

Colon Cancer Survival Analysis

Calder Glass, Kotaro Ito, Robin Zhan

2025-06-01

Introduction

We explore the *colon* dataset found in the R “survival” package, containing data observed from a clinical trial where stage B/C colon cancer patients receive adjuvant chemotherapy. 929 independent patients—484 men and 445 women—were randomly assigned between two treatments and a control group: Levamisole, Levamisole and Fluorouracil, and control (denoted as *Lev*, *Lev+5FU*, and *Obs* in the dataset). “Levamisole is a low-toxicity compound that was originally used to treat worm infestations in animals”, while “5-FU is a moderately toxic chemotherapy agent” used to treat cancer.¹

Patients were then observed until one of two events occurred: recurrence or death (denoted as “1” and “2” in its respective order under column *etype*). The time of occurrence, in days, was then recorded to later investigate and determine whether or not different treatments were effective in keeping the patients alive. Each patient in the dataset, identified by their *id*, has two rows for both recurrence and death. The status column indicates whether or not the event occurred or not (“0” indicates no and “1” indicates yes). If a patient has been recorded for 3000 days for both recurrence and death and the status remains 0 for both, it signifies that they did not experience any event for 3000 days and dropped out of the study for unknown reasons. Figure 1 and 2 below represents the Kaplan-Meier Survival Curve after splitting the dataset by *etype* (recurrence and death). The convex shape of Figure 1 conveys that many recurrences occur early on while the curve for death events show that deaths in patients are gradual and consistent.

Figure 1: KP for Recurrence Events

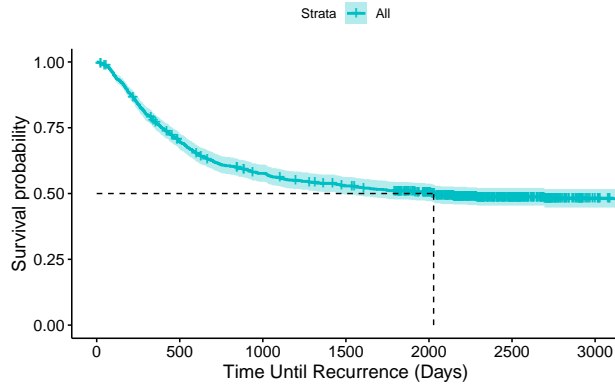
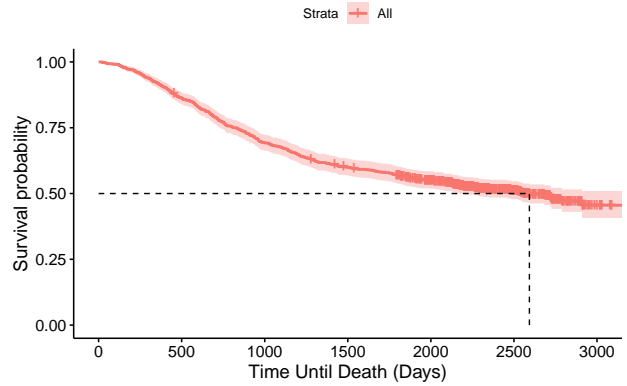


Figure 2: KP for Death Events



The dataset contains the *id*, *age* (in years), *sex* (“1” indicates man and “0” indicates woman), *rx* (the treatment type or control), *obstruct* (“1” indicates a colon obstructed by a tumor and “0” indicates no obstruction), *perfor* (“1” indicates a perforated colon and “0” indicates no perforation), *adhere* (“1” indicates cancer adhering to other organs and “0” indicates no adherence), *nodes* (the number of lymph nodes with colon cancer), *time* (time until event occurrence or censoring), *status* (whether or not the event occurred or

¹colon: Chemotherapy for Stage B/C colon cancer. Colon Cancer. Retrieved May 1st, 2025, from <https://rdrr.io/cran/survival/man/colon.html>

not), **differ** (“3” indicates quickly growing cancer, “2” indicates moderate growth, and “1” indicates slowly growing and less likely to spread), **extent** (describes the spread of the tumor and ranges from 1-4, where “1” indicates that the tumor is limited to the inner lining of the colon and “4” indicates invasion of tumor to nearby organs and tissues), **surg** (“1” indicates a long time between initial surgery and registering to the study while “0” indicates a short time), **node4** (“1” indicates a patient has more than four positive lymph nodes and “0” indicates four or less), and **etype** (recurrence or death event) of each patient.¹

Taking all of the covariates we listed above into consideration, our aim is to determine whether or not a specific treatment has a significant effect on the survival of the patient. Our secondary objective is to assess which covariate(s) have a significant effect on the hazard risk. In the course of the analysis, we omit observations with N/A values, reducing our final dataset to 888 independent patients. A five percent significance level (0.05) will be used to balance the risk of false positives to detect meaningful effects.

Model Fitting

With the clinical context established and the relevant covariates explained, we begin to evaluate the effects of treatment and other factors on the patient. Given that our dataset includes two types of events — recurrence of cancer and death — we begin by modeling these outcomes separately using marginal Cox proportional hazards models. This allows us to estimate the hazard associated with each covariate for each event type independently and by fitting separate Cox models for recurrence and death, we can assess whether specific treatments or patient characteristics are associated with an increased or decreased risk for each type of event.

Marginal Model: Recurrence

Kaplan-Meier Estimate for Recurrence

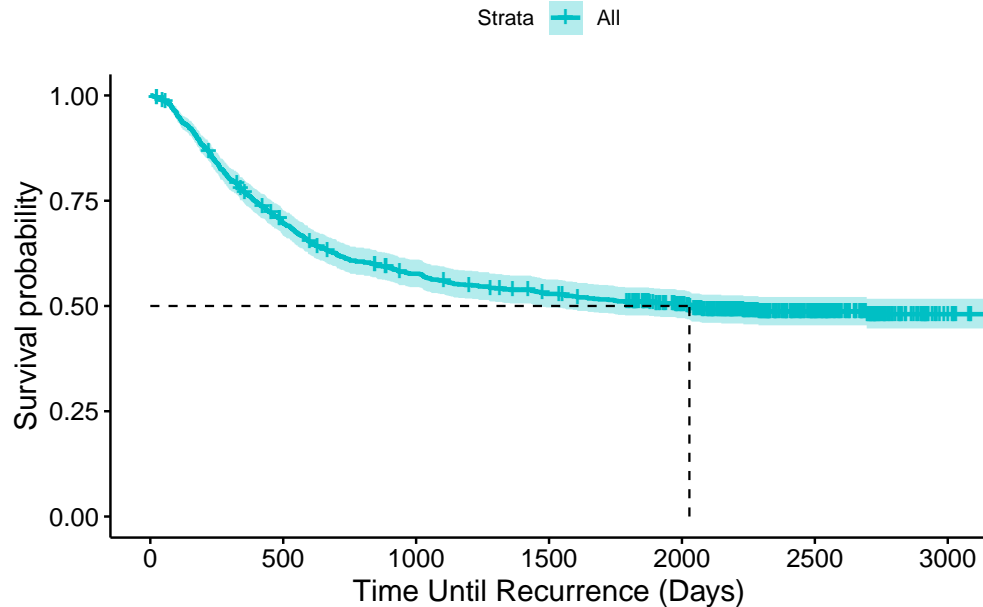
```
# data for recurrence event
recurrence_data <- colon[colon$etype == 1, ]

# survival object
surv_object_recurrence <- Surv(time = recurrence_data$time,
                               event = recurrence_data$status)

# Fit the Kaplan-Meier model for recurrence events
km_fit_recurrence <- survfit(surv_object_recurrence ~ 1)

# Plot the Kaplan-Meier survival curve for recurrence events with a title
ggsurvplot(km_fit_recurrence,
            data = recurrence_data,
            xlab = "Time Until Recurrence (Days)",
            palette = "#00BFC4", # Customize the color
            title = "Figure 1: KP for Recurrence Events",
            surv.median.line = 'hv',
            break.time.by=500)
```

Figure 1: KP for Recurrence Events



```
recurrence_median = surv_median(km_fit_recurrence)
print(recurrence_median)
```

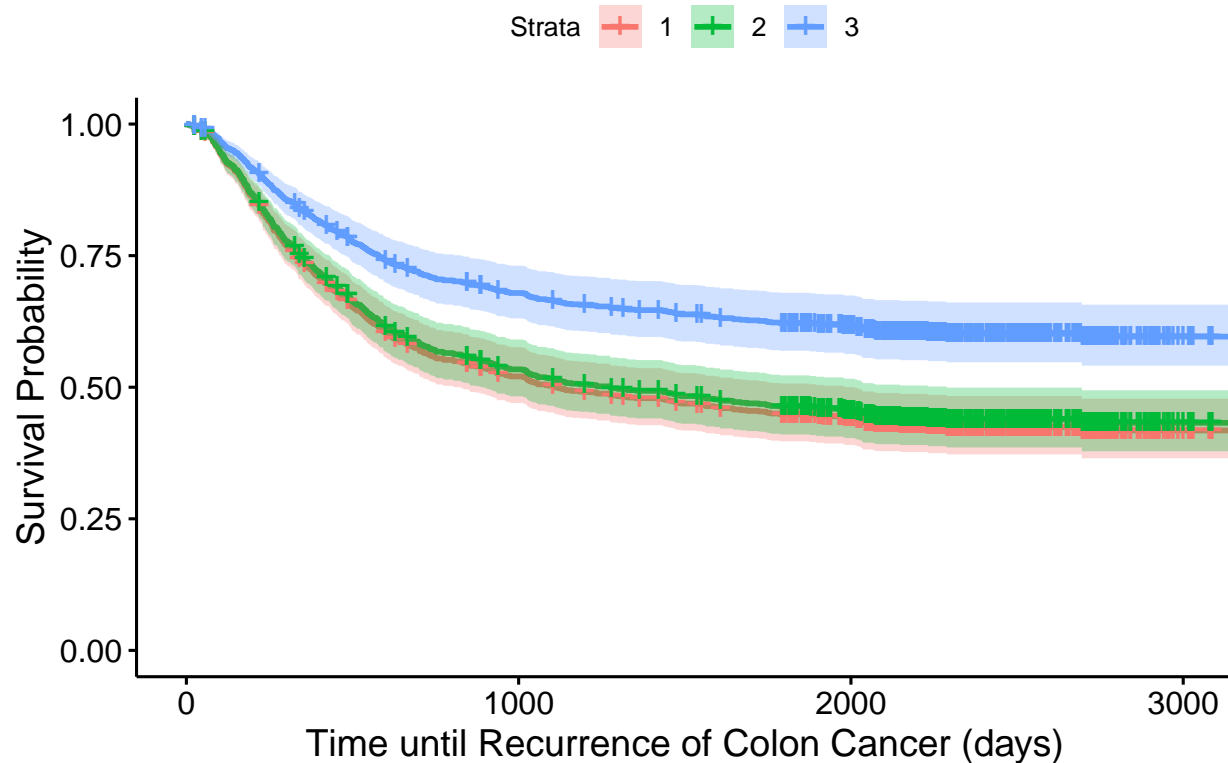
```
## strata median lower upper
## 1 All 2028 1471 NA
```

Most of the decay in the survival probability for the Kaplan-Meier curve occurs in the first 1500 days. The estimate hits the median survival probability at 2028 days.

The thick density of tick marks after approximately 1750 days suggest that many subjects were censored, either passing away after the study or dropping out before the study concluded.

```
recurrence_obj = coxph(Surv(time, status) ~ rx, data = recurrence_data)
recurrence_fit = survfit(recurrence_obj,
                          newdata = data.frame(rx= c("Obs", "Lev", "Lev+5FU")))
ggsurvplot(recurrence_fit, data = recurrence_data,
            conf.int = TRUE,
            title = "Comparison of Survival Functions for Different Treatments",
            xlab = "Time until Recurrence of Colon Cancer (days)",
            ylab = "Survival Probability")
```

Comparison of Survival Functions for Different Treatments



Median Survival Time for each Treatment

```
recurrence_cox_median = surv_median(recurrence_fit)
print(recurrence_fit)
```

```
## Call: survfit(formula = recurrence_obj, newdata = data.frame(rx = c("Obs",
##      "Lev", "Lev+5FU"))))
##
##      n events median 0.95LCL 0.95UCL
## 1 888   446   1114     805    2012
## 2 888   446   1275     891    2288
## 3 888   446    NA      NA      NA
```

```
summary(recurrence_obj)
```

```
## Call:
## coxph(formula = Surv(time, status) ~ rx, data = recurrence_data)
##
##      n= 888, number of events= 446
##
##              coef exp(coef) se(coef)      z Pr(>|z|)
## rxLev        -0.03996  0.96083  0.10967 -0.364   0.716
## rxLev+5FU    -0.52289  0.59280  0.12118 -4.315 1.6e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
```

```
##          exp(coef) exp(-coef) lower .95 upper .95
## rxLev      0.9608      1.041    0.7750    1.1912
## rxLev+5FU   0.5928      1.687    0.4675    0.7517
##
## Concordance= 0.556 (se = 0.013 )
## Likelihood ratio test= 23.15 on 2 df,  p=9e-06
## Wald test          = 21.5 on 2 df,  p=2e-05
## Score (logrank) test = 21.95 on 2 df,  p=2e-05
```

From the initial plot, the survival curve for Levamisole & Fluoracil appears to be greater than the other treatments, so the survival probabilities for the patients on Levamisole + Fluoracil are higher than those on the other treatments.

This would suggest that Levamisole & Fluoracil could improve survival probabilities for patients.

The other two survival curves, only Levamisole and observed/control, have very similar survival curves with huge overlaps between their confidence intervals. This would suggest that only taking Levamisole does not have a particular significant boost on the survival probability of colon cancer for its patients. The strong overlap is also reflected in the median survival times, which are less than 200 days apart for only Levamisole and the control group. The group taking both Levamisole and Fluoracil does not have a recorded median survival time because its survival curve doesn't decrease to the 50% mark within the allocated time of the study.

The 95% confidence interval for the hazard ratio of Levamisole + Fluoracil is (0.4675, 0.7517) and the hazard ratio itself is 0.5928, only further suggesting that the combined treatment has a significant role in increasing survival probability to colon cancer for its patients. Meanwhile, the 95% confidence interval for the hazard ratio only Levamisole is (0.775, 1.1912) and the hazard ratio itself is 0.9608, so there is very little difference between taking or not taking only Levamisole treatment in terms of the hazard rate.

Exploring Covariates for the Marginal Recurrence Model

AIC

For the AIC tests, the covariates `study` and `id` are not included. The `id` covariate is the same as the observation number, it does not have contextual significance to the event of relapse or death from colon cancer. The `study` covariate is not included as all of the subjects are from the same study.

The model with the `node4` covariate, the binary variable for whether the patient had more than 4 positive lymph nodes, had the lowest AIC.

Since the `nodes` covariate and `node4` covariate are closely related, the `nodes` covariate will be skipped.

Therefore, forward selection proceeds with the above covariate.

The model with the `extent` covariate, the description of the local spread of the tumor, had the lowest AIC.

Therefore, forward selection proceeds with the above covariate.

The model with the `surg` covariate, the time from initial surgery to registration in the study, had the lowest AIC.

Therefore, forward selection proceeds with the above covariate.

The model with the `differ` covariate, the description of the removed cancer cells from the colon, had the lowest AIC.

Therefore, forward selection proceeds with the above covariate.

```

# Level 5:
# list of covariates to put into the models
recurrence_covariates5 = c("adhere", "obstruct", "perfor", "sex", "age")

# building a model per covariate by pasting the given covariate into the formula

# the set_names function helps to clear up which AIC value corresponds to which
# model when performing the AIC function

recurrence_models5 = map(recurrence_covariates5, \(v)
  coxph(as.formula(paste("Surv(time, status) ~ rx +
                        node4 + extent + surg +
                        differ + ", v)),
  data = recurrence_data)) |>
  set_names(recurrence_covariates5)

aic_lvl5 = map_dbl(recurrence_models5, AIC) |>
  sort()

aic_lvl5

```

```

## obstruct   adhere      sex   perfor      age
## 5641.966 5643.228 5643.275 5643.552 5644.507

```

The model with the `obstruct` covariate, the binary variable for whether the cancer had adhered to other organs, had the lowest AIC.

Therefore, forward selection proceeds with the above covariate.

None of the AICs shown above are less than the previous model, so the chosen model has the following covariates: `obstruct`, `surg`, `extent`, `node4`, and `differ`.

However, `obstruct` and `differ` only decreased the AIC by less than 2, which is the standard a significant AIC gain.

Thus, the model will only use `rx`, `surg`, `extent`, and `node4`. Next, the model was summarized in order to conclude relationships between the different levels of covariates, such as the treatment covariate and the differentiation covariate.

Full Coxph Model for Recurrence

```

summary(coxph(Surv(time, status) ~ rx + node4 + extent + surg, data = recurrence_data))

```

```

## Call:
## coxph(formula = Surv(time, status) ~ rx + node4 + extent + surg,
##       data = recurrence_data)
##
##      n= 888, number of events= 446
##
##              coef exp(coef) se(coef)      z Pr(>|z|)
## rxLev      -0.03010   0.97034  0.10996 -0.274   0.7843
## rxLev+5FU -0.49267   0.61099  0.12158 -4.052 5.07e-05 ***

```

```
## node4      0.83883    2.31366    0.09877    8.493 < 2e-16 ***
## extent     0.53551    1.70833    0.11865    4.513 6.38e-06 ***
## surg       0.23389    1.26351    0.10383    2.253 0.0243 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
##          exp(coef) exp(-coef) lower .95 upper .95
## rxLev      0.9703      1.0306      0.7822      1.2037
## rxLev+5FU   0.6110      1.6367      0.4814      0.7754
## node4       2.3137      0.4322      1.9065      2.8078
## extent      1.7083      0.5854      1.3539      2.1556
## surg        1.2635      0.7914      1.0309      1.5487
##
## Concordance= 0.656 (se = 0.013 )
## Likelihood ratio test= 120 on 5 df,  p=<2e-16
## Wald test              = 121.2 on 5 df,  p=<2e-16
## Score (logrank) test = 126.6 on 5 df,  p=<2e-16
```

From the likelihood ratio test, the overall p-value is less than $2e - 16$, which is much less than the critical value/significance level of 0.05.

rx: The hazard rate for patients who took the treatment with just Levamisole is only 2.97% less hazardous than taking no treatment at all. Those who took Fluoracil in addition to Levamisole benefited with a hazard ratio of 0.611, 38.9% less hazardous than no treatment at all.

node4: Patients who had more than 4 positive lymph nodes had over double the hazard rate of those who didn't.

extent: As the spread of the tumor developed from muscles to contiguous structures, the hazard ratio to those who only had submucosa development increased to as high as 2.3137 times as likely to suffer a recurrence of colon cancer.

surg: Patients with a long time from their initial surgery to registration in the study had a 26.35% greater hazard rate than those with a shorter time interval.

Marginal Model: Death

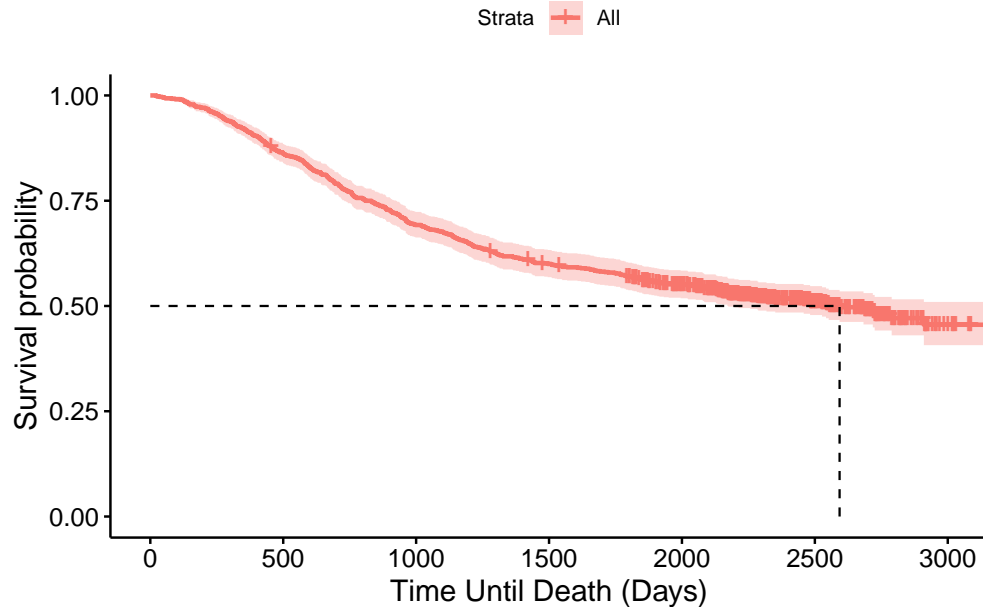
Kaplan-Meier Estimate for Death

```
# Data set for death event
# Filter the data for death event
colon_death <- colon_clean[colon_clean$etype == 2, ]

# Fit the Kaplan-Meier model for the death event
km_fit_death <- survfit(Surv(time, status) ~ 1, data = colon_death)

# Plot the Kaplan-Meier curve for death events
ggsurvplot(km_fit_death, data = colon_death,
            title = "Kaplan-Meier Survival Curve for Death Events",
            xlab = "Time Until Death (Days)",
            surv.median.line = 'hv',
            break.time.by = 500)
```

Kaplan–Meier Survival Curve for Death Events



```
# Median survival time
Death_med <- surv_median(km_fit_death)
print(Death_med)
```

```
##      strata median lower upper
## 1      All    2593    2174    NA
```

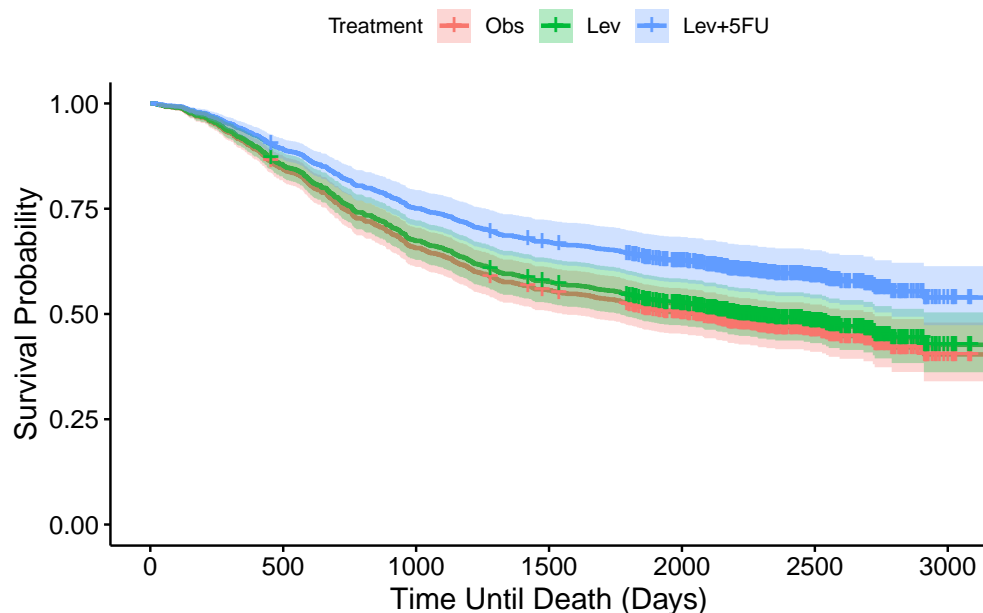
The Kaplan-Meier curve declines slowly and almost linearly over the 3000 days follow-up and the median survival time is 2593 days. At the beginning of the study before one year (365 days), the survival probability is roughly above 90%, which indicate that patients in the study begin with a near-perfect chance of remaining alive. Also, the numerous tick marks in the late tail indicate that many individuals were censored alive at the later stage of the study, which is common because it is hard to follow-up for a long period.

Proportional Hazard Model for Death

```
# Coxph model with rx(treatment) as covariate
cox_death <- coxph(Surv(time, status) ~ rx, data = colon_death)
# Create fit for different treatment
fit_death <- survfit(cox_death, newdata = data.frame(rx = c("Obs", "Lev",
                                                           "Lev+5FU")))

# Plot fit for coxph
ggsurvplot(fit_death, data = colon_death, conf.int = TRUE,
            ylab = "Survival Probability",
            xlab = "Time Until Death (Days)",
            title = "Coxph of Death Event by Treatment",
            legend.title = "Treatment",
            legend.lab = levels(colon_death$rx),
            break.time.by = 500)
```


Coxph of Death Event by Treatment



```
# median survival time
cox_med <- surv_median(fit_death)
print(cox_med)
```

```
##   strata median lower upper
## 1      1   2052  1550  2718
## 2      2   2257  1767    NA
## 3      3     NA  2789    NA
```

This plot shows the survival curves for three treatments after fitting a Cox model with only treatment (rx) as the covariate. Lev+5FU (blue) curve is located above the other two curves, which might indicate that Levamisole+5-FU (Lev+5FU) can increase patients' survival rate. And its survival probability decreases from 1 and ends above 0.5, showing that most of the patients survive after the study.

Lev (green) and obs (red) lines do not show much difference. The median survival time for obs is 2052 days and for Lev is 2257 days greater than obs. 95% confidence interval of median survival time for obs is (1550, 2718) and the lower bound for confidence interval for Lev is 1767. Since the two confidence intervals overlap, there is not statistically significant difference between the median survival time of obs and Lev.

```
# Summary of PH model
summary(cox_death)
```

```
## Call:
## coxph(formula = Surv(time, status) ~ rx, data = colon_death)
##
##   n= 888, number of events= 430
##
##              coef exp(coef) se(coef)      z Pr(>|z|)
## rxLev        -0.06269   0.93923  0.11319 -0.554  0.57967
## rxLev+5FU    -0.38280   0.68195  0.12110 -3.161  0.00157 **
```

```
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
##          exp(coef) exp(-coef) lower .95 upper .95
## rxLev      0.9392      1.065    0.7524    1.1725
## rxLev+5FU   0.6819      1.466    0.5379    0.8646
##
## Concordance= 0.537 (se = 0.013 )
## Likelihood ratio test= 11.41 on 2 df,  p=0.003
## Wald test              = 10.9 on 2 df,  p=0.004
## Score (logrank) test = 11.02 on 2 df,  p=0.004
```

Since the p -value of likelihood ratio test is 0.003 less than $\alpha = 0.05$, there is sufficient evidence to conclude that **rx** has significant impact on the survival time of patients.

Since hazard ratio for **Lev** is 0.9392, patients on **Lev** has 6.08% lower hazard rate than observation. Also the 95% confidence interval for **Lev** is (0.7524, 1.1725) including 1. Thus, Levamisole does not have significant impact on the survival probability.

The hazard ratio for **Lev+5FU** is 0.6819, **Lev+5FU** has 32% lower hazard rate than **Observation**. Also, 95 confidence interval (0.5379, 0.8646) does not include 1. Treatment **Lev+5FU** significantly lower the hazard rate and increase the survival probability of patient.

Exploring Covariates for the Marginal Death Model:

AIC

node4 had the smallest AIC; **node4** will be the first covariate.

extent had the smallest AIC and will be the second covariate.

surg had the smallest AIC, so it will be the third covariate.

differ had the smallest AIC and will be the fourth covariate.

```
# 5th Covariate
uni_vars5 <- c("obstruct", "adhere",
              "perfor", "age", "sex")

uni_models5 <- map(uni_vars5, \(v)
  coxph(as.formula(paste("Surv(time, status) ~ rx + node4 + extent + surg +
                        differ + ", v))),
  data = colon_death)
) |> set_names(uni_vars5)

aic_tbl5 <- map_dbl(uni_models5, AIC) |>
  sort() |>
  round(2)

aic_tbl5
```

```
## obstruct      age    adhere    perfor      sex
## 5425.77 5426.61 5427.38 5429.36 5429.50
```

We selected extra covariates by forward AIC while always keeping treatment (**rx**) in the model. Adding **node4**, **extent**, and **surg** each cut AIC by > 2 points, and **differ** lowered it by another 2.4; **obstruct** reduced AIC by < 2. Because 2 points is the standard threshold for a meaningful gain, we stopped at **rx + node4 + extent + surg + differ**. This captures nearly all improvement in fit without adding unnecessary parameters.

Full Coxph Model for Death

```
full_death <- coxph(Surv(time, status) ~ node4 + extent +
                    surg + differ + rx, data = colon_death)
summary(full_death)
```

```
## Call:
## coxph(formula = Surv(time, status) ~ node4 + extent + surg +
##       differ + rx, data = colon_death)
##
##      n= 888, number of events= 430
##
##              coef exp(coef) se(coef)      z Pr(>|z|)
## node41          0.90682   2.47645  0.10076  9.000 < 2e-16 ***
## extent2         0.58647   1.79764  0.60331  0.972  0.33100
## extent3         1.10438   3.01735  0.58214  1.897  0.05781 .
## extent4         1.56152   4.76605  0.61527  2.538  0.01115 *
## surglong        0.23256   1.26182  0.10625  2.189  0.02862 *
## differmoderate -0.08838   0.91541  0.16812 -0.526  0.59910
## differpoor      0.23602   1.26620  0.19489  1.211  0.22587
## rxLev          -0.04548   0.95554  0.11429 -0.398  0.69066
## rxLev+5FU      -0.37257   0.68896  0.12185 -3.058  0.00223 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
##              exp(coef) exp(-coef) lower .95 upper .95
## node41          2.4764    0.4038    2.0326    3.0172
## extent2         1.7976    0.5563    0.5510    5.8646
## extent3         3.0173    0.3314    0.9641    9.4437
## extent4         4.7661    0.2098    1.4271   15.9176
## surglong        1.2618    0.7925    1.0246    1.5540
## differmoderate  0.9154    1.0924    0.6584    1.2727
## differpoor      1.2662    0.7898    0.8642    1.8552
## rxLev           0.9555    1.0465    0.7638    1.1954
## rxLev+5FU       0.6890    1.4515    0.5426    0.8748
##
## Concordance= 0.662 (se = 0.013 )
## Likelihood ratio test= 126.3 on 9 df,  p=<2e-16
## Wald test              = 129.1 on 9 df,  p=<2e-16
## Score (logrank) test = 139.3 on 9 df,  p=<2e-16
```

After adjusting for the four strongest prognostic factors—**node4**, **extent**, **surg**, and **differ**—the overall likelihood-ratio test is highly significant ($p < 2 \times 10^{-16}$), confirming that the set of covariates is statistically significant to explain variation in survival model. From the summary of the cox proportional model, we can observe the following effect of treatment and prognostic covariate:

Treatment effect:

The combination therapy Levamisole+5-FU has statistically significant survival benefit, reducing the hazard of death by approximately 31% with (HR = 0.689, 95% CI 0.54~0.88). Levamisole alone does not show significant benefit because the 95% CI (0.7638 – 1.1954) include 1.

Prognostic covariates:

node4: having more than 4 positive lymph nodes has hazard ratio of 2.4764 and significantly increase the hazard risk by 147% compared to less than 4 lymph nodes (95% CI 2.03~3.02).

extent: Contiguous structures of local spread (**extent** = 4) raises the hazard by 377% compared to to submucosa of local spread (**extent**=1) (HR = 4.77, 95 % CI 1.96–15.9).

surg: Long time from surgery to registration (**surg** = 1) also raise the hazard rate by 26.18% compared to shorter time (**surg**= 0) (HR = 1.26, 95% CI 1.02~1.55).

Counting Process Model

The results from the marginal Cox proportional hazards models revealed that the combination treatment of Levamisole and 5-Fluorouracil (Lev+5FU) was significantly associated with reduced hazard for both recurrence and death. Additionally, time from initial surgery(**surg**), level of tumor spread(**extent**), and whether or not patients had more than four positive lymph nodes(**node4**) were important covariates associated with recurrence and death.

The data is looking at the time between recurrences of colon cancer and death. As a result, it cannot be treated as independent intervals like the gap model because there is a relationship between the recurrence of colon cancer and death from colon cancer. Thus, the counting process model is the best model to use.

```
# Preparing the split
colon$start = 0
colon$stop = colon$time

# Builds the start and stop covariates for the subset of the colon cancer data
# where only the recurrences of cancer are noted.
recurrence_data = recurrence_data %>%
  mutate(start = 0,
         stop = time)

# Established the subset of the colon cancer data where only the deaths from
# colon cancer are noted

death_data = colon[colon$etype == 2,]

# Modifying the death subset to join the recurrence subset by id.
# The mutation is so the end time for the recurrence row for a given id/subject
# is the start time for the death row for the same id/subject.

death_data = death_data %>%
  left_join(recurrence_data %>%
    select(id, recurrence_time = time), by = "id") %>%
  mutate(start = recurrence_time,
         stop = time)

# The two subsets are then combined by the id covariate from earlier.
# They are then arranged in order of id. The recurrence_time covariate,
# established in the joining process, is removed as it's just a helper variable.
# Finally, the start covariate is mutated in the case where the subject does
```

```
#not have a recurrence of colon cancer but passes away instead.
```

```
colon_counting = bind_rows(recurrence_data, death_data) %>%  
  arrange(id) %>%  
  select(-recurrence_time) %>%  
  mutate(start = if_else(start == stop, 0, start))
```

Now the colon dataset counts the time between the beginning to the first recurrence, the time between the recurrence to the next recurrence or death.

However, the goal of the counting process model is to examine the effect of the treatments before recurrence and after recurrence. To do this, a new *episode* covariate must be created. If a given subject experienced a recurrence of colon cancer and died during the study, then they would have both an episode of 0 & episode of 1 for the respective rows. The same would apply if their death was censored/outside of the study. However, if a subject did not have a recurrence of colon cancer but passed away in the study, then they would only have two episodes of 1, but the first row would be deleted. In the possibility that a patient didn't have a recurrence of colon cancer and their death was censored, there would be two episodes of 0, but the first row would remain.

```
# Recurrence + Death and Recurrence + Censored Death are achieved in the first two cases.  
# Censored Recurrence + Death is achieved in the third case  
# Censored Recurrence + Censored Death is achieved in the last case
```

```
colon_counting = colon_counting %>%  
  group_by(id) %>%  
  mutate(episode = case_when(  
    (diff(stop) != 0 & etype == 1) ~ 0,  
    (diff(stop) != 0 & etype == 2) ~ 1,  
    (diff(stop) == 0 & diff(status) != 0) ~ 1,  
    (diff(stop) == 0 & diff(status) == 0) ~ 0))
```

```
# Rows are kept if they fall in one of three conditions:  
#If the difference in the stop times is not equal to 0 ->  
# Recurrence + Death & Recurrence + Censored Death  
#If the differences for both stop times and statuses are equal to 0 ->  
# Censored Recurrence + Censored Death  
#If the difference for stop times is equal to 0 and the difference for statuses  
# is not equal to 0 -> Censored Recurrence + Death
```

```
colon_counting = colon_counting %>%  
  group_by(id) %>%  
  filter((diff(stop) != 0) | (diff(status) == 0 & diff(stop) == 0 & etype == 1) | (diff(status) != 0 & diff(stop) == 0))
```

Now that the counting process model is fully setup, it can be evaluated to see if it violates the cox proportional hazards assumption. The covariates from the marginal model will be included. Because of the episode covariate, the interaction between its levels and the treatment covariate will be tested.

```
summary(coxph(Surv(start,stop,status) ~ strata(episode)*rx + strata(node4) + extent + surg, data = colon_counting))
```

```
## Call:  
## coxph(formula = Surv(start, stop, status) ~ strata(episode) *  
##       rx + strata(node4) + extent + surg, data = colon_counting)  
##
```

```

##    n= 1328, number of events= 870
##
##
##               coef exp(coef) se(coef)      z Pr(>|z|)
## rxLev          -0.06913   0.93321  0.11011 -0.628  0.5301
## rxLev+5FU       -0.51119   0.59978  0.12151 -4.207 2.59e-05
## extent          0.32920   1.38986  0.07922  4.155 3.25e-05
## surg           0.15879   1.17209  0.07482  2.122  0.0338
## strata(episode)episode=1:rxLev  0.21254   1.23682  0.15854  1.341  0.1800
## strata(episode)episode=1:rxLev+5FU 0.73624   2.08806  0.17311  4.253 2.11e-05
##
## rxLev
## rxLev+5FU          ***
## extent             ***
## surg               *
## strata(episode)episode=1:rxLev
## strata(episode)episode=1:rxLev+5FU ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
##               exp(coef) exp(-coef) lower .95 upper .95
## rxLev          0.9332    1.0716    0.7521    1.1580
## rxLev+5FU       0.5998    1.6673    0.4727    0.7611
## extent          1.3899    0.7195    1.1900    1.6233
## surg            1.1721    0.8532    1.0122    1.3572
## strata(episode)episode=1:rxLev  1.2368    0.8085    0.9065    1.6875
## strata(episode)episode=1:rxLev+5FU 2.0881    0.4789    1.4873    2.9316
##
## Concordance= 0.59 (se = 0.012 )
## Likelihood ratio test= 45.68 on 6 df,  p=3e-08
## Wald test              = 43.22 on 6 df,  p=1e-07
## Score (logrank) test = 43.46 on 6 df,  p=9e-08

```

From the summary, the model as a whole with the additional interaction between the given episode of colon cancer (time after recurrence vs time before recurrence) and the treatment covariate is statistically significant at a p-value of $3e - 08$ vs $\alpha = 0.05$ for the critical value.

rx: The hazard rate for patients who took the treatment with just Levamisole is only 6.68% less hazardous than taking no treatment at all. Those who took Fluoracil in addition to Levamisole benefited with a hazard ratio of 0.5998, 40.02% less hazardous than no treatment at all.

extent: As the spread of the tumor developed from muscles to contiguous structures, the hazard ratio to those who only had submucosa development increased to as high as 1.3899 times as likely to suffer a recurrence of colon cancer.

surg: Patients with a long time from their initial surgery to registration in the study had a 17.21% greater hazard rate than those with a shorter time interval.

The main takeaway from the summary above is that the treatment of Levamisole + Fluoracil was 1.68 times as effective as just Levamisole in the second time interval after recurrence to death in its hazard rate. This would suggest that for subjects who've just experienced a recurrence in colon cancer, the combined treatment of Levamisole and Fluoracil poses greater hazard risk of death than taking just Levamisole.

Proportional Hazards Assumption (if you are editing this, just remove the existing assumption and put in your revised one)

To verify that the Cox models we use are valid, we assessed the proportional-hazards (PH) assumption. Our checking strategy combined (i) graphical inspection and (ii) formal statistical testing.

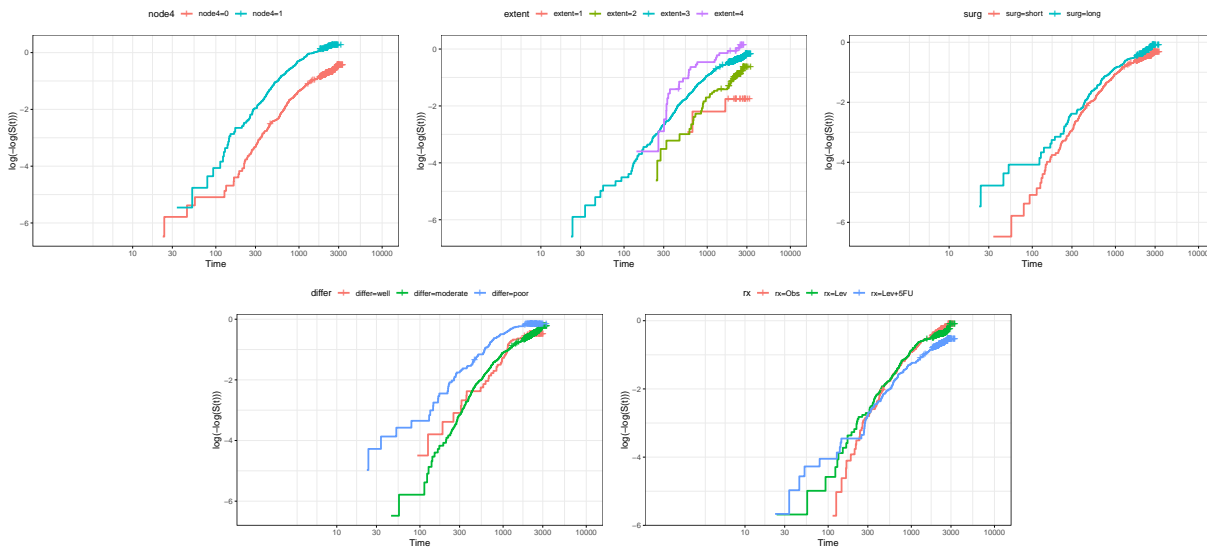
Recurrence Event

Death Event

Log-log survival plots

Under the proportional-hazard (PH) assumption the category-specific curves should be roughly straight and parallel; systematic divergence, convergence, or crossing suggests non-PH behavior.

```
for (v in c("node4", "extent", "surg", "differ", "rx")) {  
  fit <- survfit(as.formula(paste("Surv(time,status) ~", v)), data = colon_death)  
  print(ggsurvplot(fit, data = colon_death, fun = "cloglog",  
    legend.title = v, ggtheme = theme_bw()))  
}
```



- *node4*: Two curves seem parallel so *node4* satisfy the assumption.
- *extent*: At the early stage the blue curve crosses the purple curve and the green curve crosses the red curve. After that, the purple, blue, and green curves are relatively parallel but the red curve seems to deviate from them. Thus, the assumption is concerning.
- *surg*: Two curves never cross and remain parallel over time. So, *surg* satisfy the assumption.
- *differ*: The green curve crosses the other two curves four times over time. And the gap between each curve becomes smaller over time. Thus, *differ* violates the assumption.
- *rx*: Three curves cross each other at the early stage but with few observations. And the green and red curves are close to each other but the gap between each curve is relatively parallel. It is hard to make a conclusion by looking at the log-log plot so we will use statistical testing to test the assumption.

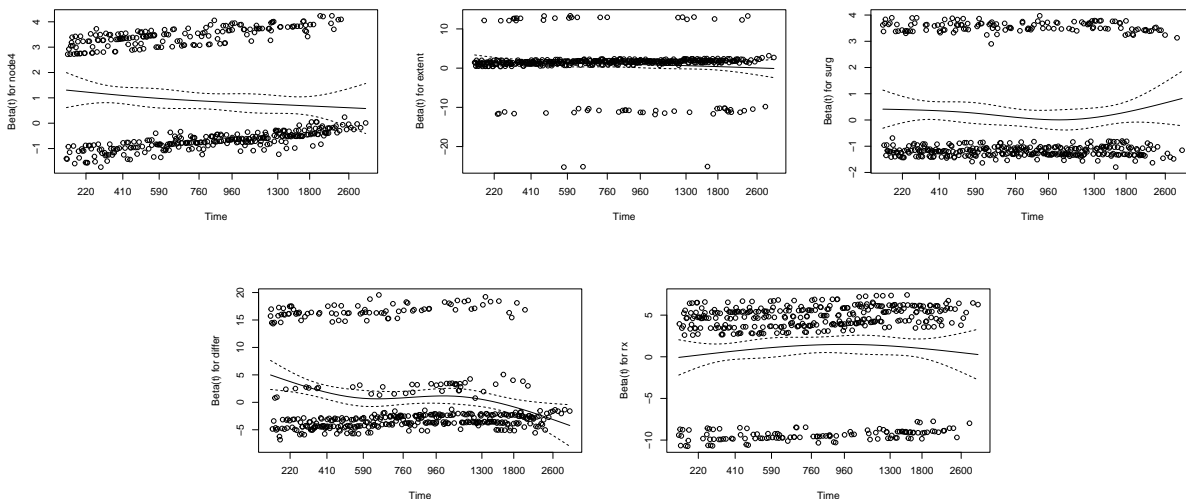
Cox ZPH

```
ph <- cox.zph(full_death)
print(ph)
```

```
##          chisq df      p
## node4    5.4346  1 0.01974
## extent   6.5521  3 0.08763
## surg      0.0205  1 0.88616
## differ   15.8513  2 0.00036
## rx        2.5680  2 0.27693
## GLOBAL   27.5832  9 0.00112
```

- *node4*: Since p -value is 0.01974 less than $\alpha = 0.05$, *node4* violate the assumption. This conclusion is different from the conclusion made from log-log plot.
- *extent*: Since p -value is 0.08763 greater than $\alpha = 0.05$, it does not violate the assumption. Although there is no sufficient evidence of non-PH, it is concerning about the assumption.
- *surg*: p -value is large so it does not violate the assumption.
- *differ*: Since p -value is 0.00036 smaller than $\alpha = 0.05$, it violate the assumption.
- *rx*: p -value is large so it does not violate the assumption.

```
plot(ph)
```



- *node4*
 - Smoothed $\beta(t)$ line is almost flat except for a slight downward bend late in follow-up.
 - **Decision**: keep *node4* in the model; add a mild log-time interaction ($\text{node4} \times \log t$) so its tiny time-trend is absorbed while preserving a clinically useful hazard-ratio.
- *differ*

- Smoothed $\beta(t)$ line declines steadily and have a clear monotone drift.
- **Decision:** stratify on differ ('strata(differ)'). This gives each grade its own baseline hazard, removes the PH requirement, and avoids reporting a misleading, time-dependent HR for grade—acceptable because grade is a control variable, not a primary target of inference.

```
full_death2 <- coxph(Surv(time, status) ~ node4 + extent +
                     surg + strata(differ) + rx, data = colon_death)
# summary(full_death2)
cox.zph(full_death2)
```

```
##           chisq df      p
## node4    3.77162  1 0.052
## extent   5.89794  3 0.117
## surg     0.00474  1 0.945
## rx       2.46820  2 0.291
## GLOBAL  12.18168  7 0.095
```

```
fit_final <- coxph(Surv(time, status) ~ extent + surg + node4 + (node4) +
                  strata(differ) + rx, data = colon_death,
                  tt = function(x,t,...) {x*log(t)})
cox.zph(fit_final)           # should now give non-significant p-values
```

```
##           chisq df      p
## extent   5.89794  3 0.117
## surg     0.00474  1 0.945
## node4    3.77162  1 0.052
## rx       2.46820  2 0.291
## GLOBAL  12.18168  7 0.095
```

References

1. colon: Chemotherapy for Stage B/C colon cancer. Colon Cancer. Retrieved May 1st, 2025, from <https://rdrr.io/cran/survival/man/colon.html>
2. CG Moertel, TR Fleming, JS MacDonald, DG Haller, JA Laurie, CM Tangen, JS Ungerleider, WA Emerson, DC Tormey, JH Glick, MH Veeder and JA Maillard, Fluorouracil plus Levamisole as an effective adjuvant therapy after resection of stage II colon carcinoma: a final report. Annals of Internal Med, 122:321-326, 1991.