

# P8108 Group 2 Survival Analysis Project

Yiming Zhao (yz3955)      Wenshan Qu (wq2160)      Tucker Morgan (t1m2152)  
Junzhe Shao (js5959)      Benjamin Goebel (bpg2118)

2022-12-1

```
library(survival)
library(tidyverse)
library(tidymodels)
library(glmnet)
library(ranger)
library(survminer)
library(arsenal)
knitr::opts_chunk$set(message = FALSE, warning = FALSE)
```

## Train Validation Test Split

```
set.seed(2022)

rotterdam_split <- initial_split(select(rotterdam, -rtime, -recur, -pid),
                                prop = 0.8, strata = death)
rotterdam_training <- training(rotterdam_split)
rotterdam_test <- testing(rotterdam_split)

rotterdam_train_val_split <- initial_split(rotterdam_training,
                                           prop = 0.8, strata = death)
rotterdam_training <- training(rotterdam_train_val_split)
rotterdam_validation <- testing(rotterdam_train_val_split)
```

## Perform 10-fold Cross-Validation

The output contains 1 row for each fold/repeat. So, 10 folds \* 5 repeats = 50 rows. The `split_analysis` column is a list column containing a data frame for each row with 9 folds combined, and the `split_assessment` column is a list column containing a data frame for each row with 1 fold.

```
set.seed(2022)

rotterdam_folds <- vfold_cv(rotterdam_training, v = 10, repeats = 5,
                           strata = death)

rotterdam_folds <- rotterdam_folds %>%
  mutate(split_analysis = map(splits, analysis),
         split_assessment = map(splits, assessment))
```

## Introduction

## Methods

The dataset of interest for this analysis comes from the Rotterdam tumor bank, including data from 2982 breast cancer patients. Follow up time for patients varied from just 1 month to as long as 231 months. Several prognostic variables are recorded including year of surgery, age at surgery, menopausal status (pre- or post-), tumor size (mm), differentiation grade, number of positive lymph nodes, progesterone receptors (fmol/l), estrogen receptors (fmol/l), and indicators for hormonal treatment and chemotherapy treatment. The outcome considered in this analysis was patient death.

(Placeholder for Cross-validation)

As part of this analysis, we consider the Cox Proportional Hazard (Cox PH) model, which allows us to model the hazard ratio based on covariates to understand their impact on the survival function. The Cox PH typically takes the form:

$$h(t|Z = z) = h_0(t)e^{\beta'z}.$$

In this application, we use the elastic net penalty, a mixture of the  $\ell_1$  and  $\ell_2$  norm regularization penalties. In the Cox PH framework, this penalty term takes the form of:

$$\lambda\left(\alpha \sum |\beta_i| + \frac{1}{2}(1 - \alpha) \sum \beta_i^2\right)$$

where  $\lambda$  represents our penalty coefficient and  $\alpha$  is the mixing parameter for the two regularization methods. This penalty helps to avoid over-fitting of our data. The algorithm used here in **glmnet** uses the Breslow approximation to handle ties. For more details on the derivation of this term and the algorithm used to fit the penalized Cox PH model, see Simon et al. (2011).

## Exploratory Data Analysis

```
print(summary(tableby(hormon~age+meno+size+grade+nodes+pgr+er+chemo+dtime+death,
                      rotterdam,numeric.simplify = TRUE, numeric.test = "kwt")))
```

	0 (N=2643)	1 (N=339)	Total (N=2982)	p value
<b>age</b>				< 0.001
Mean (SD)	54.098 (12.984)	62.549 (9.921)	55.058 (12.953)	
Range	24.000 - 90.000	28.000 - 88.000	24.000 - 90.000	
<b>meno</b>				< 0.001
Mean (SD)	0.519 (0.500)	0.879 (0.327)	0.560 (0.496)	
Range	0.000 - 1.000	0.000 - 1.000	0.000 - 1.000	
<b>size</b>				< 0.001
<=20	1283 (48.5%)	104 (30.7%)	1387 (46.5%)	
20-50	1119 (42.3%)	172 (50.7%)	1291 (43.3%)	
>50	241 (9.1%)	63 (18.6%)	304 (10.2%)	
<b>grade</b>				< 0.001
Mean (SD)	2.722 (0.448)	2.826 (0.380)	2.734 (0.442)	
Range	2.000 - 3.000	2.000 - 3.000	2.000 - 3.000	
<b>nodes</b>				< 0.001
Mean (SD)	2.327 (4.207)	5.720 (4.576)	2.712 (4.384)	
Range	0.000 - 34.000	1.000 - 24.000	0.000 - 34.000	
<b>pgr</b>				< 0.001

	0 (N=2643)	1 (N=339)	Total (N=2982)	p value
<b>er</b>				
Mean (SD)	168.706 (300.337)	108.233 (200.302)	161.831 (291.311)	
Range	0.000 - 5004.000	0.000 - 1497.000	0.000 - 5004.000	
<b>chemo</b>				0.069
Mean (SD)	164.792 (272.563)	180.608 (271.693)	166.590 (272.465)	
Range	0.000 - 3275.000	0.000 - 2444.000	0.000 - 3275.000	
<b>death</b>				< 0.001
Mean (SD)	0.209 (0.407)	0.083 (0.276)	0.195 (0.396)	
Range	0.000 - 1.000	0.000 - 1.000	0.000 - 1.000	
<b>dttime</b>				< 0.001
Mean (SD)	2679.067 (1309.178)	2030.534 (1043.971)	2605.340 (1298.078)	
Range	36.000 - 7043.000	45.000 - 6270.000	36.000 - 7043.000	
<b>death</b>				0.093
Mean (SD)	0.421 (0.494)	0.469 (0.500)	0.427 (0.495)	
Range	0.000 - 1.000	0.000 - 1.000	0.000 - 1.000	

```
print(summary(tableby(~age+meno+size+grade+nodes+pgr+er+chemo+hormon+dttime+death,
                      rotterdam,numeric.simplify = TRUE, numeric.test = "kwt")))
```

	Overall (N=2982)
<b>age</b>	
Mean (SD)	55.058 (12.953)
Range	24.000 - 90.000
<b>meno</b>	
Mean (SD)	0.560 (0.496)
Range	0.000 - 1.000
<b>size</b>	
<=20	1387 (46.5%)
20-50	1291 (43.3%)
>50	304 (10.2%)
<b>grade</b>	
Mean (SD)	2.734 (0.442)
Range	2.000 - 3.000
<b>nodes</b>	
Mean (SD)	2.712 (4.384)
Range	0.000 - 34.000
<b>pgr</b>	
Mean (SD)	161.831 (291.311)
Range	0.000 - 5004.000
<b>er</b>	
Mean (SD)	166.590 (272.465)
Range	0.000 - 3275.000
<b>chemo</b>	
Mean (SD)	0.195 (0.396)
Range	0.000 - 1.000
<b>hormon</b>	
Mean (SD)	0.114 (0.317)
Range	0.000 - 1.000
<b>dttime</b>	
Mean (SD)	2605.340 (1298.078)
Range	36.000 - 7043.000

	Overall (N=2982)
<b>death</b>	
Mean (SD)	0.427 (0.495)
Range	0.000 - 1.000

## Cross-Validation

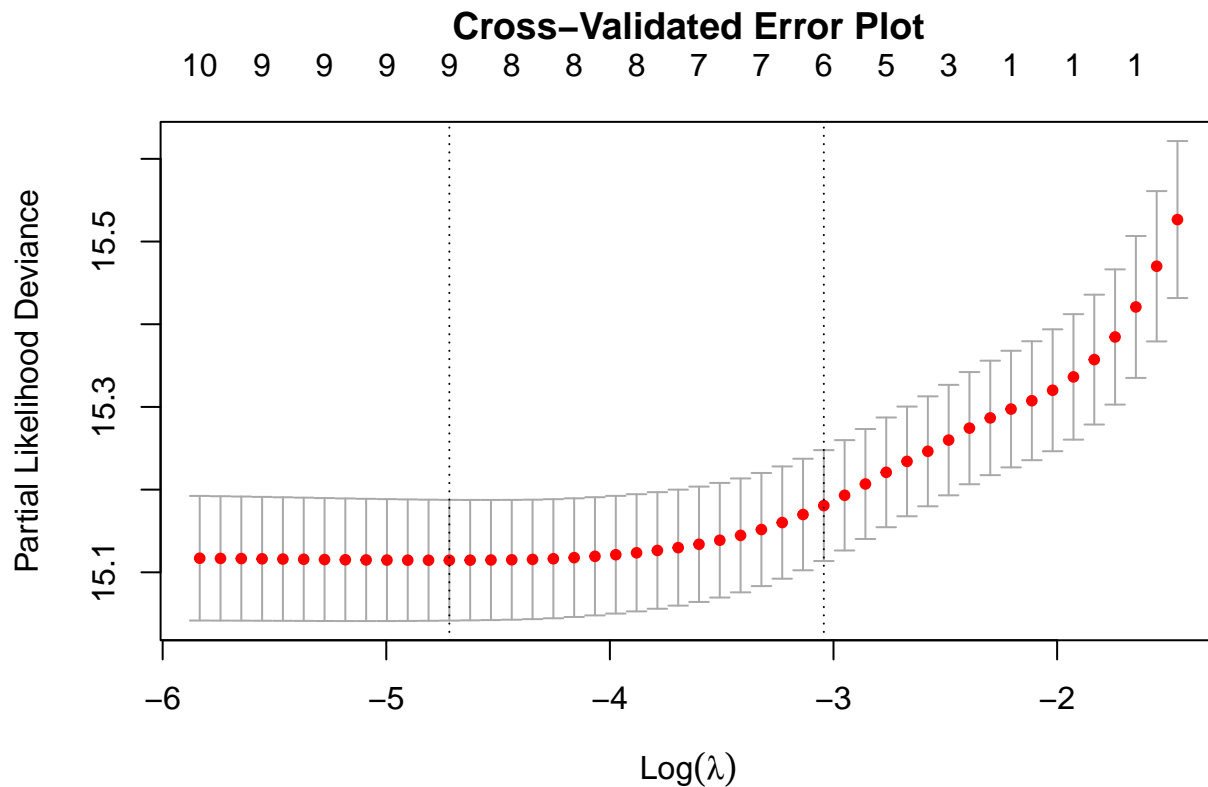
### Cox/Cox with elastic net

```
set.seed(2022)

cox_trn_x <- model.matrix(Surv(dtime, death) ~ ., rotterdam_training)[,-1]
cox_trn_y <- Surv(rotterdam_training$dtime, rotterdam_training$death)

cv_coxfit <- cv.glmnet(cox_trn_x, cox_trn_y, family = "cox", type.measure = "deviance")

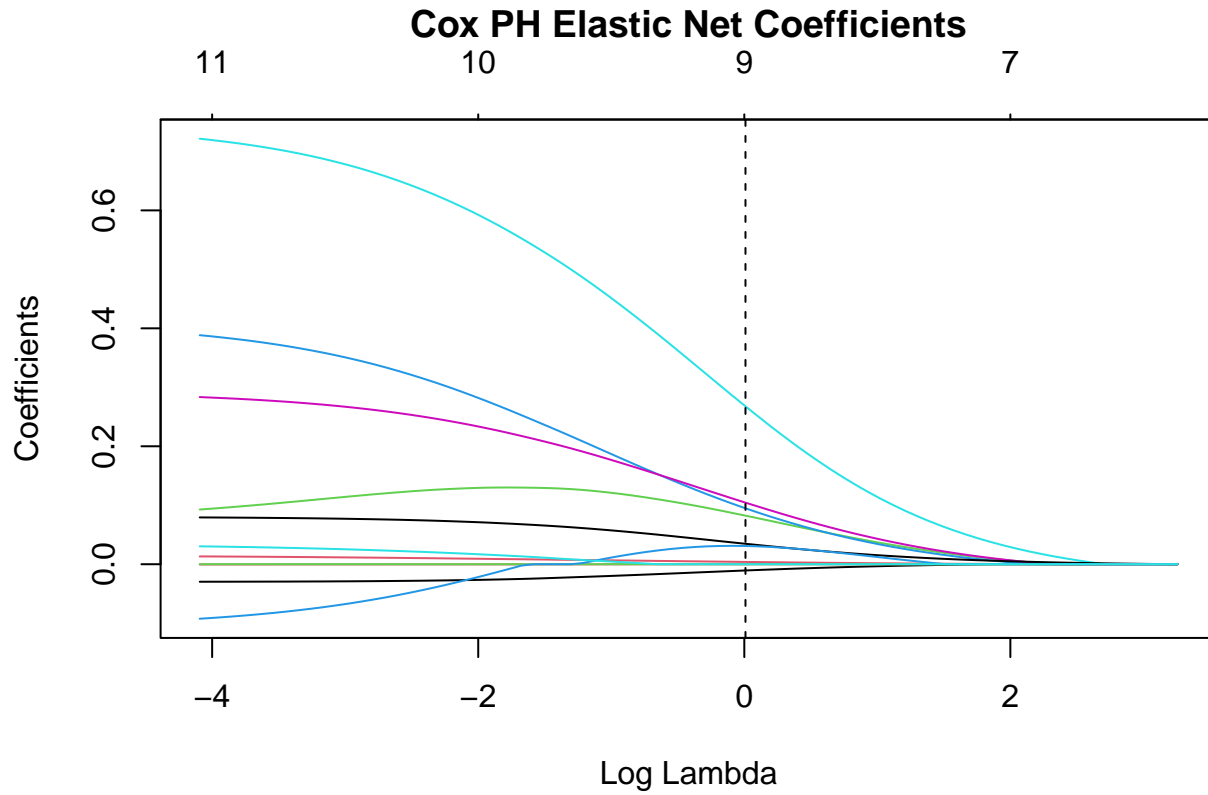
par(mar = c(4,4,5,1))
plot(cv_coxfit, main = "Cross-Validated Error Plot")
```



```
coxnetfit <- glmnet(cox_trn_x, cox_trn_y, family = "cox", alpha = cv_coxfit$lambda.min)

par(mar = c(4,4,5,1))
```

```
plot(coxnetfit, xvar = "lambda",
     main = "Cox PH Elastic Net Coefficients")
abline(v = cv_coxfit$lambda.min, lty = 2)
```



```
coxnetfit_df <-
  data.frame(
    "coef" = as.vector(coef(coxnetfit, s = cv_coxfit$lambda.min)),
    "exp_coef" = as.vector(coef(coxnetfit, s = cv_coxfit$lambda.min)) %>% exp()
  )

rownames(coxnetfit_df) <- labels(coef(coxnetfit, s = cv_coxfit$lambda.min))[[1]]

coxnetfit_df %>% round(digits = 4) %>%
  knitr::kable(caption = "Cox Proportion Hazard Elastic Net Coefficients")
```

Table 3: Cox Proportion Hazard Elastic Net Coefficients

	coef	exp_coef
year	-0.0297	0.9707
age	0.0132	1.0133
meno	0.0927	1.0972
size20-50	0.3884	1.4746
size>50	0.7215	2.0575

	coef	exp_coef
grade	0.2834	1.3276
nodes	0.0794	1.0826
pgr	-0.0005	0.9995
er	0.0000	1.0000
hormon	-0.0924	0.9117
chemo	0.0303	1.0308

In the table above, we can see that the estrogen receptors covariate is selected out with a null value of 0 or  $\exp(coef) = 1$ . We can fit a cox proportional hazard model using only the selected covariates in the `coxph` function to find unbiased estimates of the coefficients along with standard errors and confidence intervals.

```
coxfit <- coxph(Surv(dtime, death) ~ year + age + meno + size + grade +
               nodes + pgr + hormon + chemo,
               data = rotterdam_training, ties = "breslow")
coxfit %>%
  broom::tidy() %>%
  mutate(estimate = exp(estimate))
```

```
## # A tibble: 10 x 5
##   term      estimate std.error statistic  p.value
##   <chr>      <dbl>    <dbl>    <dbl>    <dbl>
## 1 year        0.970  0.0129     -2.33  1.96e- 2
## 2 age         1.01  0.00471     3.06  2.21e- 3
## 3 meno        1.07  0.124      0.568  5.70e- 1
## 4 size20-50   1.51  0.0826     4.97  6.76e- 7
## 5 size>50     2.11  0.115      6.49  8.53e-11
## 6 grade       1.34  0.0903     3.24  1.19e- 3
## 7 nodes       1.08  0.00646    12.5  1.26e-35
## 8 pgr         1.00  0.000147   -3.36  7.78e- 4
## 9 hormon      0.899  0.117     -0.910 3.63e- 1
## 10 chemo      1.03  0.104      0.327 7.44e- 1
```

```
confint(coxfit) %>% exp() %>% knitr::kable()
```

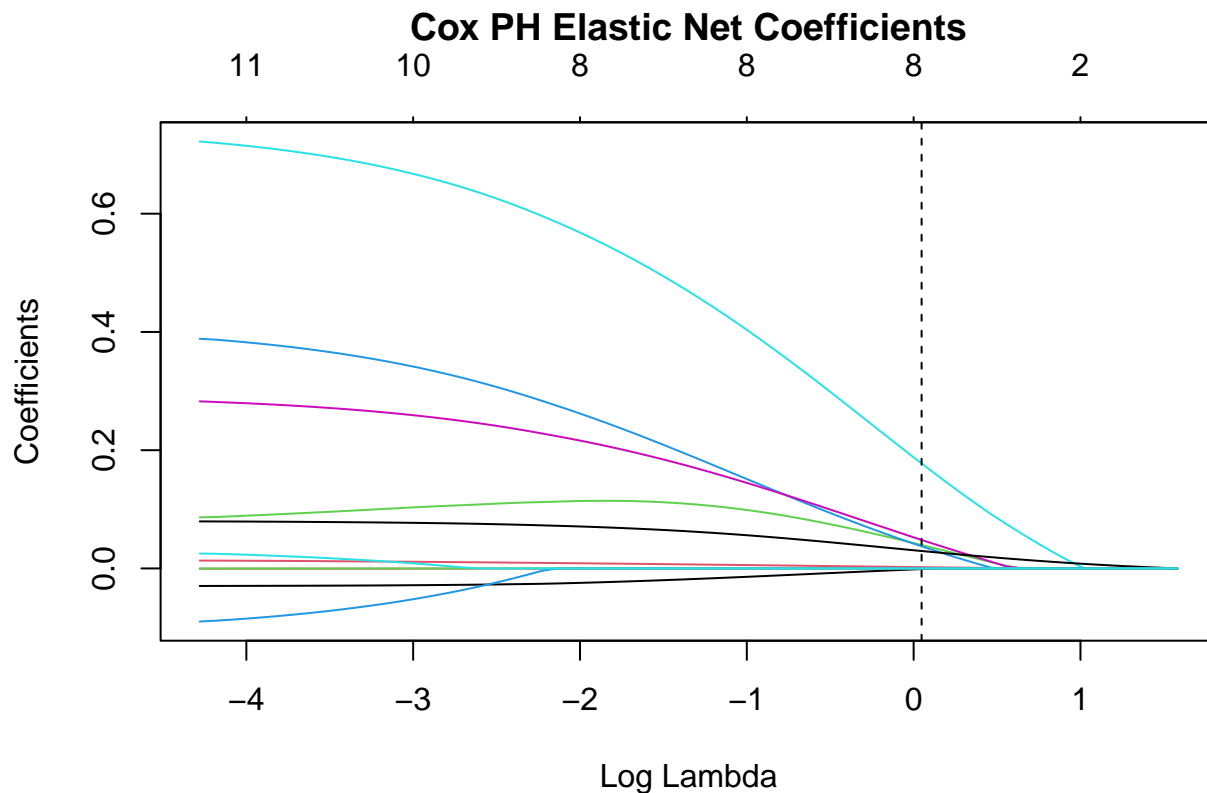
	2.5 %	97.5 %
year	0.9463073	0.9951950
age	1.0052046	1.0239543
meno	0.8412832	1.3689527
size20-50	1.2820640	1.7722756
size>50	1.6831143	2.6407952
grade	1.1225407	1.5991351
nodes	1.0702019	1.0976591
pgr	0.9992188	0.9997943
hormon	0.7145907	1.1307566
chemo	0.8437328	1.2686315

In our (minimum error) model, we find significant effects for year of surgery, age at surgery, size of tumor, differentiation grade, number of positive lymph nodes, and progesterone receptors. The largest magnitude effects come from increasing size of tumor.

Below, we see the analogous results for a “1se” rule model.

```
coxnetfit_1se <- glmnet(cox_trn_x, cox_trn_y, family = "cox", alpha = cv_coxfit$lambda.1se)

par(mar = c(4,4,5,1))
plot(coxnetfit_1se, xvar = "lambda",
     main = "Cox PH Elastic Net Coefficients")
abline(v = cv_coxfit$lambda.1se, lty = 2)
```



```
coxnetfit_1se_df <-
  data.frame(
    "coef" = as.vector(coef(coxnetfit_1se, s = cv_coxfit$lambda.1se)),
    "exp_coef" = as.vector(coef(coxnetfit_1se, s = cv_coxfit$lambda.1se)) %>% exp()
  )

rownames(coxnetfit_1se_df) <- labels(coef(coxnetfit_1se, s = cv_coxfit$lambda.1se))[[1]]

coxnetfit_1se_df %>% round(digits = 4) %>%
  knitr::kable(caption = "Cox Proportion Hazard Elastic Net Coefficients (1se)")
```

Table 5: Cox Proportion Hazard Elastic Net Coefficients (1se)

	coef	exp_coef
year	-0.0285	0.9720
age	0.0114	1.0114
meno	0.1026	1.1080
size20-50	0.3440	1.4105
size>50	0.6704	1.9551
grade	0.2605	1.2975
nodes	0.0773	1.0803
pgr	-0.0004	0.9996
er	0.0000	1.0000
hormon	-0.0541	0.9474
chemo	0.0095	1.0095

Here, we remove `er` and `chemo` and find the following results.

```
coxfit_1se <- coxph(Surv(dtime, death) ~ year + age + meno + size + grade + nodes + pgr + hormon,
  data = rotterdam_training, ties = "breslow")
coxfit_1se %>%
  broom::tidy() %>%
  mutate(estimate = exp(estimate))
```

```
## # A tibble: 9 x 5
##   term      estimate std.error statistic  p.value
##   <chr>      <dbl>     <dbl>     <dbl>   <dbl>
## 1 year      0.970    0.0129     -2.34  1.92e- 2
## 2 age      1.01    0.00464     3.05  2.27e- 3
## 3 meno     1.07    0.122      0.517 6.05e- 1
## 4 size20-50 1.51    0.0825     4.99 6.16e- 7
## 5 size>50   2.11    0.115     6.51 7.56e-11
## 6 grade     1.34    0.0903     3.24 1.19e- 3
## 7 nodes     1.08    0.00633    12.8 2.11e-37
## 8 pgr       1.00    0.000147   -3.35 8.01e- 4
## 9 hormon    0.897    0.117     -0.928 3.53e- 1
```

```
confint(coxfit_1se) %>% exp() %>% knitr::kable()
```

	2.5 %	97.5 %
year	0.9462185	0.9951111
age	1.0050778	1.0235067
meno	0.8384713	1.3529343
size20-50	1.2837050	1.7740258
size>50	1.6860598	2.6446630
grade	1.1227352	1.5993840
nodes	1.0709185	1.0978328
pgr	0.9992211	0.9997958
hormon	0.7133152	1.1282404

Here we again find significant effects for year of surgery, age at surgery, size of tumor, differentiation grade, number of positive lymph nodes, and pgr.



## Random survival forest

The survival tree and the corresponding random survival forest (RSF) are highly favorable non-parametric methods when studying survival data. Generally, for a single survival tree, it will assign subjects to groups based on certain splitting rules regarding their covariates, and the subjects in each group will share a similar survival behavior.

```
set.seed(2023)
## Random Survival Forest
rsf <- ranger(Surv(time = dtime, event = death) ~ .,
              data = rotterdam_training,
              num.trees = 300,
              min.node.size = 15)

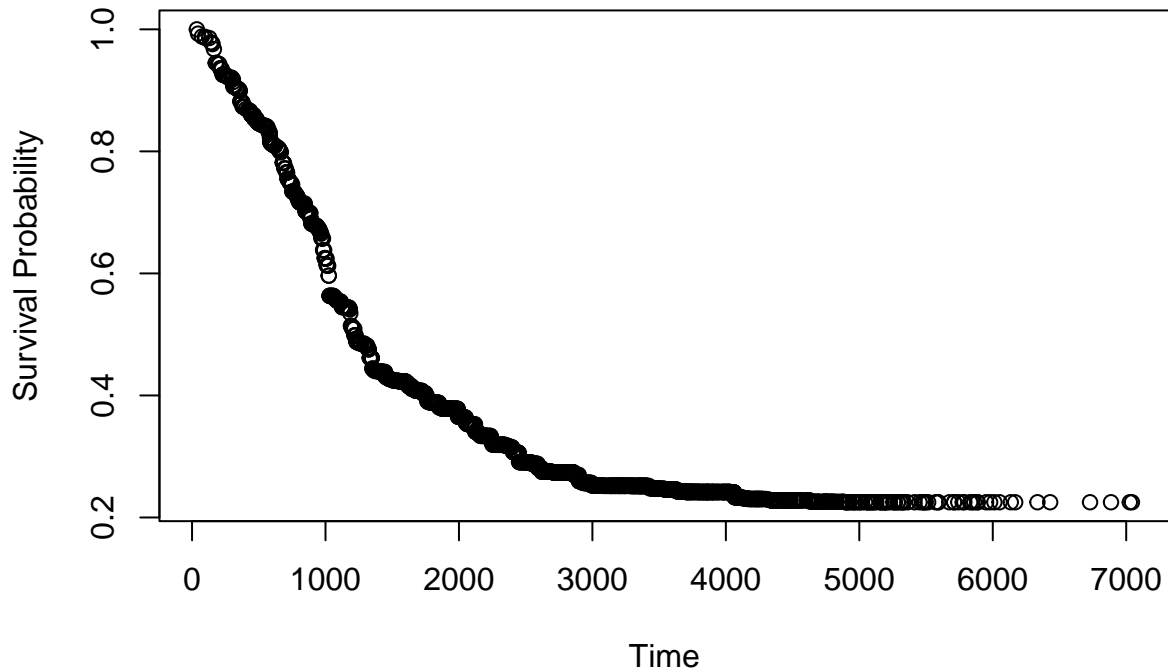
## Remove variables not for prediction, and the outcome
rotterdam_test_d <-
  rotterdam_test %>%
  select(-death)

## Make prediction on all the test data points
pred_rsf <- predict(rsf, rotterdam_test_d, type = "response")
# Look at individual 7
pred_ref_7 <- data.frame(
  time = pred_rsf$unique.death.times,
  survival = pred_rsf$survival[7,])
head(pred_ref_7) %>% knitr::kable(align = "c")
```

time	survival
36	1.0000000
45	0.9929877
74	0.9881553
97	0.9877437
101	0.9857489
129	0.9857489

```
plot(pred_ref_7$time, pred_ref_7$survival,
     xlab = "Time", ylab = "Survival Probability",
     main = "Survival Prediction for Patient 7")
```

## Survival Prediction for Patient 7



```
# Find estimated median survival time for individual 7
head(pred_ref_7[pred_ref_7$survival <= 0.5,]) %>% knitr::kable(align = "c") #1163
```

	time	survival
306	1217	0.4984613
307	1218	0.4984613
308	1222	0.4984613
309	1226	0.4984613
310	1229	0.4927672
311	1231	0.4927672

```
# See the truth of individual 7
rotterdam_test[7,] %>% knitr::kable(align = "c")
```

	year	age	meno	size	grade	nodes	pgr	er	hormon	chemo	dtime	death
2463	1992	69	1	20-50	2	8	5	6	1	0	1869	0

With **ranger** package, we trained the random survival forest with training dataset used for survival prediction. As a non-parametric method, there is no parameters in RSF that could be interpreted. The ultimate goal of RSF is to predict the survival probability function of a given data point based on its covariate vector. Compared to semi-parametric Cox-PH model which forces the outcome and the covariates to have a special

connection, the RSF makes prediction based on the survival time of training data points that shares similar propensity with the given input data point.

Since the “truth” of test data point (a single survival time) and the prediction we made here (a survival probability function) are not comparable, here we show the prediction result of the 7th test data point (pid = 58). The survival curve has been shown above, and the median survival time is 1163 days.

## Comparison of Cox Proportional-Hazards Elastic Net with Random Survival Forest

We compared the Cox proportional-hazards elastic net model with the random survival forest by calculating the Brier score for each model on the validation set. The formula for the Brier score is as follows.

$$BS = \frac{1}{n} \sum_{i=1}^n (p_i - o_i)^2$$

The Brier score is used to evaluate the accuracy of probabilistic predictions from a model; its value ranges from 0 to 1 with 0 being perfect and 1 being the opposite. We calculated the Brier score using the validation set. For each observation in the validation set, predictions were made at the observed time of censoring or event. Our analysis proceeds as follows.

First, we calculated the Brier score for the Cox proportional-hazards elastic net model.

```
# Purpose: Calculates the Brier score for the Cox proportional-hazards elastic
#           net model.
# Arguments: fit: The Cox proportional-hazards elastic net model.
#            train: A dataframe, the training data used to fit the model.
#            test: A dataframe, the data to use to calculate the Brier score.
# Returns: A double, the Brier score.
brier_coxnet <- function(fit, train, test) {
  train_x <- model.matrix(Surv(dtime, death) ~ ., train)[,-1]
  train_y <- Surv(pull(train, dtime), pull(train, death))
  test_x <- model.matrix(Surv(dtime, death) ~ ., test)[,-1]
  test_y <- pull(test, death)
  num_obs <- nrow(test_x)
  p <- vector(mode = "double", length = num_obs)
  for(i in 1:num_obs) {
    surv_fit <- survival::survfit(fit, s = cv_coxfit$lambda.min,
                                  x = train_x,
                                  y = train_y,
                                  newx = test_x[i, ])
    time_index <- tail(which(surv_fit$time <= test[i, "dtime"]), n = 1)
    p[i] <- 1 - surv_fit$surv[time_index]
  }
  return(DescTools::BrierScore(resp = test_y, pred = p))
}
(brier_coxnet <- round(brier_coxnet(coxnetfit, rotterdam_training, rotterdam_validation), 3))

## [1] 0.329
```

The Brier score for the Cox elastic net model is 0.329.

Second, let's calculate the Brier score for the random survival forest model.

```

# Purpose: Calculates the Brier score for the random survival forest model.
# Arguments: fit: The random survival forest model.
#            df: A dataframe, the data to use to calculate the Brier score.
# Returns: A double, the Brier score.
brier_ranger <- function(fit, df) {
  x <- df
  pred <- predict(fit, data = x)
  num_obs <- nrow(df)
  p <- vector(mode = "double", length = num_obs)
  for(i in 1:num_obs) {
    time_index <- tail(which(pred$unique.death.times <= x[i, "dtime"]), n = 1)
    p[i] <- 1 - pred$survival[i, time_index]
  }
  return(DescTools::BrierScore(resp = df$death, pred = p))
}
(brier_ranger <- round(brier_ranger(rsf, rotterdam_validation), 3))

```

```
## [1] 0.322
```

The Brier score for the random survival forest model is 0.322. The two models have very similar Brier scores.

## Conformalized survival analysis

### Supplemental analyses

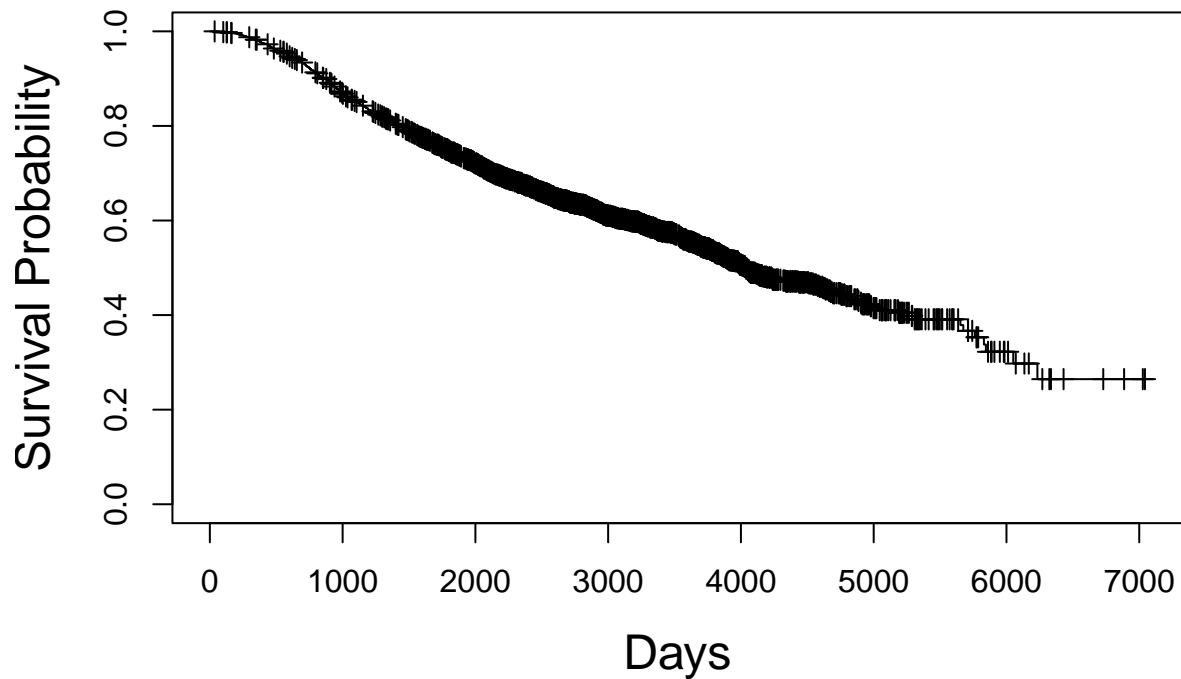
#### Kaplan-Meier Survival Estimate

```

KM = survfit(Surv(dtime, death) ~ 1, data = rotterdam)
plot(KM, conf.int = FALSE, mark.time = TRUE,
     xlab = "Days", ylab = "Survival Probability",
     main = "Kaplan-Meier Survival Estimate", cex.lab = 1.5, cex.main = 1.5)

```

## Kaplan–Meier Survival Estimate

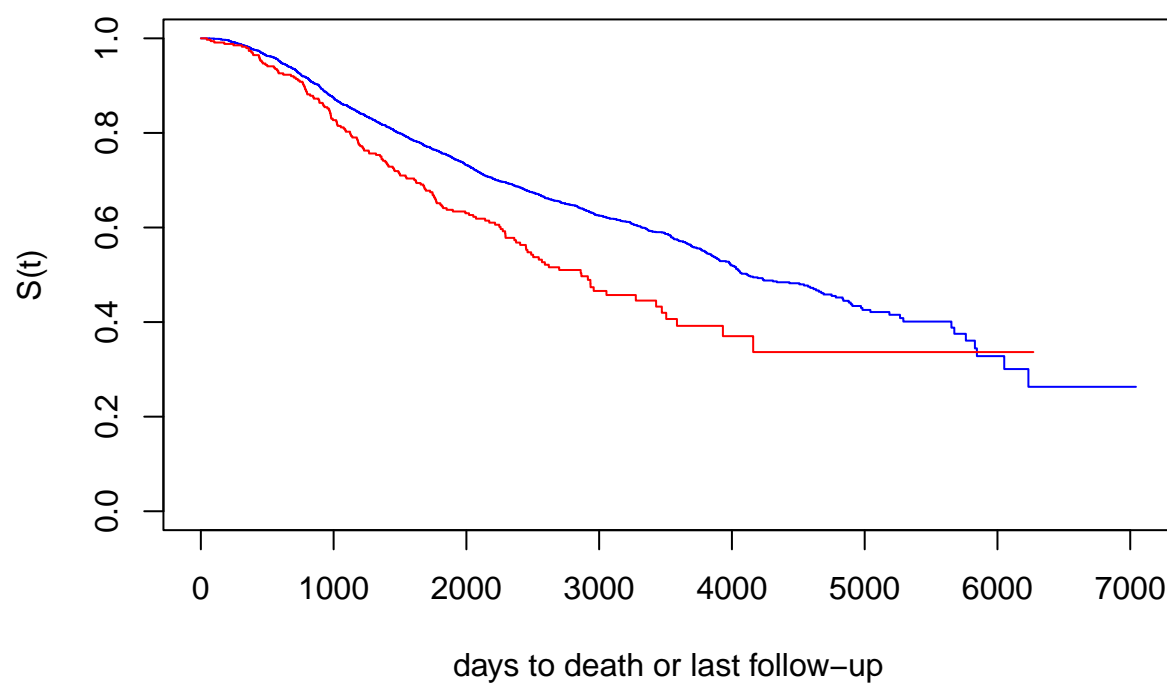


```
# make Kaplan-Meier estimates
kmfit <- survfit(Surv(dtime, death) ~ hormon, data = rotterdam, type=c("kaplan-meier"))
# print Kaplan-Meier table
#summary(kmfit)
```

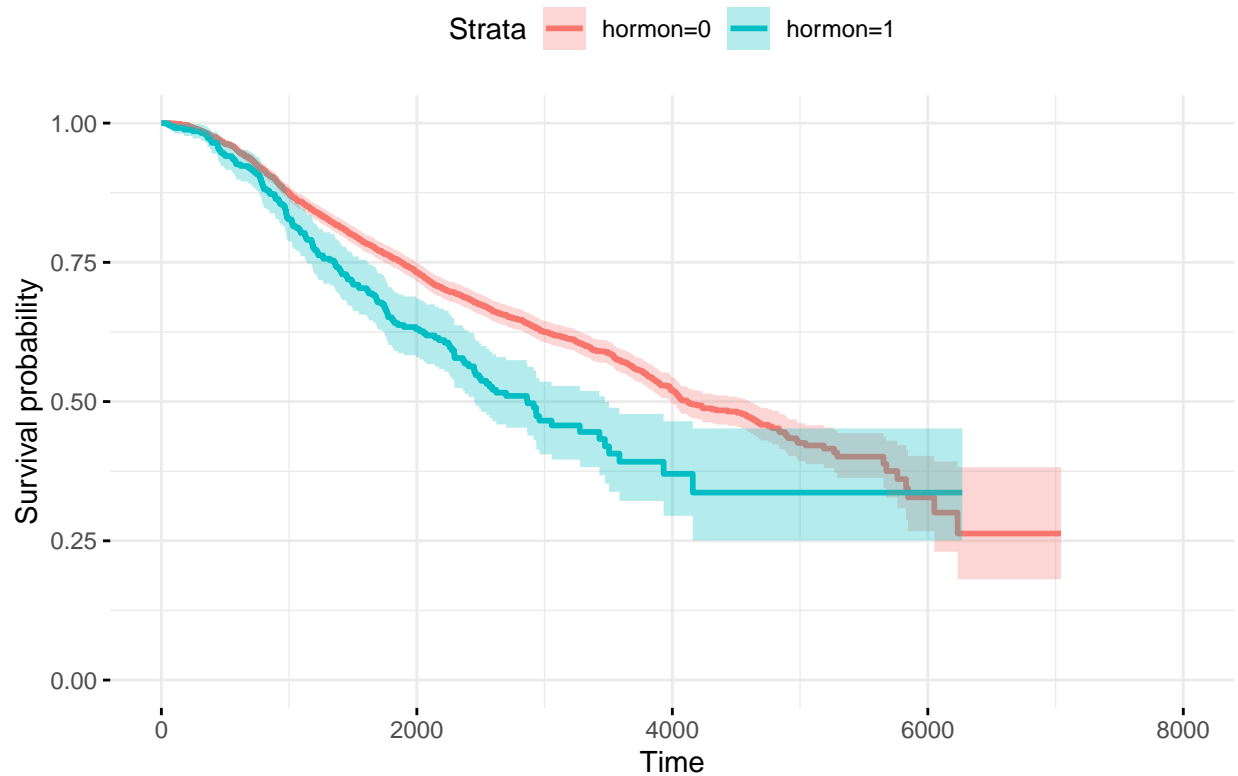
```
plot(kmfit,
     ylab="S(t)",
     xlab="days to death or last follow-up",
     main = "Kaplan Meier estimates of Breast cancer survival by hormonal treatment assignments for rotterdam",
     col = c("blue", "red"))
```

```
ggsurvplot(kmfit, conf.int = 0.95, censor=F, title = " KM survival by hormonal treatment assignments",
           ggtheme = theme_minimal())
```

## estimates of Breast cancer survival by hormonal treatment assignments



## KM survival by hormonal treatment assignments



### Log-rank Test

The null hypothesis of our log-rank test is:  $H_0 : S_1(t) = S_0(t)$ , where  $S_1(t)$  is the survival function of hormon treatment group,  $S_0(t)$  is the survival function of control group.

```
logrank <- survdiff(Surv(dtime, death) ~ hormon, data = rotterdam)
logrank
```

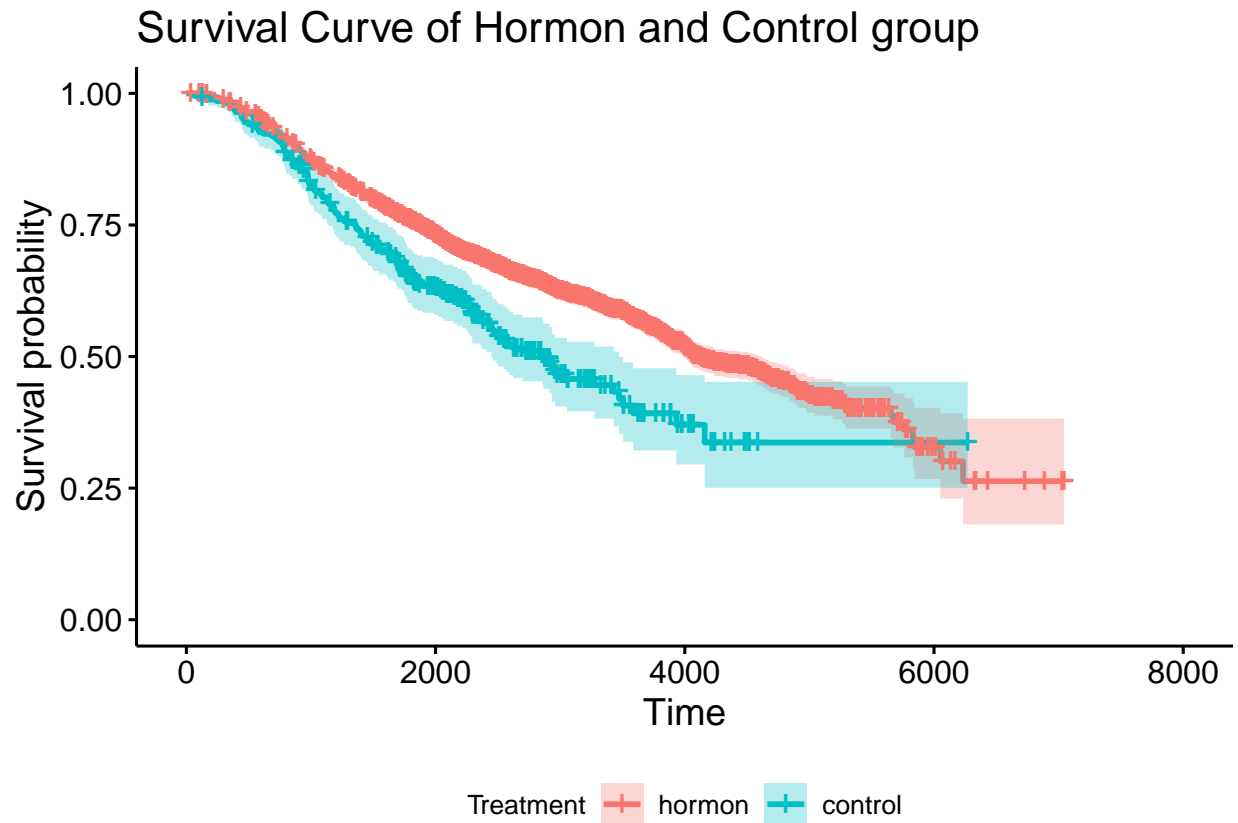
```
## Call:
## survdiff(formula = Surv(dtime, death) ~ hormon, data = rotterdam)
##
##           N Observed Expected (O-E)^2/E (O-E)^2/V
## hormon=0 2643     1113     1162      2.04      23.7
## hormon=1  339      159      110     21.43      23.7
##
##  Chisq= 23.7  on 1 degrees of freedom, p= 1e-06
```

```
logrank$pvalue
```

```
## [1] 1.133649e-06
```

The test statistic is 23.7, and the corresponding p-value is  $1.133^{-6} \ll 0.05$ , thus we reject the null and conclude that we are 95% confident that  $S_1(t) \neq S_0(t)$ . And since the test statistic is positive, we can conclude that the hormon treatment is significantly effective to breast cancer.

```
ggsurvplot(survfit(Surv(dtime,death) ~ hormon, data = rotterdam),
  conf.int = TRUE,
  legend = c("bottom"),
  legend.title = c("Treatment"),
  legend.labs = c("hormon", "control")) +
  ggtitle("Survival Curve of Hormon and Control group")
```



Results

Discussion

How our results compare with past research

Conclusion



## References

—Note this reference is in MLA format—

Simon, Noah et al. “Regularization Paths for Cox’s Proportional Hazards Model via Coordinate Descent.”  
Journal of statistical software vol. 39,5 (2011): 1-13. doi:10.18637/jss.v039.i05