatment' and numerical covariate 'number' are included in the model matrix.

Assumptions of the Fine and Gray Model

Key assumptions of the Fine and Gray model include:

- **Non-informative censoring:** Assumes that censoring is unrelated to the risk of the event.
(Note: This assumption is difficult to verify directly.)
- **Log-linearity:** Assumes that covariate effects are linear on the log scale.
(This can also be challenging to confirm directly.)
- **Proportional hazards:**
 - **Graphical Method:** Schoenfeld residuals can be plotted against time. Proportionality holds if the mean residuals are constant over time and close to 0.
 - **Lin's Test:** Based on the cumulative sum of residuals. Proportionality holds if the test yields a p-value > 0.05 .

Testing the PHA: In practice

Approaches for testing the proportional hazards assumption

Using Schoenfeld Residuals:

```
1 par(mfrow=c(1,2))
2
3 mod2 <- crr(bladder_v2$stop, bladder_v2$recidive, cov1 = cov1)
4 for(j in 1:ncol(mod2$res)){
5     scatter.smooth(mod2$uft, mod2$res[,j],
6                   main = names(mod2$coef)[j],
7                   xlab = "Time",
8                   ylab = "Schoenfeld Residuals")
9 }
```

Testing the PHA: In practice

Using Lin's Test (Package: crskdiag):

```
1 diag_crr(Crsk(stop, recidive) ~ treatment + number,  
2         data = bladder_v2, test = "prop", seed = 1234)
```

Question: Draw conclusions from both approaches.

Handling Violations of the Proportional Hazards Assumption

If the proportional hazards assumption is not satisfied:

- Include a **time-dependent covariate** in the model.

Example: Model the covariate `number` as a quadratic function of time:

$$\beta_1 \cdot \text{number} + \beta_2 \cdot \text{number} \cdot t + \beta_3 \cdot \text{number} \cdot t^2$$

Question: Adapt the R code to include this time-dependent variable and interpret the results.

Competing Risks: A Research Frontier

Competing risks analysis is an active area of research, with many tools and methodologies available. Examples include:

- **riskRegression** Package: For advanced risk prediction using regression models.
- **Multi-state Models:** Comprehensive analysis of complex event histories.

Modeling Approaches

1. Fine-Gray Model (FG)

- Models subdistribution hazard:

$$\alpha_1(t|Z) = \alpha_{01}(t)\exp(\beta Z)$$

- Censoring = true censoring only
- Event = event of interest

2. Cause-Specific Hazards (CS)

- Uses classical Cox model
- Censoring = competing events + true censoring
- Models net survival

Summary

	Single Event	Competing Risks
Function	Survival	Cumulative Incidence
Test	Log-rank	Gray's test
Model	Cox	Fine-Gray/CS
Measure	Hazard Ratio	Subdistribution HR

Key Takeaways

- Consider competing events in analysis
- Choose appropriate method based on research question
- Verify model assumptions

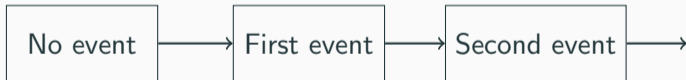
Recurrent events

Context

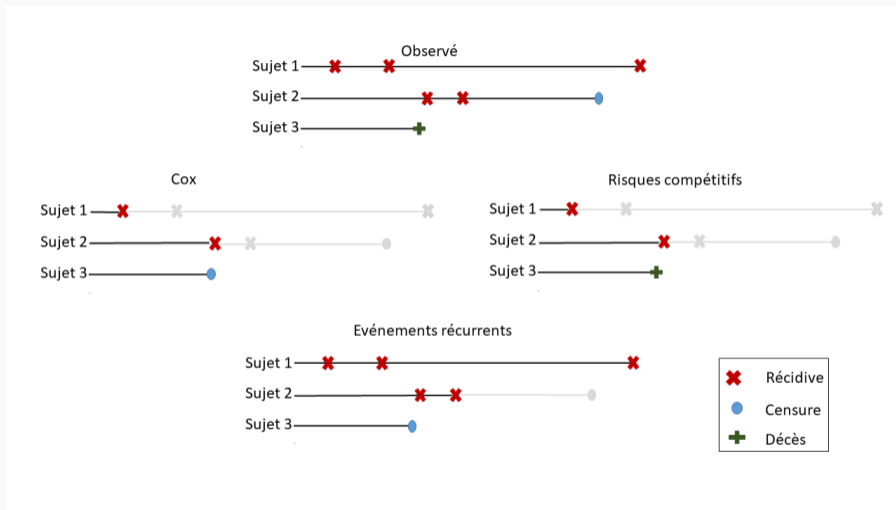
In clinical studies, the primary outcome may involve the occurrence of an event. This event can happen multiple times for the same patient.

- Asthma attacks, disease progression, (re)hospitalizations, etc.

Definition Stochastic processes that generate events of the same type repeatedly over time.



Context



Challenges

- **Intra-subject correlation:** Multiple events within the same patient can weaken their health status, potentially altering (increasing or decreasing) the likelihood of subsequent events.
- **Inter-patient heterogeneity:** Patients followed for a longer period are at higher risk of experiencing more events compared to those followed for a shorter period.

Cox models are no longer appropriate:

- They only focus on the time to the first event → loss of information.
- Covariates may have different effects depending on whether it is the first, second, or subsequent event.

→ Recurrent events require dedicated methods.

When studying recurrent events, the key scientific questions include:

- Does the treatment reduce the number of events compared to the control?
- How many events does the treatment prevent compared to the control?
- What is the treatment's effect on the number of subsequent events (e.g., starting from the 3rd event) compared to the control?
- What is the treatment's effect on the number of future occurrences for patients who have already experienced the event?

Measures of effect

- The cumulative number of events over the study period.
- An event rate, per unit of time.
- Time to the first event.
- Time between successive events.

Mean Cumulative Function (MCF)

Let $N(t)$ denote a counting process, the *MCF* is defined as: $\mu(t) = \mathbb{E}\{N(t)\}$

- $dN_i(t)$: An increment of N_i over the small interval $[t, t + dt]$.
- $Y_i(t)$: An indicator of whether individual i is at risk during $[t, t + dt]$.
- $Y(t) = \sum_{i=1}^n Y_i(t)$: Total number of individuals at risk during $[t, t + dt]$, with n as the number of subjects.
- $dN(t) = \sum_{i=1}^n Y_i(t)dN_i(t)$: Total number of events observed during $[t, t + dt]$.
- $t_{(1)}, t_{(2)}, \dots, t_{(T)}$: The T distinct event times for the n individuals.

The Nelson-Aalen estimator of the MCF is:

$$\hat{\mu}(t) = \sum_{h|t_h \leq t} \frac{dN(t_h)}{Y(t_h)} \quad (3)$$

Time Intervals

- **Gaptime**

- Time between two successive events.
- Lower bound is 0.
- Upper bound corresponds to the time between two successive events.

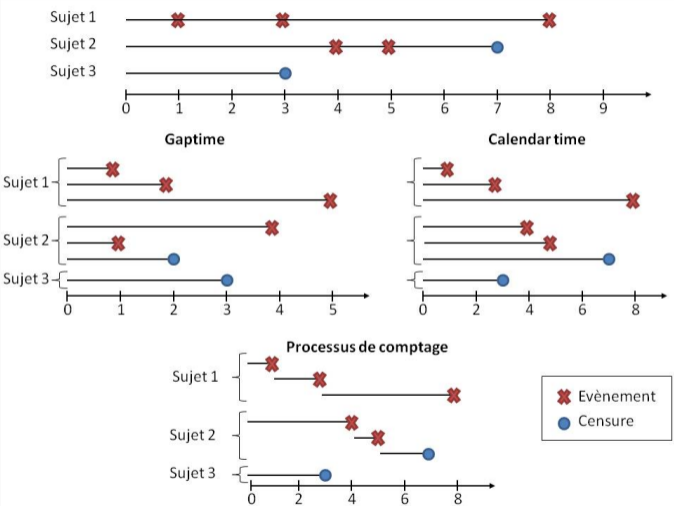
- **Calendar Time**

- Time from the start of the study until the occurrence of the event.
- Lower bound is 0.
- Upper bound is the time from the start of the study until the event occurs.

- **Counting Process**

- Similar to calendar time, but a subject may enter the study late or be censored.
- Lower bound is the time of the previous event.
- Upper bound corresponds to the time of the current event.

Time intervals



At-risk sets

- **Unrestricted Set:** All time intervals for each subject are considered at risk for any event, regardless of how many events the subject has had.
- **Restricted Set:** Only includes time intervals corresponding to the k^{th} event for subjects who have already had $k - 1$ events.
- **Semi-Restricted Set:** Related to an instantaneous baseline risk for each event, i.e., the semi-restricted set contains subjects who have had $k - 1$ or fewer events for the k^{th} event.

$$r_i(t) = Y_i(t)r_0(t) \exp(\beta Z_i(t)) \quad (4)$$

- Semi-parametric approach (no assumptions on the baseline hazard function).
- Extension of the Cox model (proportional intensity model).
- Time scale: Counting process.
- Risk set: Unrestricted.
- Baseline hazard function: Not stratified.
- Parameter estimation: Maximum likelihood.
- Outputs event intensity or hazard ratio (HR) for recurrent events.
- Assumes independence between events.

Data needs to be reshaped:

- Patient with no events: One row.
- Patient with multiple recurrent events and no follow-up after the last event: One row per event.
- Patient with multiple recurrent events and follow-up after the last event: One row per event plus one row for follow-up.

Conditional Models – Andersen-Gill in practice

```
1 summary(coxph(Surv(start, stop, event) ~ rx + number + cluster(id), data = bladder2))
```

The 'cluster' option uses robust variance, which corrects for the dependency between events from the same subject.

Question: What is the effect of the variable `rx`? How should it be interpreted?

Conditional Models – Andersen-Gill

Advantages:

- Inter-event dependence can be addressed using a time-dependent covariate.
- Focuses on the overall treatment effect.
- Allows analysis of all hospitalizations across all subjects.

Limitations:

- Assumes independence between events within the same subject, although robust variance can be used.
- Correlation within the same subject is explained by covariates.
- The proportional hazards assumption may be too strong in practice: constant HR over time and across events.

$$r_i(t) = Y_i(t)r_{0j}(t) \exp(\beta Z_i(t)) \quad (5)$$

- Extension of the Cox model.
- Time scale: Counting process or gaptime.
- A subject is not at risk for the k^{th} event until they have had the previous event.
- Stratified baseline hazard function.
- Parameter estimation: Maximum likelihood.
- Analyzes the time between events via statistical tests or HR for each ordered event time.

Using the same data as for the Andersen-Gill model:

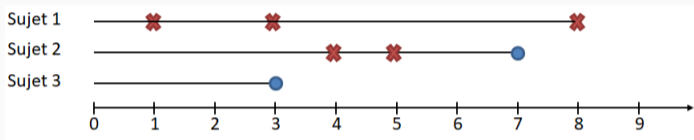
```
1 summary(coxph(Surv(start, stop, event) ~ rx + number + cluster(id) + strata(enum), data = bladder2))
```

$$r_i(t) = Y_i(t)r_{0j}(t) \exp(\beta Z_i(t)) \quad (6)$$

- Extension of the Cox model.
- Time scale: Calendar time.
- Risk set: Semi-restricted.
- Stratified baseline hazard function.
- Parameter estimation: Maximum likelihood.

Marginal Models – WLW in practice

Data needs to be reshaped: If the maximum number of events per patient is 3, then each patient will have 3 rows.



Id	start	stop	recurrence	num observation
1	0	1	1	1
1	0	3	1	2
1	0	8	1	3
2	0	4	1	1
2	0	5	1	2
2	0	7	0	3
3	0	3	0	1
3	0	3	0	2
3	0	3	0	3

Marginal Models – WLW in practice

```
1 summary(coxph(Surv(stop, event) ~ rx + number + cluster(id) +  
  strata(enum), data = bladder))
```

Interpretation: The probability of the k^{th} recurrence from randomization in all patients (whether or not they had the $(k - 1)^{\text{th}}$ recurrence) is lower in $rx = 1$ compared to $rx = 0$. i.e., the delay of the k^{th} recurrence from randomization is longer for $rx = 1$ than for $rx = 0$ among those who had the $(k - 1)^{\text{th}}$ recurrence.

Advantages:

- Preserves randomization.
- Accounts for intra-patient dependence in the variance estimation.
- Aims to assess the cumulative effect of treatment on events from randomization onwards.

Limitations:

- Cannot analyze all events, as very few patients experience a high number of events.
- Requires defining a maximum number of events per patient.
- More suitable for analyzing multiple types of events rather than recurrent events.
- Need to limit the number of events per patient.

- **Generalized Linear Model:**
- Focuses on the number of events per patient.

$$\log(Y) \sim \alpha + \beta Z$$

where

- Y is the random variable representing the number of events, $Y \sim P(\lambda)$
- β is the parameter vector to be estimated.
- Z is the matrix of covariates.

Parametric approaches – Poisson in practice

```
1 summary(glm(recur ~ treatment + number + offset(log(stop)),  
             data = bladder_poisson, family = "poisson"))
```

Interpretation of coefficients:

- Intercept: $\exp(\text{coef}) =$ average number of events per month in the placebo group with 0 initial tumors.
- rx: Relative difference in the average number of events per month between placebo and thiotepa groups. The average number of events per month in the thiotepa group is $\exp(-0.93) = 60\%$ lower than in the placebo group.
- number: For each unit increase in the `number` variable, the average number of events per month increases by $\exp(0.34) = 40\%$.

Parametric approaches – Poisson in practice

The "offset" option accounts for the varying follow-up durations; patients followed for longer periods can have more recurrences than those followed for shorter periods.

The interpretation of results differs with and without the offset:

- With offset: rate (average number of events per time unit).
- Without offset: average number of events.

Checking for overdispersion assumption: mean = variance.

```
1 require(AER)
2 dispersiontest(glm(recur ~ treatment + number + offset(log(stop
   )), data = bladder_poisson, family = "poisson"), alternative
   = c("greater"))
```


Parametric approaches – Poisson

- **Advantages:**

- Focuses on event rates.
- Simple: total number of events / total follow-up time in each group.
- Assumes all events are independent.
- Aims to assess the overall effect of treatment on the average number of events.

- **Limitations:**

- Assumes that past events do not influence the occurrence of future events.
- Overdispersion assumption: mean = variance, leading to the use of a negative binomial model if violated.

Parametric approaches – Negative Binomial

- Main assumption: Events within a single individual are dependent.
- Each individual has their own event rate following a Poisson distribution.
- The set of Poisson rates follows a Gamma distribution.

```
1 require(MASS)
2 summary(glm.nb(recur ~ treatment + number + offset(log(stop)),
  data = bladder_poisson))
```

The interpretation remains the same as for the Poisson model.

Parametric approaches – Negative Binomial

- **Advantages:**

- Focuses on event rates.
- Does not assume overdispersion.
- Aims to assess the overall effect of treatment on the average number of events.

- **Limitations:**

- Assumes that past events do not impact the occurrence of future events.

Further survival models for survival problems

Frailty Models – Context

Traditional methods assume that **populations are homogeneous and observations are independent** of each other. However:

- This assumption is often unrealistic, as there may be unobserved important covariates
Example: environmental or genetic factors that influence survival times.
- Data may exhibit correlation structures
Example: effects associated with a specific center in multicenter studies, or repeated events for the same patient (e.g., recurrent hospitalizations).

Frailty models address these challenges by incorporating both heterogeneity in the data and dependencies between event times.

Frailty Models

Frailty models are commonly used to account for dependence between event times of certain individuals.

Consider $T_{ik} = \min(X_{ik}, C_{ik})$, where i indexes the i th individual in group k (with $k = 1, \dots, G$). The hazard function for individual i in group k is given by:

$$h_{ik}(t|Z_{ik}, \omega_k) = h_0(t)\omega_k \exp\left(\sum_{j=1}^p \beta_j Z_{jik}\right) \quad (7)$$

- $h_0(t)$: baseline hazard function
- Z_{jik} : vector of covariates for individual i in group k
- β_j : regression coefficient for covariate j
- ω_k : frailty term for group k , modeling shared random effects within the group.

Interpretation of ω_k :

- If $\omega_k > 1$: Individuals within the same group tend to experience the event **more quickly** than predicted by a model without dependence.
- If $\omega_k < 1$: Individuals within the same group tend to experience the event **less quickly** than predicted by a model without dependence.

Frailty Models – Assumptions and properties

- **Independence:** Observations are conditionally independent given the frailty term ω_k .
- **Distribution of random effects:** Frailty terms are often assumed to follow a Gamma or Gaussian distribution, with variance θ :
 - θ measures the heterogeneity between groups.
 - **Interpretation:** A larger θ indicates greater variability between groups.
- **Proportional hazards:** The proportional hazards assumption holds conditionally on the frailty term ω_k . This means:
 - The regression coefficients β are interpreted conditionally on the frailty.
 - For example, if Z_{ik} is binary (0 or 1), then e^β represents the risk ratio between an individual coded 1 and an individual coded 0, **within the same group**.

Frailty Models – In practice

```
1 bladder_v1$centre <- ifelse(bladder_v1$stop < 20, 1,  
2   ifelse(bladder_v1$stop >= 20 & bladder_v1$stop < 40, 2,  
3   ifelse(bladder_v1$stop >= 40 & bladder_v1$stop < 60, 3, 4))  
   )
```

In this example, a new variable `centre` is created to define subgroups based on the time variable `stop`. Four centers are defined:

- Center 1: `stop < 20`
- Center 2: `20 <= stop < 40`
- Center 3: `40 <= stop < 60`
- Center 4: `stop >= 60`

Frailty Models – In practice

```
1 summary(coxph(Surv(stop, recidive) ~ treatment + number +  
frailty(centre), data = bladder_v1))
```

Interpreting the results:

- The "gamma" values in the model output represent the estimated random effects for each subgroup (center).
- **Question:** Do we observe a significant center effect?

Model selection:

- Models with and without the random effect can be compared using the Akaike Information Criterion (AIC).
- The preferred model minimizes the AIC, indicating the best trade-off between goodness-of-fit and model complexity.

Joint Models – Context

- Joint models for time-to-event and longitudinal data combine two types of outcomes:
 - **Time-to-event (survival)** outcome: e.g., time to death, time to recurrence.
 - **Longitudinal outcome**: e.g., repeated measurements of a biomarker, symptom scores, etc.
- These models allow for the analysis of the relationship between the survival process and the longitudinal outcome.
- The 'frailtypack' package in R implements joint models that use a shared latent process (frailty) to account for correlation between the two outcomes.
- They are particularly useful when longitudinal measurements are used to inform the survival outcome.

Joint Models

- A joint model consists of two main components:
 - **Longitudinal submodel:** Models the evolution of the longitudinal outcome over time.
 - **Survival submodel:** Models the time-to-event data (e.g., time to death or disease progression).
- Both submodels share a **latent random effect (frailty)**, which accounts for the correlation between the two outcomes.
- Example: Let $Y(t)$ be the longitudinal outcome at time t , and let T be the time-to-event outcome.
- The survival model is typically a Cox proportional hazards model or a parametric survival model, while the longitudinal model can be a linear mixed model or a nonlinear mixed model.

Joint Models – Longitudinal and Survival Submodels

- The **longitudinal submodel** usually takes the form:

$$Y_i(t) = \beta_0 + \beta_1 X_i(t) + b_i + \epsilon_i(t)$$

where:

- $Y_i(t)$: Longitudinal outcome for individual i at time t .
 - $X_i(t)$: Covariates (e.g., time-varying treatments) for individual i .
 - b_i : Random effect (frailty) that links the longitudinal and survival outcomes.
 - $\epsilon_i(t)$: Measurement error.
- The **survival submodel** often assumes a Cox model or parametric survival model:

$$h_i(t) = h_0(t) \exp(\gamma Z_i + \delta b_i)$$

where:

- $h_i(t)$: Hazard function for individual i at time t .
- Z_i : Covariates influencing the survival outcome.
- b_i : Latent frailty term shared between the longitudinal and survival models.

Joint Models – In practice

```
1 library(frailtypack)
2
3 # Fit a joint model with a Cox survival model and a linear
  longitudinal model
4 fit <- jointModel(
5   lme(fixed = Y ~ time + treatment, random = ~1|patientID,
6       data = longitudinal_data),
7   coxph(Surv(time, status) ~ age + treatment, data = survival
8         _data)
9 )
10 # Summary of the joint model fit
summary(fit)
```

- The lme function is used to fit the longitudinal submodel, and coxph fits the survival submodel.

Joint Models – Interpretation

- **Latent frailty term:** The shared frailty term, denoted by b_i , captures the correlation between the longitudinal and survival outcomes.
 - A significant frailty term suggests that the longitudinal outcome has an impact on the survival outcome.
- **Cox model output:** The coefficients of the Cox model (in the survival submodel) show the effect of covariates on the survival hazard.
- **Longitudinal model output:** The coefficients of the longitudinal model show how covariates (e.g., time-varying treatments) affect the longitudinal outcome.
- If the frailty term is significant, it indicates that the longitudinal marker (e.g., biomarker level) is predictive of the time-to-event outcome (e.g., time to disease progression).

Joint Models

Advantages:

- **Captures dependence** by accounting for the correlation between survival and longitudinal outcomes.
- **Time-varying effects**

Limitations:

- **Complexity** to specify, estimate, and interpret compared to separate models.
- **Assumptions:** Shared frailty assumption might not always hold, requiring careful modeling decisions.
- **Interpretation:** Frailty terms and the relationship between outcomes can be challenging to interpret.

But there are many more...

- Combinations of all models we have seen
- Multi-state models
- Interval censoring
- ...

Conclusion

Conclusion

- Always seek more information about the study objectives: What questions are we trying to answer?
- Explore your data before starting the analysis:
 - Do we have all the necessary data to answer the research questions?
 - What are the characteristics of the data?
 - Do these characteristics match the ones described in the study?
- Check the assumptions of the models.
- Verify the sample size, and do this **at every stage** of the analysis.
- Compare your results with the literature / consult clinical expertise.