

Logistic Regression

2023-12-06

Clean the data

```
1 heart_df =
2   read_csv("./Project_2_data.csv") |>
3   janitor::clean_names() |>
4   relocate(survival_months, status) |>
5   mutate(
6     race = as.numeric(factor(race, levels = c("white", "Black", "Other"))),
7     marital_status = as.numeric(factor(marital_status, levels = c("Married",
8 "Divorced", "Single", "Widowed", "Separated"))),
9     t_stage = as.numeric(factor(t_stage, levels = c("T1", "T2", "T3",
10 "T4"))),
11     n_stage = as.numeric(factor(n_stage, levels = c("N1", "N2", "N3"))),
12     x6th_stage = as.numeric(factor(x6th_stage, levels = c("IIA", "IIIA",
13 "IIIC", "IIB", "IIIB"))),
14     differentiate = as.numeric(factor(differentiate, levels = c("Poorly
15 differentiated", "Moderately differentiated", "Well differentiated",
16 "Undifferentiated"))),
17     grade = as.numeric(factor(grade, levels = c("3", "2", "1", "anaplastic;
18 Grade IV"))),
19     a_stage = as.numeric(factor(a_stage, levels = c("Regional",
20 "Distant"))),
21     estrogen_status = as.numeric(factor(estrogen_status, levels =
22 c("Positive", "Negative"))),
23     progesterone_status = as.numeric(factor(progesterone_status, levels =
24 c("Positive", "Negative"))),
25     status = as.numeric(factor(status, levels = c("Alive", "Dead")))
26   ) |>
27   rename(ms = "marital_status",
28     t_s = "t_stage",
29     n_s = "n_stage",
30     x6_s = "x6th_stage",
31     a_s = "a_stage",
32     diff = "differentiate",
33     est = "estrogen_status",
34     pro = "progesterone_status",
35     rne = "regional_node_examined",
36     rnp = "regional_node_positive"
37   )
```

```

1 ## Rows: 4024 Columns: 16
2 ## — Column specification
3
4 ## Delimiter: ","
5 ## chr (11): Race, Marital Status, T Stage, N Stage, 6th Stage,
6 ## differentiate, ...
7 ## dbl (5): Age, Tumor Size, Regional Node Examined, Regino1 Node Positive,
8 ## Su...
9 ##
10 ## I use `spec()` to retrieve the full column specification for this data.
11 ## I specify the column types or set `show_col_types = FALSE` to quiet this
12 ## message.

```

```

1 variablesummary =
2   lapply(heart_df[,3:16], table)

```

1. outcome: 1 continuous; 1 binary
 2. predictors: 4 continuous; 10 categorical

Examine predictors

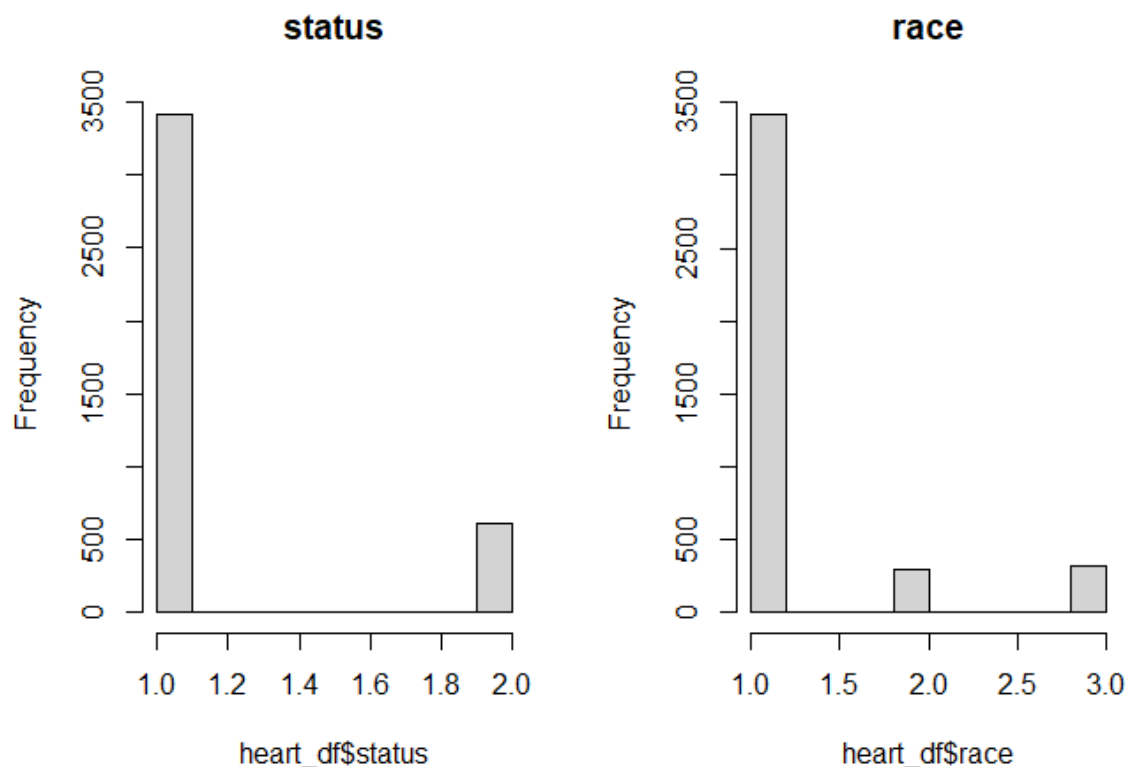
For categorical variable

Graphic plot

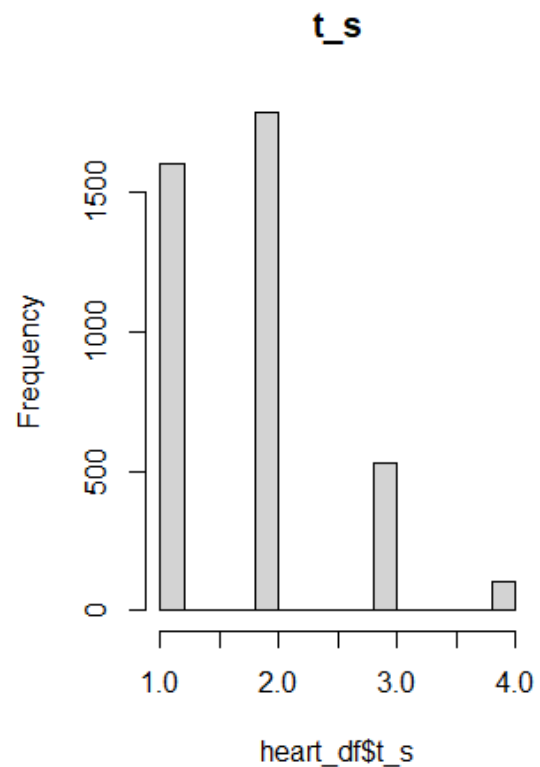
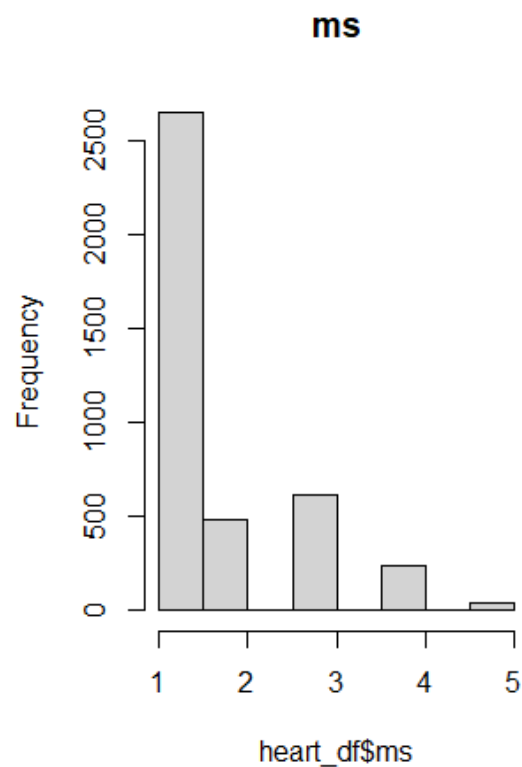
```

1 # Histograms for each categorical variables
2 par(mfrow=c(1,2))
3 hist(heart_df$status, main='status')
4 hist(heart_df$race, main='race')

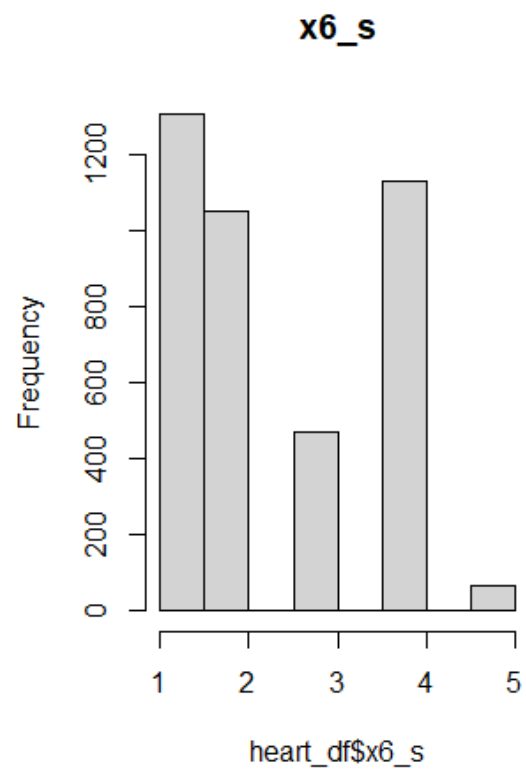
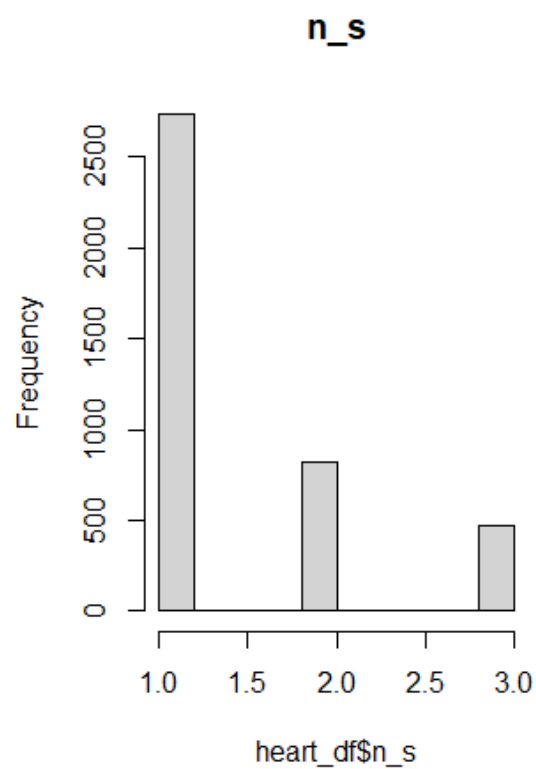
```



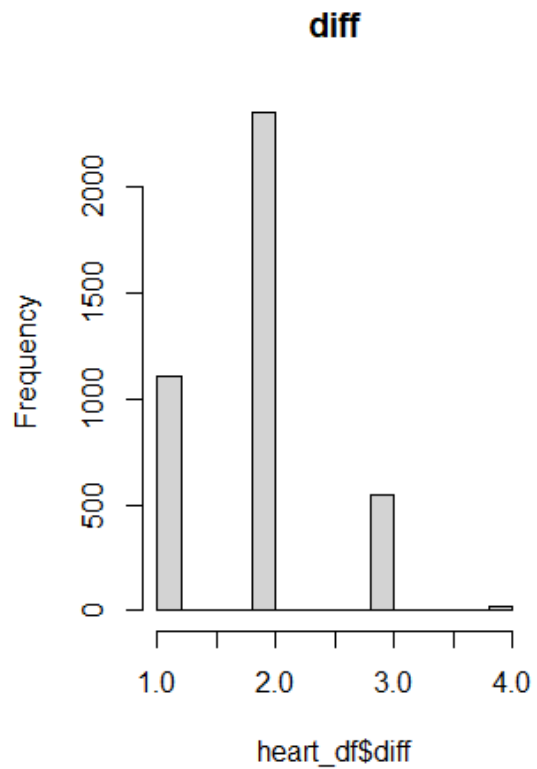
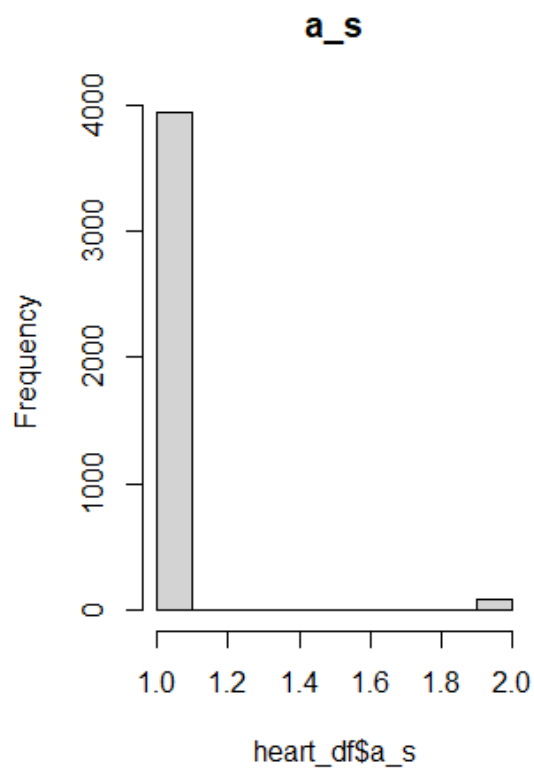
```
1 hist(heart_df$ms, main='ms')
2 hist(heart_df$t_s, main='t_s')
```



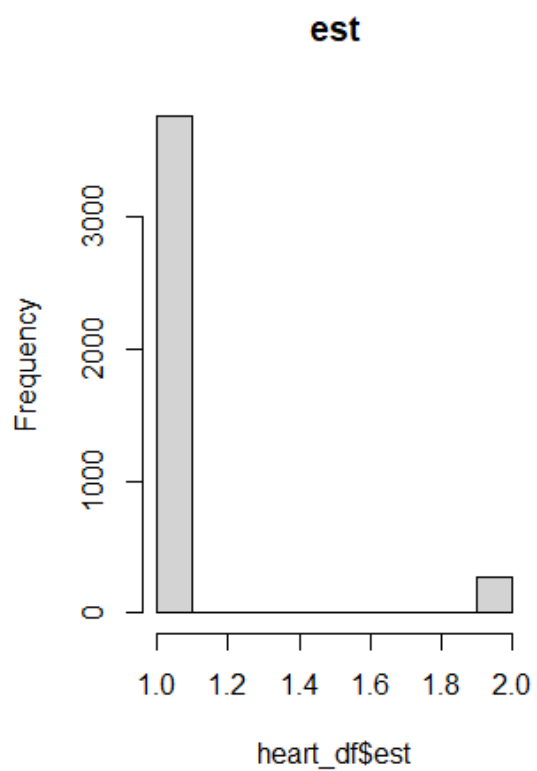
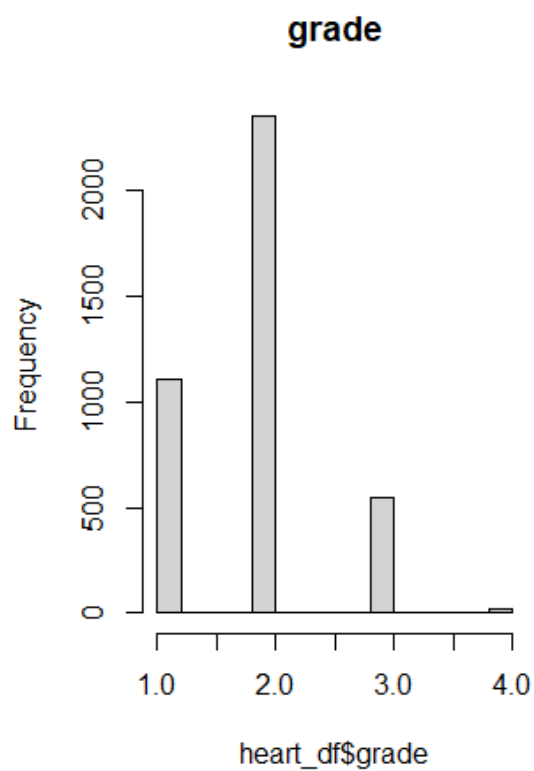
```
1 hist(heart_df$n_s, main='n_s')
2 hist(heart_df$x6_s, main='x6_s')
```



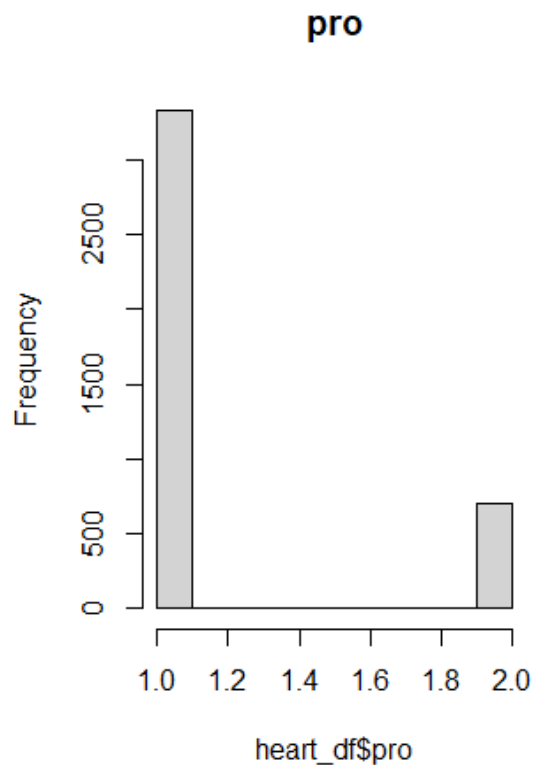
```
1 hist(heart_df$a_s, main='a_s')
2 hist(heart_df$diff, main='diff')
```



```
1 hist(heart_df$grade, main='grade')
2 hist(heart_df$est, main='est')
```



```
1 hist(heart_df$pro, main='pro')
```



most categorical variables have a extremely skewed.
 And I find that some variables like `grade` has a very small level. I
 will focus on that in the following steps.

Char-squared test

```
1 heart_df_2 =
2   read_csv("./Project_2_data.csv") |>
3   janitor::clean_names() |>
4   relocate(survival_months, status) |>
5   dplyr::select("race", "marital_status", "t_stage", "n_stage", "x6th_stage",
6     "differentiate", "grade", "a_stage", "estrogen_status",
7     "progesterone_status") |>
8   mutate_all(as.factor)
```

```
1 ## Rows: 4024 Columns: 16
2 ## — Column specification
3
4 ## Delimiter: ","
5 ## chr (11): Race, Marital Status, T Stage, N Stage, 6th Stage,
6   differentiate, ...
7 ## dbl (5): Age, Tumor Size, Regional Node Examined, Regional Node Positive,
8   Su...
9 ##
10 ## I use `spec()` to retrieve the full column specification for this data.
11 ## I specify the column types or set `show_col_types = FALSE` to quiet this
12   message.
```

```
1 result_table =
2   data.frame(Variable1 = character(),
3             Variable2 = character(),
```

```

4         ChiSquare = numeric(),
5         PValue = numeric(),
6         stringsAsFactors = TRUE)
7
8
9 for (col1 in names(heart_df_2)[1:(ncol(heart_df_2) - 1)]) {
10   for (col2 in names(heart_df_2)[(match(col1, names(heart_df_2)) +
11     1):ncol(heart_df_2)]) {
12
13     contingency_table <- table(heart_df_2[[col1]], heart_df_2[[col2]])
14
15     chi_sq_test_result <- chisq.test(contingency_table,
16                                     correct = T)
17
18
19     variable1 <- col1
20     variable2 <- col2
21     chi_square <- chi_sq_test_result$statistic
22     p_value <- chi_sq_test_result$p.value
23
24
25     result_table <- rbind(result_table, data.frame(variable1, variable2,
26     chi_square, p_value))
27   }
28 }
29
30 original_row_names = rownames(result_table)
31 new_row_names <- as.character(1:nrow(result_table))
32 rownames(result_table) <- new_row_names
33
34 result_table =
35   result_table |>
36   arrange(desc(p_value))
37
38 knitr::kable(result_table)

```

	variable1	variable2	chi_square	p_value
7	race	a_stage	3.069776e-01	0.8577104
4	race	x6th_stage	8.839562e+00	0.3560107
2	race	t_stage	8.462431e+00	0.2061430
3	race	n_stage	6.079684e+00	0.1932759
10	marital_status	t_stage	1.712041e+01	0.1451239
44	a_stage	progesterone_status	2.382806e+00	0.1226770
15	marital_status	a_stage	7.658489e+00	0.1049203
16	marital_status	estrogen_status	7.690336e+00	0.1036033
13	marital_status	differentiate	1.915564e+01	0.0848409
14	marital_status	grade	1.915564e+01	0.0848409
9	race	progesterone_status	5.043148e+00	0.0803331
12	marital_status	x6th_stage	2.810804e+01	0.0306927
17	marital_status	progesterone_status	1.104690e+01	0.0260420
37	differentiate	a_stage	1.057708e+01	0.0142471
40	grade	a_stage	1.057708e+01	0.0142471
11	marital_status	n_stage	2.235252e+01	0.0043030
24	t_stage	progesterone_status	1.380823e+01	0.0031781
8	race	estrogen_status	1.340900e+01	0.0012254
23	t_stage	estrogen_status	1.954986e+01	0.0002104
5	race	differentiate	2.790280e+01	0.0000980
6	race	grade	2.790280e+01	0.0000980
43	a_stage	estrogen_status	1.558922e+01	0.0000787
35	x6th_stage	progesterone_status	4.248544e+01	0.0000000
30	n_stage	progesterone_status	3.684600e+01	0.0000000
29	n_stage	estrogen_status	4.252308e+01	0.0000000
34	x6th_stage	estrogen_status	5.200461e+01	0.0000000
20	t_stage	differentiate	9.095689e+01	0.0000000
21	t_stage	grade	9.095689e+01	0.0000000
26	n_stage	differentiate	1.155011e+02	0.0000000
27	n_stage	grade	1.155011e+02	0.0000000

	variable1	variable2	chi_square	p_value
1	race	marital_status	1.379574e+02	0.0000000
31	x6th_stage	differentiate	1.580412e+02	0.0000000
32	x6th_stage	grade	1.580412e+02	0.0000000
39	differentiate	progesterone_status	1.478031e+02	0.0000000
42	grade	progesterone_status	1.478031e+02	0.0000000
38	differentiate	estrogen_status	2.174167e+02	0.0000000
41	grade	estrogen_status	2.174167e+02	0.0000000
18	t_stage	n_stage	3.234137e+02	0.0000000
28	n_stage	a_stage	3.555764e+02	0.0000000
22	t_stage	a_stage	5.832589e+02	0.0000000
33	x6th_stage	a_stage	7.291926e+02	0.0000000
45	estrogen_status	progesterone_status	1.054843e+03	0.0000000
19	t_stage	x6th_stage	6.784079e+03	0.0000000
25	n_stage	x6th_stage	6.686834e+03	0.0000000
36	differentiate	grade	1.207200e+04	0.0000000

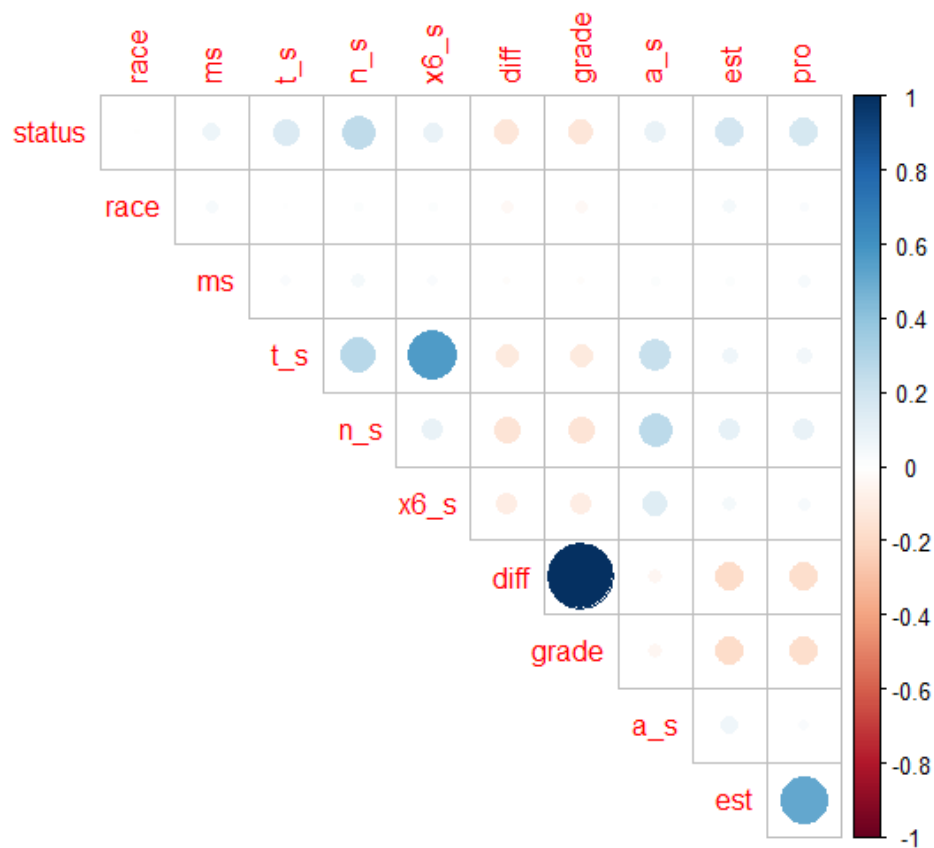
There're many categorical variables correlate with others. Particularly, difference and grade should preserve one. I delete `differentiate`. What's more, delete `x6th_stage`, `est` on account of their extremely small p-values.

Correlation for categorical variables

```

1 par(mfrow = c(1, 1), mar = c(5, 4, 4, 2) + 0.1)
2 cor_cate = cor(heart_df[,c(2, 4:11, 13:14)])
3
4 corrpplot(cor_cate, type = "upper", diag = FALSE, mar = c(0, 0, 0, 0))

```

There is apparent linear correlation among variables between diff and grade; t_s and x6_s; est and pro. with the outcome we get from chi-squared test, I delete x6th_stage, est and diff variable for the further study.

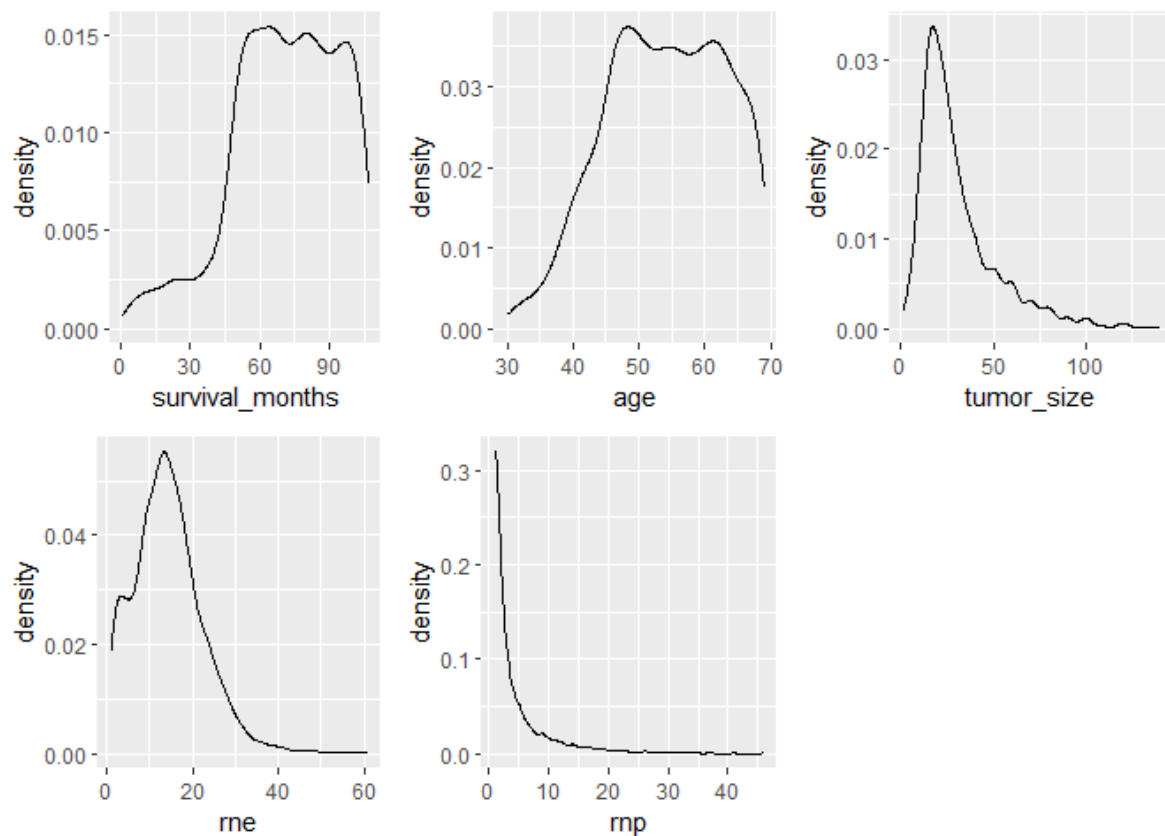
For continuous variable

Graphic plot

```

1 # Boxplots for each continuous variable
2 par(mfrow=c(2,3))
3 boxplot(heart_df$survival_months, main='survival_months')
4 boxplot(heart_df$age, main='age')
5 boxplot(heart_df$tumor_size, main='tumor_size')
6 boxplot(heart_df$rne, main='rne')
7 boxplot(heart_df$rnp, main='rnp')

```

Among all the 5 continuous variables, we can find that `age` is the only predictor which is normally distributed. The other three variables are not normally distributed. Because the logistic regression doesn't need the predictors to be normally distributed, I didn't transform these predictors.

Rank-sum test and t-test

```
1 result1 = t.test(age ~ status, data = heart_df)
2 result2 = wilcox.test(tumor_size ~ status, data = heart_df)
3 result3 = wilcox.test(rne ~ status, data = heart_df)
4 result4 = wilcox.test(rnp ~ status, data = heart_df)
5
6 broom::tidy(result1)
```

```
1 ## # A tibble: 1 × 10
2 ##   estimate estimate1 estimate2 statistic p.value parameter conf.low
3 ##   <dbl>     <dbl>     <dbl>     <dbl>   <dbl>     <dbl>     <dbl>
4 ## 1    -1.39      53.8      55.2     -3.32 0.000931    809.     -2.21
5 ##   -0.570
6 ## # i 2 more variables: method <chr>, alternative <chr>
```

```
1 broom::tidy(result2)
```

```

1 ## # A tibble: 1 × 4
2 ##   statistic p.value method
3 ##         <dbl>   <dbl> <chr>          <chr>
4 ## 1      813236 4.85e-19 wilcoxon rank sum test with continuity correct...
   two.sided

```

```

1 broom::tidy(result3)

```

```

1 ## # A tibble: 1 × 4
2 ##   statistic p.value method
3 ##         <dbl>   <dbl> <chr>          <chr>
4 ## 1    1001152 0.0673 wilcoxon rank sum test with continuity correcti...
   two.sided

```

```

1 broom::tidy(result4)

```

```

1 ## # A tibble: 1 × 4
2 ##   statistic p.value method
3 ##         <dbl>   <dbl> <chr>          <chr>
4 ## 1    692562. 6.82e-44 wilcoxon rank sum test with continuity correct...
   two.sided

```

I use two-sample t test to check age and rank-sum test check other three variables. `rne` has a different median compared between two groups stratified by `status`. But, we can't delete this variable just on account of its rank-sum test outcome, I will focus on this variable in the further study.

Check the dataset again

```

1 heart_df_log =
2   read_csv("./Project_2_data.csv") |>
3   janitor::clean_names() |>
4   relocate(survival_months, status) |>
5   filter(!grade %in% c("anaplastic; Grade IV")) |>
6   mutate(
7     race = factor(race, levels = c("white", "Black", "Other")),
8     marital_status = factor(marital_status, levels = c("Married",
9 "Divorced", "Single", "Widowed", "Separated")),
10    t_stage = factor(t_stage, levels = c("T1", "T2", "T3", "T4")),
11    n_stage = factor(n_stage, levels = c("N1", "N2", "N3")),
12    grade = factor(grade, levels = c("3", "2", "1")),
13    a_stage = factor(a_stage, levels = c("Regional", "Distant")),
14    progesterone_status = factor(progesterone_status, levels = c("Positive",
15 "Negative")),
16    status = factor(status, levels = c("Alive", "Dead"))) |>

```

```
15 dplyr::select(-survival_months, -differentiate, -x6th_stage, -
    estrogen_status)
```

```
1 ## Rows: 4024 Columns: 16
2 ## — Column specification
3 ## Delimiter: ","
4 ## chr (11): Race, Marital Status, T Stage, N Stage, 6th Stage,
   differentiate, ...
5 ## dbl (5): Age, Tumor Size, Regional Node Examined, Regional Node Positive,
   Su...
6 ##
7 ## I use `spec()` to retrieve the full column specification for this data.
8 ## I specify the column types or set `show_col_types = FALSE` to quiet this
   message.
```

```
1 summary(heart_df_log)
```

```
1 ##      status      age      race      marital_status t_stage
   n_stage
2 ## Alive:3398   Min.    :30.00   white:3397   Married   :2632   T1:1598
   N1:2722
3 ## Dead : 607   1st Qu.:47.00   Black: 288   Divorced  : 485   T2:1781   N2:
   816
4 ##              Median :54.00   Other: 320   Single    : 610   T3: 526   N3:
   467
5 ##              Mean    :53.98              widowed  : 234   T4: 100
6 ##              3rd Qu.:61.00              separated:  44
7 ##              Max.    :69.00

8 ## grade      a_stage      tumor_size      progesterone_status
9 ## 3:1111   Regional:3913   Min.    : 1.00   Positive:3311
10 ## 2:2351   Distant : 92   1st Qu.: 16.00   Negative: 694
11 ## 1: 543              Median : 25.00
12 ##              Mean     : 30.41
13 ##              3rd Qu.: 38.00
14 ##              Max.     :140.00
15 ## regional_node_examined regional_node_positive
16 ## Min.    : 1.00      Min.    : 1.000
17 ## 1st Qu.: 9.00      1st Qu.: 1.000
18 ## Median :14.00      Median : 2.000
19 ## Mean    :14.36      Mean    : 4.149
20 ## 3rd Qu.:19.00      3rd Qu.: 5.000
21 ## Max.    :61.00      Max.    :46.000
```

```
1 str(heart_df_log)
```

```

1 ## tibble [4,005 × 12] (S3: tbl_df/tbl/data.frame)
2 ## $ status : Factor w/ 2 levels "Alive","Dead": 1 1 1 1 1 1
  1 2 1 1 ...
3 ## $ age : num [1:4005] 68 50 58 58 47 51 51 40 40 69 ...
4 ## $ race : Factor w/ 3 levels "White","Black",...: 1 1 1 1
  1 1 1 1 1 1 ...
5 ## $ marital_status : Factor w/ 5 levels "Married","Divorced",...: 1
  1 2 1 1 3 1 1 2 1 ...
6 ## $ t_stage : Factor w/ 4 levels "T1","T2","T3",...: 1 2 3 1
  2 1 1 2 4 4 ...
7 ## $ n_stage : Factor w/ 3 levels "N1","N2","N3": 1 2 3 1 1 1
  1 1 3 3 ...
8 ## $ grade : Factor w/ 3 levels "3","2","1": 1 2 2 1 1 2 3
  2 1 3 ...
9 ## $ a_stage : Factor w/ 2 levels "Regional","Distant": 1 1 1
  1 1 1 1 1 1 2 ...
10 ## $ tumor_size : num [1:4005] 4 35 63 18 41 20 8 30 103 32 ...
11 ## $ progesterone_status : Factor w/ 2 levels "Positive","Negative": 1 1
  1 1 1 1 1 1 1 1 ...
12 ## $ regional_node_examined: num [1:4005] 24 14 14 2 3 18 11 9 20 21 ...
13 ## $ reginol_node_positive : num [1:4005] 1 5 7 1 1 2 1 1 18 12 ...

```

according to the procedure, I get the final dataset to make a logistic regression. 1, delete 4 level grade data for its anaplastic mode. 2, delete three categorical variables, namely `differentiate`, `x6th_stage`, `estrogen_status`. 3, delete the `survival_months` cause it's not the variable of our interest. 4, I'll attach more importance on the variable `regional_node_examined`, so it will be tested by Wald test to check whether it should be deleted.

Examine interaction

```

1 glm_result_1 = glm(data = heart_df_log,
2                   status ~ a_stage * reginol_node_positive,
3                   binomial(link = 'logit'))
4 broom::tidy(glm_result_1)

```

```

1 ## # A tibble: 4 × 5
2 ##   term                                estimate std.error statistic
  p.value
3 ##   <chr>                                <dbl>      <dbl>      <dbl>
  <dbl>
4 ## 1 (Intercept)                        -2.28      0.0629     -36.3  6.51e-
  288
5 ## 2 a_stageDistant                      1.81      0.397      4.55  5.25e-
  6
6 ## 3 reginol_node_positive                0.107     0.00774     13.8  1.48e-
  43
7 ## 4 a_stageDistant:reginol_node_positive -0.108     0.0286     -3.79  1.50e-
  4

```

```

1 glm_result_2 = glm(data = heart_df_log,
2                     status ~ regional_node_examined * reginol_node_positive,
3                     binomial(link = 'logit'))
4 broom::tidy(glm_result_2)

```

```

1 ## # A tibble: 4 × 5
2 ##   term                                estimate std.error statistic
3 ##   <chr>                                <dbl>      <dbl>      <dbl>
4 ## 1 (Intercept)                        -2.13      0.120      -17.8
5 ## 2 regional_node_examined             -0.0218    0.00790     -2.76
6 ## 3 reginol_node_positive               0.199     0.0196     10.2
7 ## 4 regional_node_examined:reginol_node_pos... -0.00300  0.000744    -4.02

```

I check all the 2 variable group to find whether there exists a correlation effect, and find `a_stage` and `reginol_node_positive`'s interaction p-value is 1.50e-4. `regional_node_examined` and `reginol_node_positive`'s interaction p-value is 5.71e-5. Their interaction events play vital roles in our model. Therefore, I will put these two interaction event in my following models.

Model variable select

1. Stepwise

build model

```

1 model =
2   glm(status ~ . + a_stage * reginol_node_positive + regional_node_examined *
3       reginol_node_positive,
4       data = heart_df_log,
5       family = binomial(link = 'logit'))
6 summary(model)

```

```

1 ##
2 ## Call:
3 ## glm(formula = status ~ . + a_stage * reginol_node_positive +
4 ##     regional_node_examined * reginol_node_positive, family =
5 ##     binomial(link = "logit"),
6 ##     data = heart_df_log)
7 ## Coefficients:
8 ##                                Estimate Std. Error z
9 ## (Intercept)                    -3.3157333   0.3565548
10 ##                                -9.299

```

10	## age	0.0214016	0.0055957
	3.825		
11	## raceBlack	0.5285528	0.1617692
	3.267		
12	## raceOther	-0.4105070	0.2015324
	-2.037		
13	## marital_statusDivorced	0.1979062	0.1413430
	1.400		
14	## marital_statusSingle	0.1233095	0.1354483
	0.910		
15	## marital_statuswidowed	0.2141907	0.1925412
	1.112		
16	## marital_statusSeparated	0.8660280	0.3712285
	2.333		
17	## t_stageT2	0.3805790	0.1313527
	2.897		
18	## t_stageT3	0.5016913	0.2663068
	1.884		
19	## t_stageT4	0.9958936	0.3197436
	3.115		
20	## n_stageN2	0.3312383	0.1474584
	2.246		
21	## n_stageN3	0.3904381	0.2757280
	1.416		
22	## grade2	-0.4487602	0.1030829
	-4.353		
23	## grade1	-0.9931654	0.1916765
	-5.181		
24	## a_stageDistant	0.7656274	0.4833319
	1.584		
25	## tumor_size	0.0003454	0.0039406
	0.088		
26	## progesterone_statusNegative	0.8338434	0.1082163
	7.705		
27	## regional_node_examined	-0.0316316	0.0084370
	-3.749		
28	## reginol_node_positive	0.1125016	0.0324023
	3.472		
29	## a_stageDistant:reginol_node_positive	-0.0542828	0.0319926
	-1.697		
30	## regional_node_examined:reginol_node_positive	-0.0007918	0.0008745
	-0.905		
31	##	Pr(> z)	
32	## (Intercept)	< 2e-16	***
33	## age	0.000131	***
34	## raceBlack	0.001086	**
35	## raceOther	0.041657	*
36	## marital_statusDivorced	0.161458	
37	## marital_statusSingle	0.362622	
38	## marital_statuswidowed	0.265949	
39	## marital_statusSeparated	0.019655	*
40	## t_stageT2	0.003763	**
41	## t_stageT3	0.059581	.
42	## t_stageT4	0.001842	**
43	## n_stageN2	0.024684	*
44	## n_stageN3	0.156768	
45	## grade2	1.34e-05	***
46	## grade1	2.20e-07	***


```

47 ## a_stageDistant 0.113180
48 ## tumor_size 0.930156
49 ## progesterone_statusNegative 1.30e-14 ***
50 ## regional_node_examined 0.000177 ***
51 ## reginol_node_positive 0.000517 ***
52 ## a_stageDistant:reginol_node_positive 0.089748 .
53 ## regional_node_examined:reginol_node_positive 0.365255
54 ## ---
55 ## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
56 ##
57 ## (Dispersion parameter for binomial family taken to be 1)
58 ##
59 ## Null deviance: 3407.5 on 4004 degrees of freedom
60 ## Residual deviance: 2946.1 on 3983 degrees of freedom
61 ## AIC: 2990.1
62 ##
63 ## Number of Fisher Scoring iterations: 5

```

```

1 | model_step = step(model, direction = "both")

```

```

1 ## Start: AIC=2990.1
2 ## status ~ age + race + marital_status + t_stage + n_stage + grade +
3 ## a_stage + tumor_size + progesterone_status + regional_node_examined +
4 ## reginol_node_positive + a_stage * reginol_node_positive +
5 ## regional_node_examined * reginol_node_positive
6 ##
7 ##
8 ## Df Deviance AIC
9 ## - tumor_size 1 2946.1 2988.1
10 ## - regional_node_examined:reginol_node_positive 1 2946.9 2988.9
11 ## - marital_status 4 2953.5 2989.5
12 ## <none> 2946.1 2990.1
13 ## - a_stage:reginol_node_positive 1 2949.0 2991.0
14 ## - n_stage 2 2951.1 2991.1
15 ## - t_stage 3 2958.9 2996.9
16 ## - race 2 2961.9 3001.9
17 ## - age 1 2960.9 3002.9
18 ## - grade 2 2981.9 3021.9
19 ## - progesterone_status 1 3002.6 3044.6
20 ##
21 ## Step: AIC=2988.11
22 ## status ~ age + race + marital_status + t_stage + n_stage + grade +
23 ## a_stage + progesterone_status + regional_node_examined +
24 ## reginol_node_positive + a_stage:reginol_node_positive +
25 ## regional_node_examined:reginol_node_positive
26 ##
27 ## Df Deviance AIC
28 ## - regional_node_examined:reginol_node_positive 1 2946.9 2986.9
29 ## - marital_status 4 2953.5 2987.5
30 ## <none> 2946.1 2988.1
31 ## - a_stage:reginol_node_positive 1 2949.0 2989.0
32 ## - n_stage 2 2951.2 2989.2
33 ## + tumor_size 1 2946.1 2990.1
34 ## - race 2 2962.0 3000.0
35 ## - age 1 2960.9 3000.9
36 ## - t_stage 3 2969.1 3005.1
37 ## - grade 2 2982.0 3020.0

```

```

36 ## - progesterone_status 1 3002.7 3042.7
37 ##
38 ## Step: AIC=2986.93
39 ## status ~ age + race + marital_status + t_stage + n_stage + grade +
40 ## a_stage + progesterone_status + regional_node_examined +
41 ## reginol_node_positive + a_stage:reginol_node_positive
42 ##
43 ## Df Deviance AIC
44 ## - marital_status 4 2954.4 2986.4
45 ## <none> 2946.9 2986.9
46 ## - a_stage:reginol_node_positive 1 2949.6 2987.6
47 ## + regional_node_examined:reginol_node_positive 1 2946.1 2988.1
48 ## + tumor_size 1 2946.9 2988.9
49 ## - n_stage 2 2956.3 2992.3
50 ## - race 2 2962.8 2998.8
51 ## - age 1 2961.7 2999.7
52 ## - t_stage 3 2971.3 3005.3
53 ## - regional_node_examined 1 2973.0 3011.0
54 ## - grade 2 2982.9 3018.9
55 ## - progesterone_status 1 3003.4 3041.4
56 ##
57 ## Step: AIC=2986.4
58 ## status ~ age + race + t_stage + n_stage + grade + a_stage +
59 ## progesterone_status +
60 ## regional_node_examined + reginol_node_positive +
61 ## a_stage:reginol_node_positive
62 ## Df Deviance AIC
63 ## <none> 2954.4 2986.4
64 ## + marital_status 4 2946.9 2986.9
65 ## - a_stage:reginol_node_positive 1 2956.9 2986.9
66 ## + regional_node_examined:reginol_node_positive 1 2953.5 2987.5
67 ## + tumor_size 1 2954.4 2988.4
68 ## - n_stage 2 2963.6 2991.6
69 ## - age 1 2970.9 3000.9
70 ## - race 2 2973.2 3001.2
71 ## - t_stage 3 2979.5 3005.5
72 ## - regional_node_examined 1 2981.0 3011.0
73 ## - grade 2 2989.9 3017.9
74 ## - progesterone_status 1 3012.5 3042.5

```

```
1 summary(model_step)
```

```

1 ##
2 ## Call:
3 ## glm(formula = status ~ age + race + t_stage + n_stage + grade +
4 ## a_stage + progesterone_status + regional_node_examined +
5 ## reginol_node_positive + a_stage:reginol_node_positive, family =
6 ## binomial(link = "logit"),
7 ## data = heart_df_log)
8 ##
9 ## Coefficients:
10 ## Estimate Std. Error z value Pr(>|z|)
11 ## (Intercept) -3.207333 0.332597 -9.643 < 2e-
12 ## 16 ***

```

```

11 ## age 0.021862 0.005419 4.034 5.48e-
12 05 ***
13 ## raceBlack 0.580034 0.158385 3.662
14 0.000250 ***
15 ## raceOther -0.418754 0.201019 -2.083
16 0.037237 *
17 ## t_stageT2 0.406565 0.112837 3.603
18 0.000314 ***
19 ## t_stageT3 0.532044 0.148621 3.580
20 0.000344 ***
21 ## t_stageT4 1.042480 0.268893 3.877
22 0.000106 ***
23 ## n_stageN2 0.389858 0.129955 3.000
24 0.002700 **
25 ## n_stageN3 0.505349 0.241470 2.093
26 0.036367 *
27 ## grade2 -0.449129 0.102872 -4.366 1.27e-
28 05 ***
29 ## grade1 -0.980279 0.190720 -5.140 2.75e-
30 07 ***
31 ## a_stageDistant 0.726733 0.479490 1.516
32 0.129611
33 ## progesterone_statusNegative 0.844014 0.107933 7.820 5.29e-
34 15 ***
35 ## regional_node_examined -0.036050 0.007232 -4.985 6.20e-
36 07 ***
37 ## reginol_node_positive 0.088394 0.015965 5.537 3.08e-
38 08 ***
39 ## a_stageDistant:reginol_node_positive -0.050505 0.031625 -1.597
40 0.110274
41 ## ---
42 ## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
43 ##
44 ## (Dispersion parameter for binomial family taken to be 1)
45 ##
46 ## Null deviance: 3407.5 on 4004 degrees of freedom
47 ## Residual deviance: 2954.4 on 3989 degrees of freedom
48 ## AIC: 2986.4
49 ##
50 ## Number of Fisher Scoring iterations: 5

```

```

1 | waldtest(model_step, "regional_node_examined")

```

```

1  ## wald test
2  ##
3  ## Model 1: status ~ age + race + t_stage + n_stage + grade + a_stage +
  progesterone_status +
4  ##      regional_node_examined + reginol_node_positive +
  a_stage:reginol_node_positive
5  ## Model 2: status ~ age + race + t_stage + n_stage + grade + a_stage +
  progesterone_status +
6  ##      reginol_node_positive + a_stage:reginol_node_positive
7  ##   Res.Df Df      F    Pr(>F)
8  ## 1    3989
9  ## 2    3990 -1 24.848 6.467e-07 ***
10 ## ---
11 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

1. stepwise: delete the variable `tumor_size`,
`regional_node_examined:reginol_node_positive`, `marital_status` 2. Wald
test shows a p-value 6.467e-07, which demonstrates that we should
preserve the variable `regional_node_examined` though it has a different
median between two groups stratified by status.

Attention here please! In this process, the
stepwise help me delete the interaction event
`regional_node_examined:reginol_node_positive` but preserve the
interaction `a_stage:reginol_node_positive` which means the latter one
has a larger influence upon the whole model.

the model

```

1  formula1 =
2    as.formula(
3      status ~ age + race + t_stage + n_stage + grade + a_stage +
      progesterone_status + regional_node_examined + reginol_node_positive +
      a_stage:reginol_node_positive)
4
5  model = glm(formula1,
6              data = heart_df_log,
7              family = binomial(link = 'logit'))
8
9  summary(model)

```

```

1  ##
2  ## Call:
3  ## glm(formula = formula1, family = binomial(link = "logit"), data =
  heart_df_log)
4  ##
5  ## Coefficients:
6  ##
  Estimate Std. Error z value
Pr(>|z|)
7  ## (Intercept)      -3.207333    0.332597  -9.643  < 2e-
  16 ***
8  ## age              0.021862    0.005419   4.034 5.48e-
  05 ***
9  ## raceBlack        0.580034    0.158385   3.662
  0.000250 ***

```

```

10 ## raceOther -0.418754 0.201019 -2.083
0.037237 *
11 ## t_stageT2 0.406565 0.112837 3.603
0.000314 ***
12 ## t_stageT3 0.532044 0.148621 3.580
0.000344 ***
13 ## t_stageT4 1.042480 0.268893 3.877
0.000106 ***
14 ## n_stageN2 0.389858 0.129955 3.000
0.002700 **
15 ## n_stageN3 0.505349 0.241470 2.093
0.036367 *
16 ## grade2 -0.449129 0.102872 -4.366 1.27e-
05 ***
17 ## grade1 -0.980279 0.190720 -5.140 2.75e-
07 ***
18 ## a_stageDistant 0.726733 0.479490 1.516
0.129611
19 ## progesterone_statusNegative 0.844014 0.107933 7.820 5.29e-
15 ***
20 ## regional_node_examined -0.036050 0.007232 -4.985 6.20e-
07 ***
21 ## reginol_node_positive 0.088394 0.015965 5.537 3.08e-
08 ***
22 ## a_stageDistant:reginol_node_positive -0.050505 0.031625 -1.597
0.110274
23 ## ---
24 ## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
25 ##
26 ## (Dispersion parameter for binomial family taken to be 1)
27 ##
28 ## Null deviance: 3407.5 on 4004 degrees of freedom
29 ## Residual deviance: 2954.4 on 3989 degrees of freedom
30 ## AIC: 2986.4
31 ##
32 ## Number of Fisher Scoring iterations: 5

```

```
1 check_collinearity(model)
```

```

1 ## # Check for Multicollinearity
2 ##
3 ## Low Correlation
4 ##
5 ##          Term VIF   VIF 95% CI Increased SE Tolerance
6 ##          age 1.03 [1.01, 1.09]         1.01      0.97
7 ##          race 1.02 [1.00, 1.10]         1.01      0.98
8 ##          t_stage 1.36 [1.31, 1.41]         1.16      0.74
9 ##          n_stage 4.10 [3.88, 4.33]         2.02      0.24
10 ##          grade 1.08 [1.05, 1.12]         1.04      0.93
11 ##          a_stage 4.27 [4.04, 4.51]         2.07      0.23
12 ##          progesterone_status 1.04 [1.02, 1.09]         1.02      0.96
13 ##          regional_node_examined 1.48 [1.43, 1.55]         1.22      0.67
14 ##          reginol_node_positive 4.68 [4.43, 4.94]         2.16      0.21
15 ##          a_stage:reginol_node_positive 3.90 [3.70, 4.12]         1.98      0.26
16 ## Tolerance 95% CI
17 ##          [0.92, 0.99]

```

```

18 ##      [0.91, 1.00]
19 ##      [0.71, 0.76]
20 ##      [0.23, 0.26]
21 ##      [0.89, 0.95]
22 ##      [0.22, 0.25]
23 ##      [0.92, 0.98]
24 ##      [0.65, 0.70]
25 ##      [0.20, 0.23]
26 ##      [0.24, 0.27]

```

We find our first model and check its collinearity, the model's performance is good.

2. Random forest

build model

```

1 set.seed(123)
2 model =
3   randomForest(status ~ . + a_stage * reginol_node_positive +
4     regional_node_examined * reginol_node_positive,
5     data = heart_df_log)
6 summary(model)

```

```

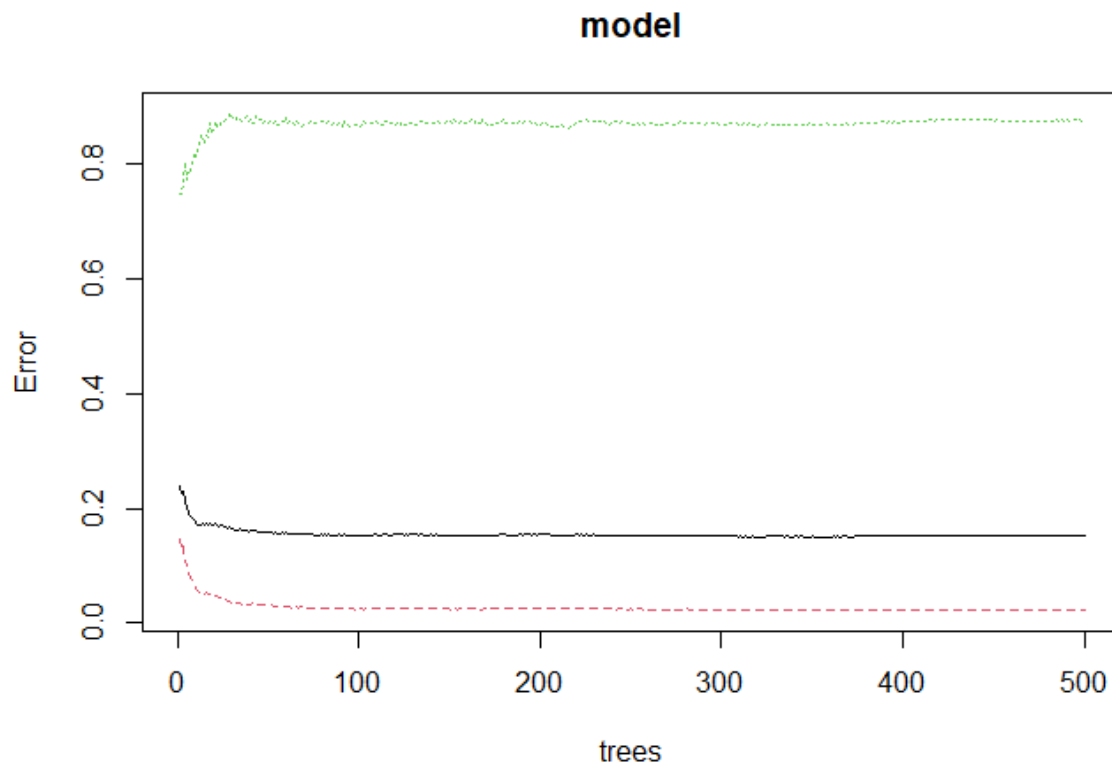
1 ##              Length Class  Mode
2 ## call              3  -none- call
3 ## type              1  -none- character
4 ## predicted         4005  factor numeric
5 ## err.rate          1500  -none- numeric
6 ## confusion          6  -none- numeric
7 ## votes             8010  matrix numeric
8 ## oob.times          4005  -none- numeric
9 ## classes           2  -none- character
10 ## importance         11  -none- numeric
11 ## importancesD        0  -none- NULL
12 ## localImportance     0  -none- NULL
13 ## proximity           0  -none- NULL
14 ## ntree              1  -none- numeric
15 ## mtry               1  -none- numeric
16 ## forest             14  -none- list
17 ## y                  4005  factor numeric
18 ## test               0  -none- NULL
19 ## inbag              0  -none- NULL
20 ## terms              3  terms  call

```

```

1 plot(model)

```



```
1 forest_min = which.min(model$err.rate[,1])
```

the forest_min which is the best parameter in the regression is 320

```
1 model_forest =
2   randomForest(status ~ . + a_stage * reginol_node_positive +
3     regional_node_examined * reginol_node_positive,
4     data = heart_df_log,
5     ntree = forest_min,
6     mtry = 4)
7 knitr::kable(model_forest$confusion)
```

	Alive	Dead	class.error
Alive	3303	95	0.0279576
Dead	531	76	0.8747941

in the simulation process, we get a outcome: the model prediction ability towards alive is good ,but when it comes to the dead, the class.error reaches nearly 0.87 which is a large number.

```
1 importance(model_forest)
```

```

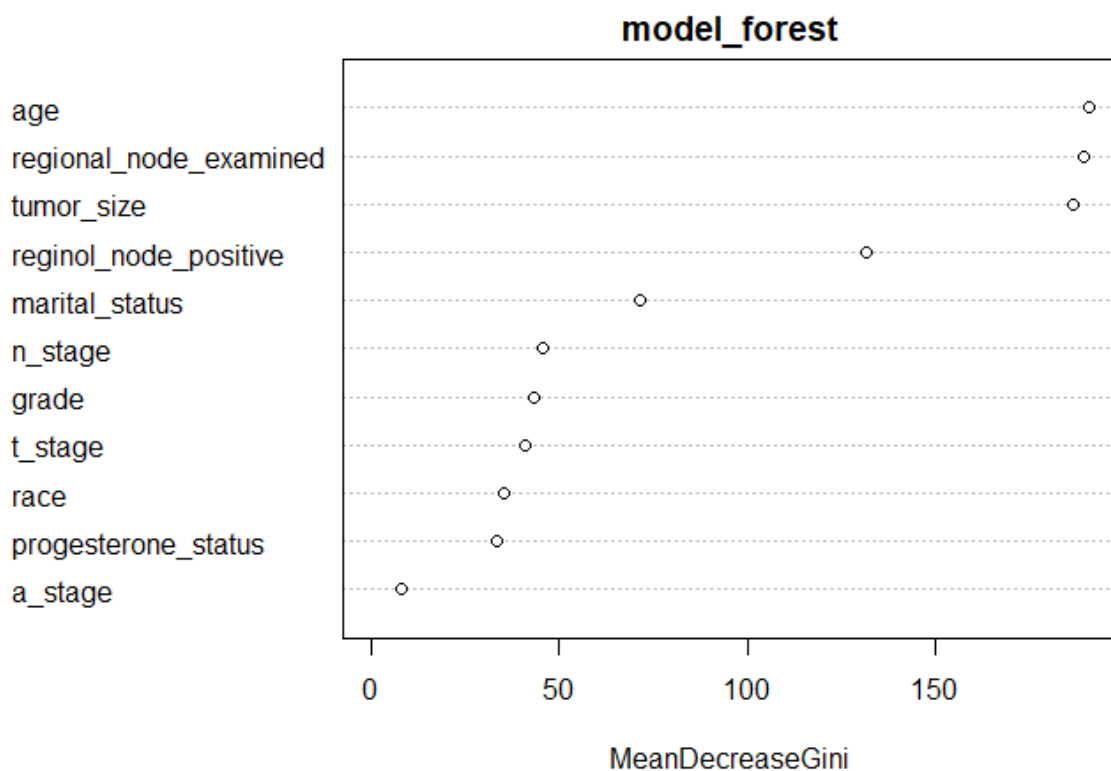
1  ##                                MeanDecreaseGini
2  ## age                            190.67222
3  ## race                           35.25882
4  ## marital_status                 71.36831
5  ## t_stage                        40.70758
6  ## n_stage                        45.74057
7  ## grade                         43.35355
8  ## a_stage                        8.17414
9  ## tumor_size                     186.30685
10 ## progesterone_status             33.25173
11 ## regional_node_examined         189.38867
12 ## reginol_node_positive          131.57251

```

```

1  par(mfrow = c(1, 1), mar = c(4, 4, 2, 1), oma = c(0, 0, 2, 0))
2  varImpPlot(model_forest)

```



```

1  par(mfrow = c(1, 1), mar = c(5, 4, 4, 2) + 0.1, oma = c(0, 0, 0, 0))

```

according to the picture, I find the first 5 most important variables have a mean decrease gini over 50, which is significantly larger than other variables, so I use these variables to build my final model.

the model


```

1 formula2 =
2   as.formula(
3     status ~ age + regional_node_examined + tumor_size + reginol_node_positive
4     + marital_status)
5
6 model = glm(formula2,
7             data = heart_df_log,
8             family = binomial(link = 'logit'))
9
10 summary(model)

```

```

1 ##
2 ## Call:
3 ## glm(formula = formula2, family = binomial(link = "logit"), data =
4 ## heart_df_log)
5 ## Coefficients:
6 ##              Estimate Std. Error z value Pr(>|z|)
7 ## (Intercept)    -3.248738   0.325798  -9.972  < 2e-16 ***
8 ## age              0.018089   0.005398   3.351 0.000805 ***
9 ## regional_node_examined -0.034997  0.006990  -5.007 5.52e-07 ***
10 ## tumor_size       0.010230   0.001977   5.175 2.28e-07 ***
11 ## reginol_node_positive  0.117507  0.008906  13.194 < 2e-16 ***
12 ## marital_statusDivorced  0.275010  0.136973   2.008 0.044668 *
13 ## marital_statusSingle   0.258083  0.129079   1.999 0.045563 *
14 ## marital_statusWidowed   0.318806  0.184381   1.729 0.083799 .
15 ## marital_statusSeparated  0.971231  0.353981   2.744 0.006074 **
16 ## ---
17 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
18 ##
19 ## (Dispersion parameter for binomial family taken to be 1)
20 ##
21 ##    Null deviance: 3407.5  on 4004  degrees of freedom
22 ## Residual deviance: 3124.2  on 3996  degrees of freedom
23 ## AIC: 3142.2
24 ##
25 ## Number of Fisher Scoring iterations: 5

```

```

1 check_collinearity(model)

```

```

1 ## # Check for Multicollinearity
2 ##
3 ## Low Correlation
4 ##
5 ##              Term  VIF   VIF 95% CI Increased SE Tolerance
6 ##              age  1.08 [1.05, 1.12]         1.04      0.93
7 ## regional_node_examined 1.47 [1.41, 1.53]         1.21      0.68
8 ##              tumor_size 1.05 [1.03, 1.10]         1.03      0.95
9 ## reginol_node_positive 1.51 [1.45, 1.57]         1.23      0.66
10 ## marital_status 1.07 [1.04, 1.12]         1.03      0.94
11 ## Tolerance 95% CI
12 ## [0.89, 0.95]
13 ## [0.65, 0.71]
14 ## [0.91, 0.97]
15 ## [0.64, 0.69]

```

```
16 ## [0.90, 0.96]
```

We find our second model and check its collinearity, the model's performance is good.

3. LASSO

build model

```
1 x = as.matrix(heart_df_log[,2:12])
2 y = as.matrix(heart_df_log[,1])
3
4 model_lasso =
5   glmnet(x, y, alpha = 1, family = "binomial")
6
7 model_lasso
```

```
1 ##
2 ## Call:  glmnet(x = x, y = y, family = "binomial", alpha = 1)
3 ##
4 ##      Df %Dev   Lambda
5 ##  1    0 0.00 0.091580
6 ##  2    1 1.23 0.083440
7 ##  3    1 2.16 0.076030
8 ##  4    1 2.88 0.069270
9 ##  5    1 3.45 0.063120
10 ##  6    1 3.91 0.057510
11 ##  7    1 4.27 0.052400
12 ##  8    2 4.72 0.047750
13 ##  9    2 5.26 0.043510
14 ## 10    2 5.70 0.039640
15 ## 11    2 6.07 0.036120
16 ## 12    3 6.40 0.032910
17 ## 13    3 6.75 0.029990
18 ## 14    3 7.04 0.027320
19 ## 15    3 7.28 0.024900
20 ## 16    3 7.48 0.022680
21 ## 17    4 7.71 0.020670
22 ## 18    4 7.94 0.018830
23 ## 19    5 8.24 0.017160
24 ## 20    5 8.52 0.015640
25 ## 21    5 8.75 0.014250
26 ## 22    5 8.95 0.012980
27 ## 23    5 9.11 0.011830
28 ## 24    5 9.25 0.010780
29 ## 25    5 9.36 0.009819
30 ## 26    5 9.46 0.008947
31 ## 27    5 9.54 0.008152
32 ## 28    5 9.61 0.007428
33 ## 29    5 9.67 0.006768
34 ## 30    5 9.71 0.006167
35 ## 31    5 9.75 0.005619
36 ## 32    5 9.79 0.005120
37 ## 33    5 9.81 0.004665
38 ## 34    5 9.84 0.004251
39 ## 35    5 9.86 0.003873
```

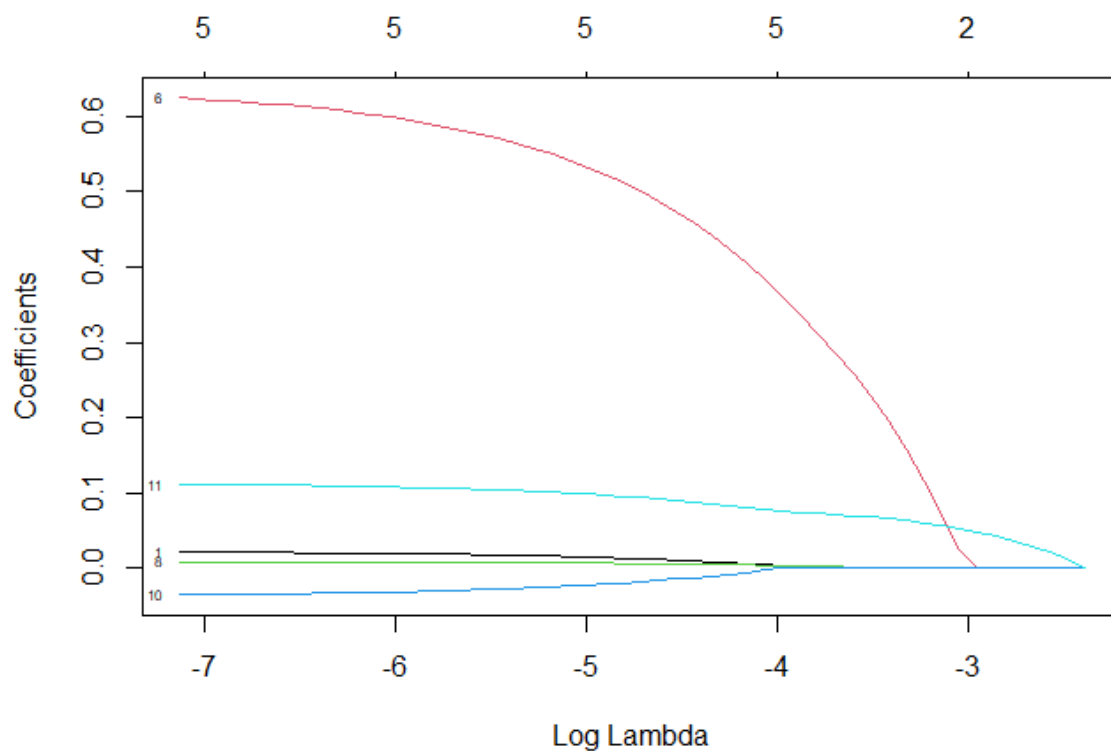
```

40 ## 36 5 9.87 0.003529
41 ## 37 5 9.89 0.003215
42 ## 38 5 9.90 0.002930
43 ## 39 5 9.91 0.002669
44 ## 40 5 9.92 0.002432
45 ## 41 5 9.92 0.002216
46 ## 42 5 9.93 0.002019
47 ## 43 5 9.93 0.001840
48 ## 44 5 9.94 0.001676
49 ## 45 5 9.94 0.001528
50 ## 46 5 9.94 0.001392
51 ## 47 5 9.94 0.001268
52 ## 48 5 9.95 0.001156
53 ## 49 5 9.95 0.001053
54 ## 50 5 9.95 0.000959
55 ## 51 5 9.95 0.000874
56 ## 52 5 9.95 0.000796

```

lambda = 0.000796 only 5 predictors were preserved

```
1 plot(model_lasso, xvar = "lambda", label = TRUE)
```



```

1 lasso.coef = coef(model_lasso, s = 0.000796)
2 lasso.coef

```

```

1 ## 12 x 1 sparse Matrix of class "dgCMatrix"
2 ##                                     s1
3 ## (Intercept)                      -4.65972880
4 ## age                               0.02173778
5 ## race                               .

```

```

6  ## marital_status      .
7  ## t_stage             .
8  ## n_stage             .
9  ## grade                0.62443173
10 ## a_stage             .
11 ## tumor_size           0.00894751
12 ## progesterone_status .
13 ## regional_node_examined -0.03496491
14 ## reginol_node_positive 0.11203866

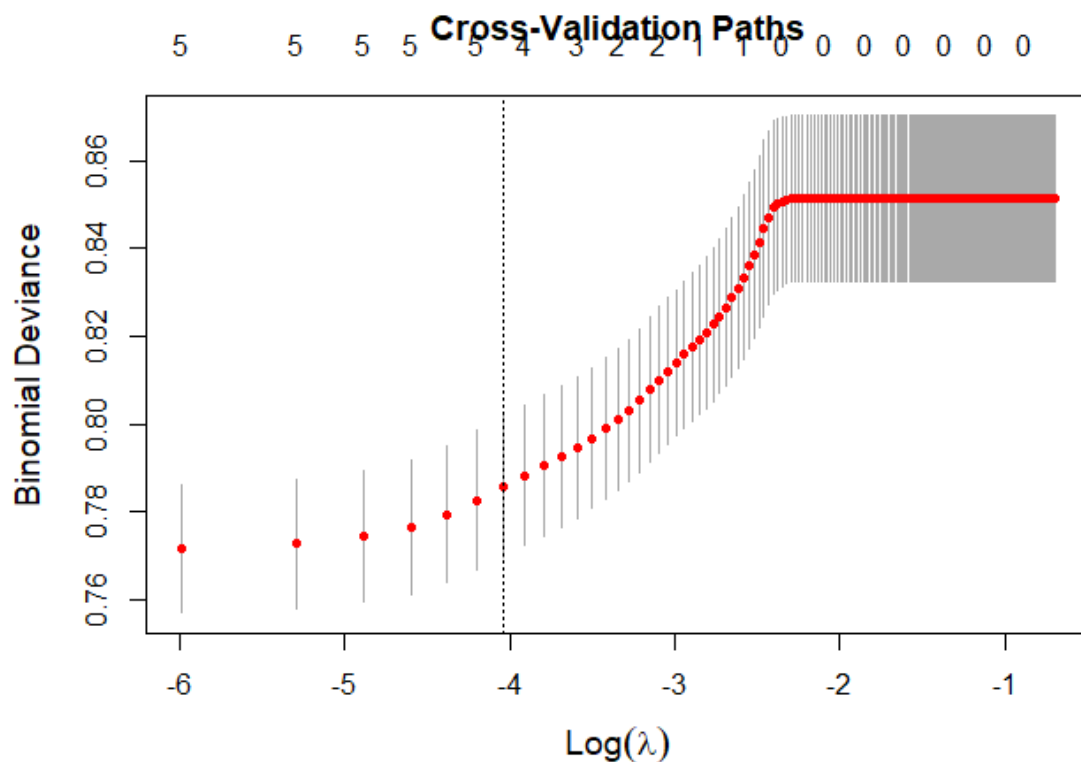
```

we can get the coefficients under this situation.
(lambda = 0.000796). Now, we want to find the best lambda.

```

1  set.seed(123)
2  lambdas = seq(0, 0.5, length.out = 200)
3  cv.lasso =
4    cv.glmnet(x,
5              y,
6              alpha = 1,
7              lambda = lambdas,
8              nfolds = 5,
9              family = "binomial")
10
11 par(mfrow = c(1, 1), mar = c(5, 5, 4, 2) + 0.1)
12 plot(cv.lasso,
13       main = "Cross-Validation Paths",
14       cex.lab = 1.2, cex.main = 1.2)

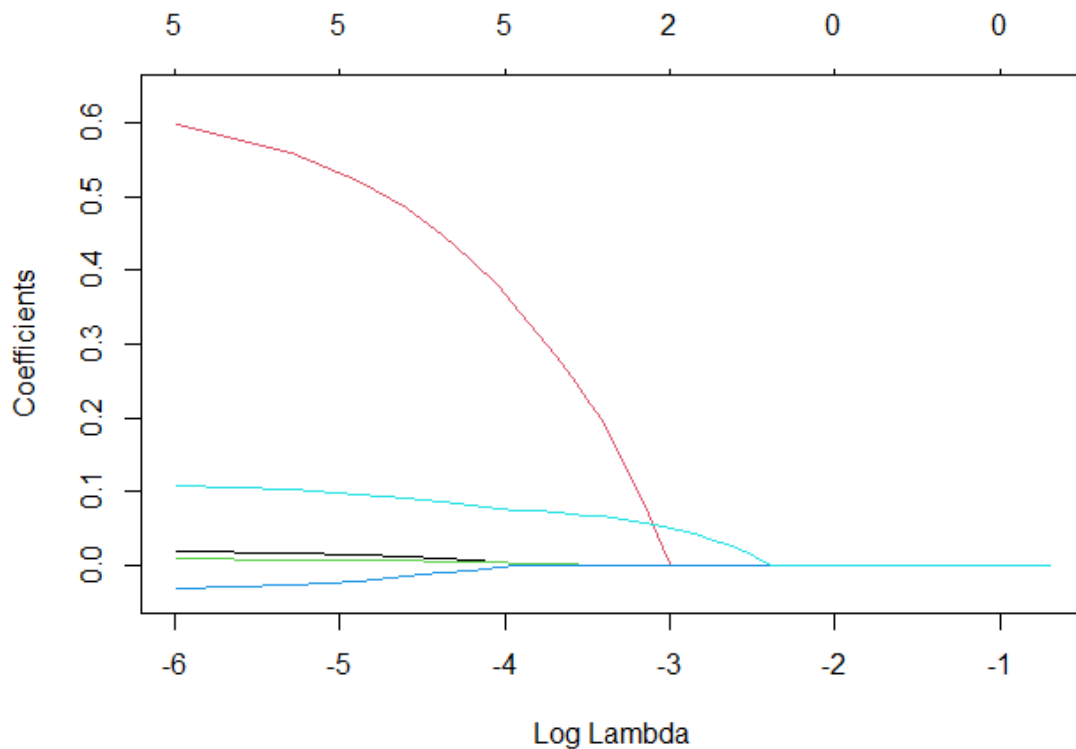
```



```

1  plot(cv.lasso$glmnet.fit,
2        xvar = "lambda",
3        label = T)

```



```
1 lasso_1se =
2   cv.lasso$lambda.1se
3 lasso_1se
```

```
1 ## [1] 0.01758794
```

```
1 lasso.coef =
2   coef(cv.lasso$glmnet.fit,
3       s = lasso_1se,
4       exact = F)
5 lasso.coef
```

```
1 ## 12 x 1 sparse Matrix of class "dgCMatrix"
2 ##                                     s1
3 ## (Intercept)                      -3.288096024
4 ## age                               0.004409007
5 ## race                             .
6 ## marital_status                    .
7 ## t_stage                           .
8 ## n_stage                           .
9 ## grade                             0.377930648
10 ## a_stage                           .
11 ## tumor_size                        0.004255330
12 ## progesterone_status               .
13 ## regional_node_examined           -0.001651846
14 ## reginol_node_positive             0.077769351
```

according to the outcome, the lasso process help me
select these variables: `age`, `grade`, `tumor_size`,
`regional_node_examined`, `reginol_node_positive`.

the model

```
1 formula3 =
2   as.formula(
3     status ~ age + grade + tumor_size + regional_node_examined +
4     reginol_node_positive)
5 model = glm(formula3,
6             data = heart_df_log,
7             family = binomial(link = 'logit'))
8
9 summary(model)
```

```
1 ##
2 ## Call:
3 ## glm(formula = formula3, family = binomial(link = "logit"), data =
4   heart_df_log)
5 ##
6 ## Coefficients:
7 ##              Estimate Std. Error z value Pr(>|z|)
8 ## (Intercept)    -2.816484   0.319834  -8.806 < 2e-16 ***
9 ## age              0.022642   0.005281   4.288 1.81e-05 ***
10 ## grade2         -0.646807   0.098600  -6.560 5.38e-11 ***
11 ## grade1         -1.257807   0.186371  -6.749 1.49e-11 ***
12 ## tumor_size      0.009180   0.002015   4.556 5.22e-06 ***
13 ## regional_node_examined -0.036780  0.007073  -5.200 1.99e-07 ***
14 ## reginol_node_positive  0.113925  0.008944  12.737 < 2e-16 ***
15 ## ---
16 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
17 ##
18 ## (Dispersion parameter for binomial family taken to be 1)
19 ##
20 ## Null deviance: 3407.5  on 4004  degrees of freedom
21 ## Residual deviance: 3068.3  on 3998  degrees of freedom
22 ## AIC: 3082.3
23 ##
24 ## Number of Fisher Scoring iterations: 5
```

```
1 check_collinearity(model)
```

```
1 ## # Check for Multicollinearity
2 ##
3 ## Low Correlation
4 ##
5 ##              Term  VIF   VIF 95% CI Increased SE Tolerance
6 ##              age  1.02 [1.00, 1.10]          1.01      0.98
7 ##              grade 1.02 [1.00, 1.10]          1.01      0.98
8 ##              tumor_size 1.05 [1.03, 1.10]          1.03      0.95
9 ## regional_node_examined 1.47 [1.42, 1.54]          1.21      0.68
10 ## reginol_node_positive 1.51 [1.45, 1.58]          1.23      0.66
11 ## Tolerance 95% CI
12 ##      [0.91, 1.00]
13 ##      [0.91, 1.00]
14 ##      [0.91, 0.97]
15 ##      [0.65, 0.71]
```

```
16 | ## [0.63, 0.69]
```

We find our third model and check its collinearity, the model's performance is good.

Diagnosis

```
1 | fit1 = lrm(formula1, data = heart_df_log, x = T, y = T)
2 | fit2 = lrm(formula2, data = heart_df_log, x = T, y = T)
3 | fit3 = lrm(formula3, data = heart_df_log, x = T, y = T)
```

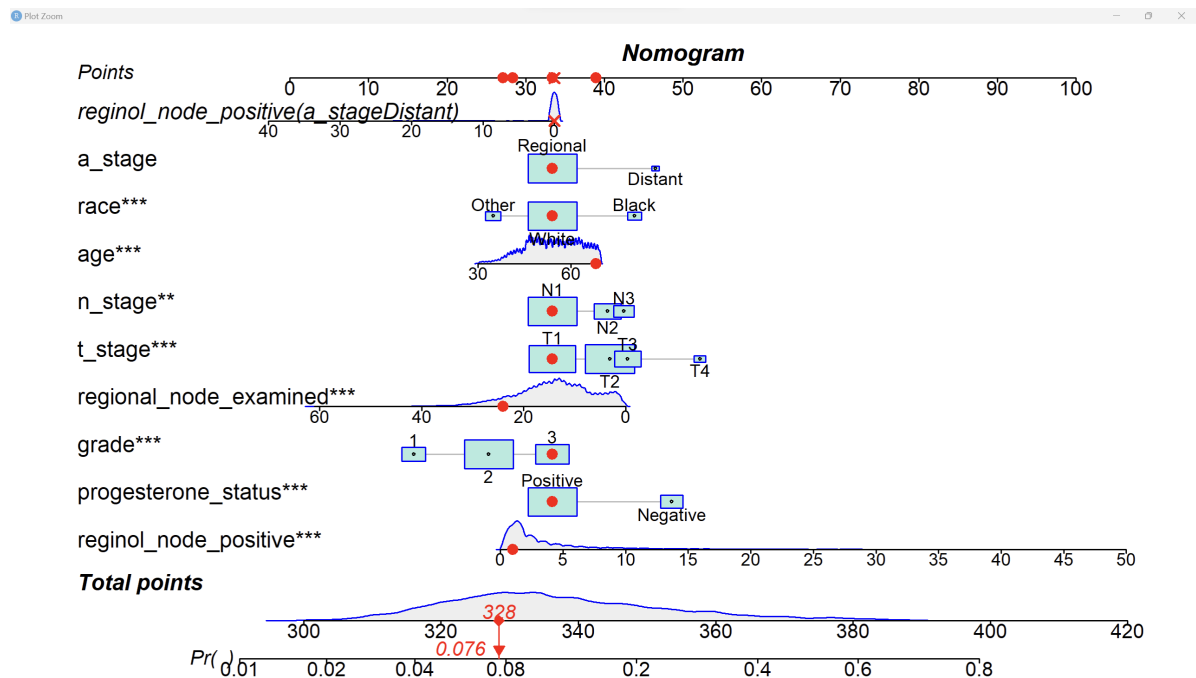
Nomogram

```
1 | dd = datadist(heart_df_log)
2 | options(datadist = "dd")
3 |
4 | nomogram_plot = function(input, number){
5 |   nomo_plot =
6 |     regplot(input,
7 |             observation = heart_df_log[number,],
8 |             center = TRUE,
9 |             title = "Nomogram",
10 |            points = TRUE,
11 |            odds = FALSE,
12 |            showP = TRUE,
13 |            rank= "sd",
14 |            clickable = FALSE)
15 | }
```

```
1 | nomogram_plot(fit1, 1)
```

```
1 | ## Regression input lrm formula:
2 |
3 | ## status ~` age + race + t_stage + n_stage + grade + a_stage +
  progesterone_status + regional_node_examined + reginol_node_positive +
  a_stage:reginol_node_positive
```

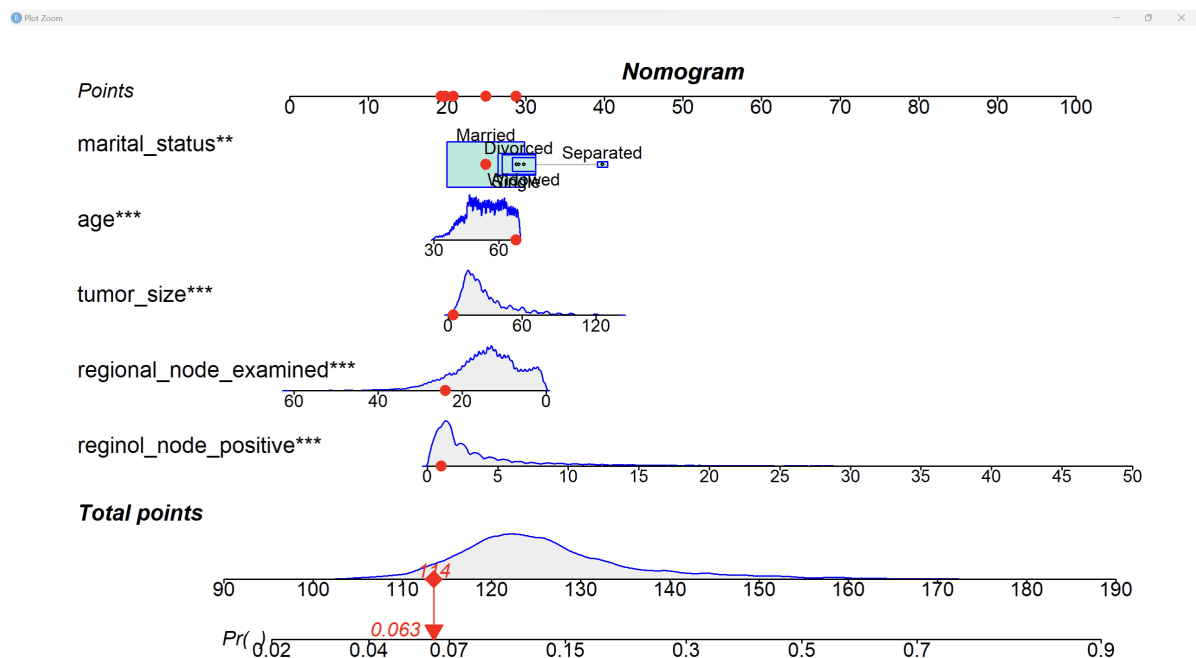
```
1 | knitr::include_graphics("nomogram_plot_1.png")
```



```
1 nomogram_plot(fit2, 1)
```

```
1 ## Regression input lrm formula:
2
3 ## status ~` age + regional_node_examined + tumor_size +
  reginol_node_positive + marital_status
```

```
1 knitr::include_graphics("nomogram_plot_2.png")
```



```
1 nomogram_plot(fit3, 1)
```



```

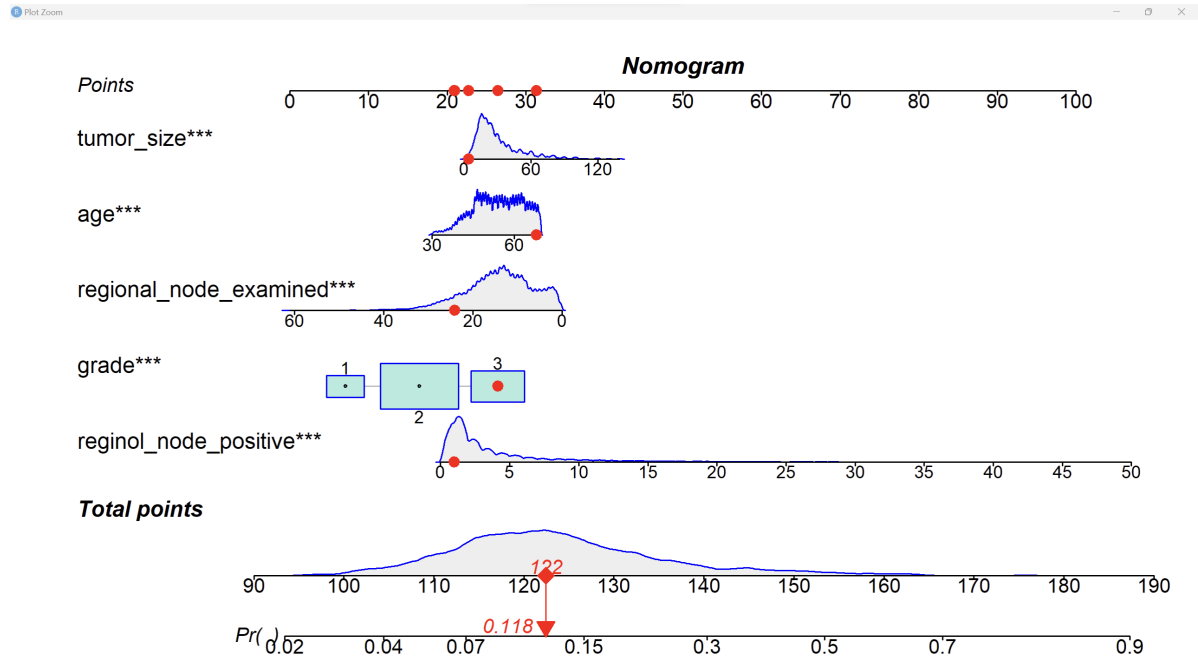
1 ## Regression input lrm formula:
2
3 ## status ~ age + grade + tumor_size + regional_node_examined +
  reginol_node_positive

```

```

1 knitr::include_graphics("nomogram_plot_3.png")

```



According to the monogram graph, we can compare three model by random choosing am observation from our dataset. In my code, I choose the first observation. We find that the probability to dead calculated in three models are 7.6%, 6.3% and 11.8%, the truth is that this person is alive, so the model 2 has the best simulation effect.

Calibration curve

```

1 cali_plot = function(input){
2   cal = calibrate(input, method = "boot", B = 500)
3
4   plot(cal,
5     xlim = c(0,1),
6     xlab = "Predicted Probability",
7     ylab = "Observed Probability",
8     legend = FALSE,
9     subtitles = TRUE)
10  abline(0,1,col = "black",
11    lty = 2,
12    lwd = 2)
13  lines(cal[,c("predy", "calibrated.orig")],
14    type = "l",
15    lwd = 2,
16    col = "red",
17    pch = 16)
18  lines(cal[,c("predy", "calibrated.corrected")],
19    type = "l",

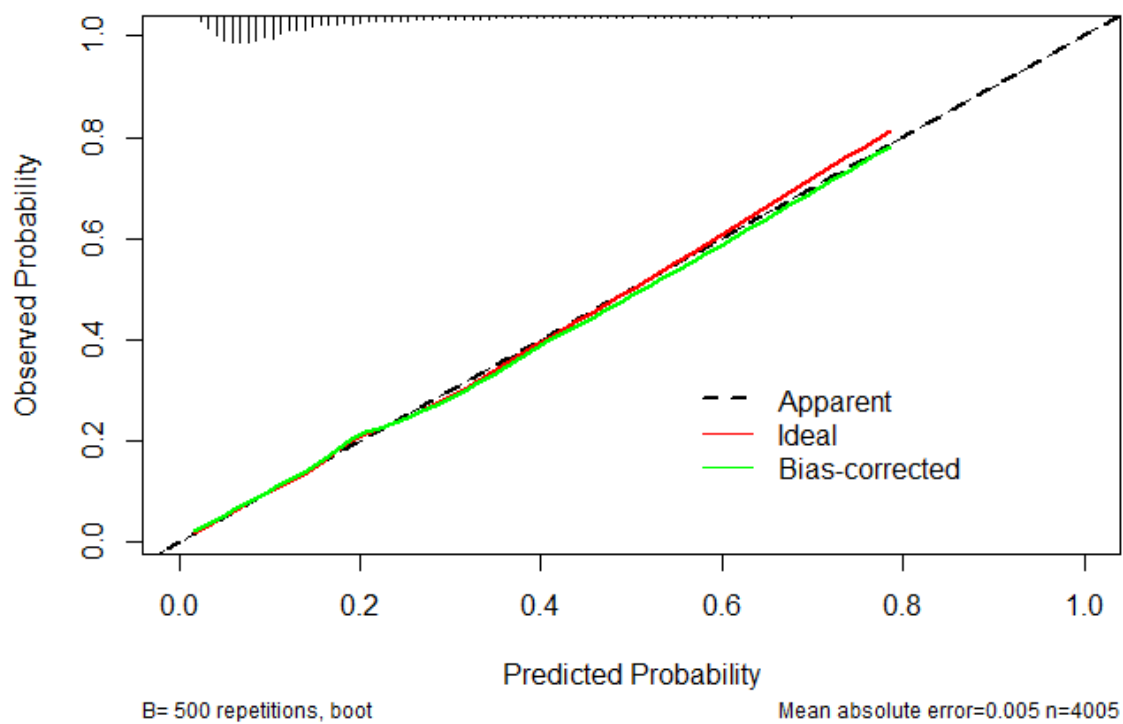
```

```

20     lwd = 2,
21     col = "green",
22     pch = 16)
23   legend(0.55, 0.35,
24         c("Apparent", "Ideal", "Bias-corrected"),
25         lty = c(2, 1, 1),
26         lwd = c(2, 1, 1),
27         col = c("black", "red", "green"),
28         bty = "n")
29 }

```

```
1 cali_plot(fit1)
```

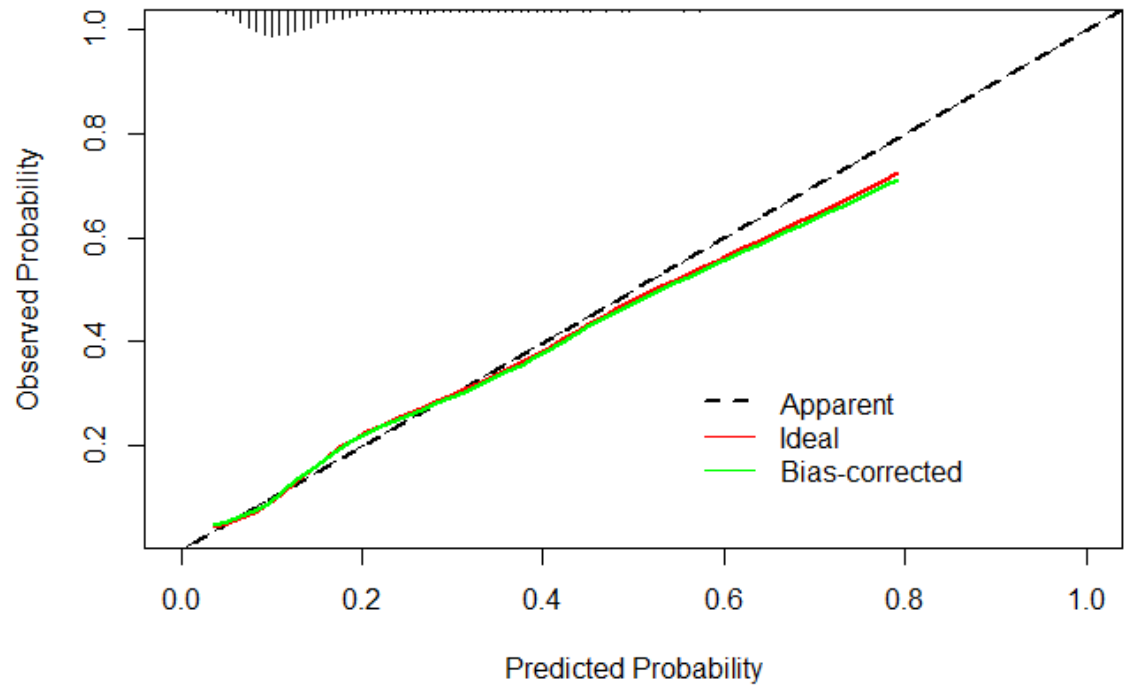


```

1 ##
2 ## n=4005   Mean absolute error=0.005   Mean squared error=4e-05
3 ## 0.9 quantile of absolute error=0.012

```

```
1 cali_plot(fit2)
```

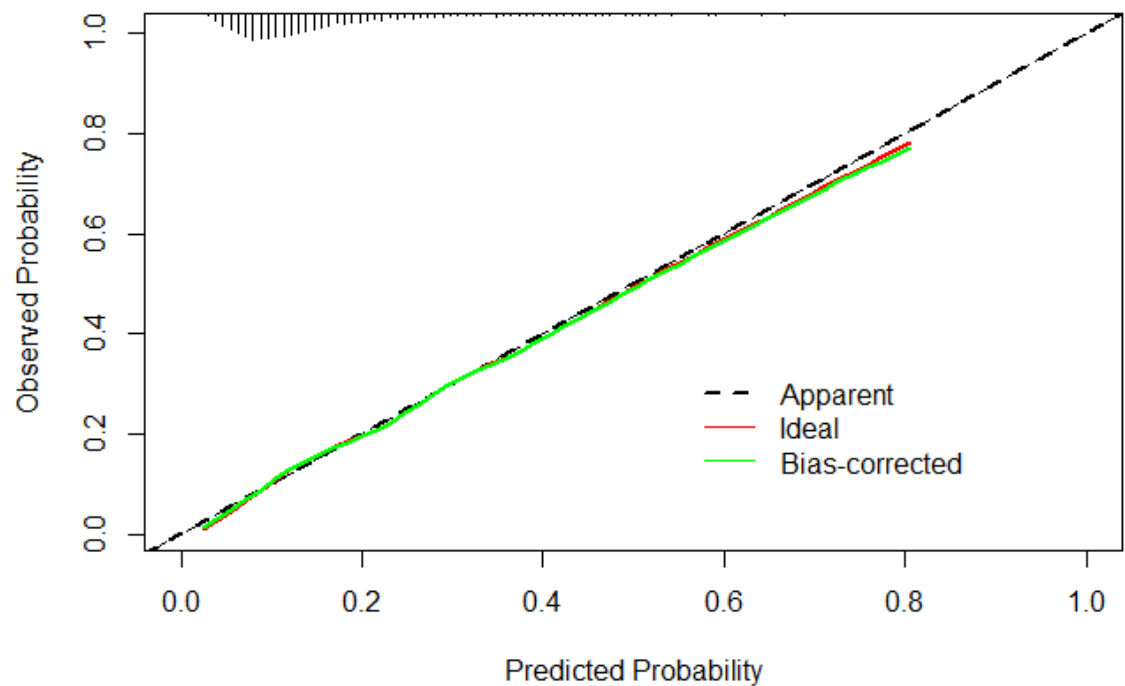


B= 500 repetitions, boot

Mean absolute error=0.009 n=4005

```
1 ##
2 ## n=4005 Mean absolute error=0.009 Mean squared error=0.00014
3 ## 0.9 quantile of absolute error=0.019
```

```
1 cali_plot(fit3)
```



B= 500 repetitions, boot

Mean absolute error=0.006 n=4005

```
1 ##
2 ## n=4005   Mean absolute error=0.006   Mean squared error=5e-05
3 ## 0.9 Quantile of absolute error=0.01
```

According to the calibration, all of three models have a good performance. Their calibration curves closely follow the diagonal line (45-degree line).

Hosmer and Lemeshow Goodness-of-Fit Test

```
1 fit1 = glm(formula1,
2           data = heart_df_log,
3           family = binomial(link = logit))
4 fit2 = glm(formula2,
5           data = heart_df_log,
6           family = binomial(link = logit))
7 fit3 = glm(formula3,
8           data = heart_df_log,
9           family = binomial(link = logit))
10
11 h11 = hoslem.test(fit1$y, fitted(fit1), g = 10)
12 h12 = hoslem.test(fit2$y, fitted(fit2), g = 10)
13 h13 = hoslem.test(fit3$y, fitted(fit3), g = 10)
14
15 h11$p.value
```

```
1 ## [1] 0.6940802
```

```
1 h12$p.value
```

```
1 ## [1] 0.2399465
```

```
1 h13$p.value
```

```
1 ## [1] 0.7272005
```

according to the outcome, we can find that all of three model's p-value are bigger than 5% which means they all have a good simulation performance. However, the model3 is the best and the model2 is the worst.

Validation

```
1 cro_validation = function(input, num){
2   train =
3   trainControl(method = "cv", number = num)
4
5   model_caret =
6   train(input,
7         data = heart_df_log,
8         trControl = train,
```

```
9         method = "glm",
10         na.action = na.pass)
11
12 model_caret$finalModel
13 print(model_caret)
14 }
```

```
1 set.seed(123)
2 result1 = cro_validation(formula1, 5)
```

```
1 ## Generalized Linear Model
2 ##
3 ## 4005 samples
4 ##    9 predictor
5 ##    2 classes: 'Alive', 'Dead'
6 ##
7 ## No pre-processing
8 ## Resampling: Cross-Validated (5 fold)
9 ## Summary of sample sizes: 3204, 3205, 3203, 3205, 3203
10 ## Resampling results:
11 ##
12 ##    Accuracy    Kappa
13 ##    0.8566796  0.1874696
```

```
1 result2 = cro_validation(formula2, 5)
```

```
1 ## Generalized Linear Model
2 ##
3 ## 4005 samples
4 ##    5 predictor
5 ##    2 classes: 'Alive', 'Dead'
6 ##
7 ## No pre-processing
8 ## Resampling: Cross-Validated (5 fold)
9 ## Summary of sample sizes: 3203, 3204, 3205, 3204, 3204
10 ## Resampling results:
11 ##
12 ##    Accuracy    Kappa
13 ##    0.8524347  0.1152586
```

```
1 result3 = cro_validation(formula3, 5)
```

```

1 ## Generalized Linear Model
2 ##
3 ## 4005 samples
4 ##    5 predictor
5 ##    2 classes: 'Alive', 'Dead'
6 ##
7 ## No pre-processing
8 ## Resampling: Cross-Validated (5 fold)
9 ## Summary of sample sizes: 3205, 3203, 3205, 3203, 3204
10 ## Resampling results:
11 ##
12 ##    Accuracy    Kappa
13 ##    0.8521874  0.1191757

```

```

1 result_tb =
2   tibble(model = c("model1", "model2", "model3"),
3           cro_val_accuracy = c(result1[1,1], result2[1,1], result3[1,1]),
4           cro_val_kappa = c(result1[1,2], result2[1,2], result3[1,2]))
5
6 knitr::kable(result_tb)

```

model	cro_val_accuracy	cro_val_kappa
model1	0.8566796	0.1874696
model2	0.8524347	0.1152586
model3	0.8521874	0.1191757

To test the probability to predict a new observation, I use the method cross-validation with 5 folds. We can find no matter from a accuracy perspective or a kappa perspective, the model 1 is the best among these three models.

Conclusion

```

1 result_tb =
2   result_tb |>
3   mutate(
4     nomo_plot = c("7.6%", "6.3%", "11.8%"),
5     hl_pvalue = c(hl1$p.value, hl2$p.value, hl3$p.value)
6   )
7
8 knitr::kable(result_tb)

```

model	cro_val_accuracy	cro_val_kappa	nomo_plot	hl_pvalue
model1	0.8566796	0.1874696	7.6%	0.6940802
model2	0.8524347	0.1152586	6.3%	0.2399465
model3	0.8521874	0.1191757	11.8%	0.7272005

According to the outcome form, I will choose the model1 as my final model. Though it has a worse performance in monogram plot procedure, It performs perfectly in other procedures.

Interpretation

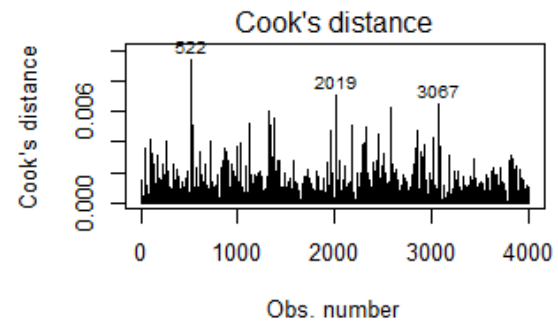
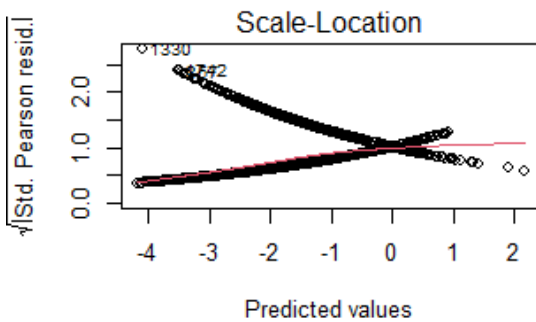
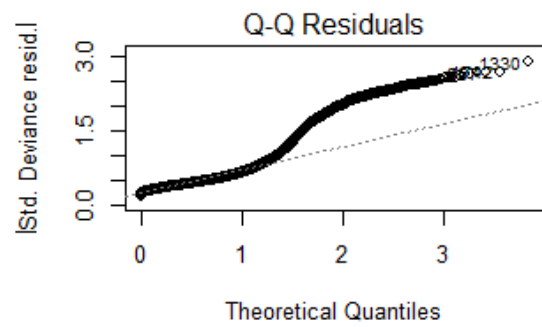
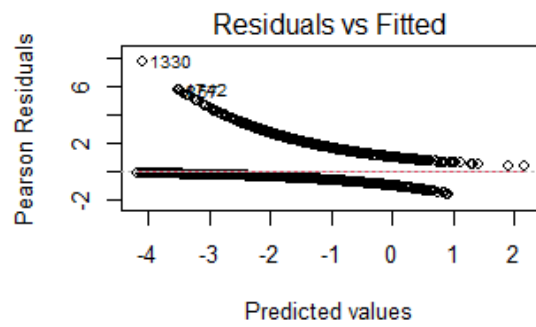
```
1 | coef_model = coef(fit1)
```

$\log\left(\frac{p}{1-p}\right) = -3.205321 + 0.02186205 \cdot \text{age} + 0.5803038 \cdot \text{f}(raceBlack) + 0.4375129 \cdot \text{f}(raceOther) + 0.0000004 \cdot \text{f}(a_stageT2) + 0.522447 \cdot \text{f}(a_stageT3) + 0.0200008 \cdot \text{f}(a_stageT4) + 0.3880029 \cdot \text{f}(a_stageT5) + 0.330344 \cdot \text{f}(a_stageT6) + 0.0401202 \cdot \text{f}(regnode0) + 0.0462770 \cdot \text{f}(regnode1) + 0.7207728 \cdot \text{f}(a_stageDistant) + 0.0483142 \cdot \text{f}(regnodeareaNegative) + 0.0000002 \cdot \text{regnode_adj_continued} + 0.0000002 \cdot \text{regnode_adj_positive} + 0.0000003 \cdot \text{f}(a_stageDistant) + \text{regnode_adj_positive}$

1, continuous: controlling other conditions unchanged, increasing by one year in age is associated with a 0.02186205-fold increase in the risk of death. 2, categorical variable: All other conditions being constant, the probability of death for Black individuals is 0.58 times higher than that for White individuals. 3, interaction: if two observations are each in a distant a_stage and a regional a_stage, then one unit increase in regional node positive will make a distant one 0.6762288 more times probability to die.

Check model's effect with data stratified by race

```
1 | par(mfrow = c(2, 2), mar = c(5, 4, 4, 2) + 0.1)
2 | # residual vs fitted plot
3 | plot(fit1, which = 1)
4 |
5 | # QQ plot
6 | plot(fit1, which = 2)
7 |
8 | plot(fit1, which = 3)
9 | plot(fit1, which = 4)
```



```

1 heart_df_white =
2   heart_df_log |>
3   filter(race == "white")
4
5 heart_df_other =
6   heart_df_log |>
7   filter(!race == "white")

```

```

1 cal_predict_value = function(input, data){
2
3   predicted_values =
4     predict(input,
5             newdata = data,
6             type = "response")
7
8   predicted_classes =
9     ifelse(predicted_values > 0.5, "Dead", "Alive")
10
11
12   accuracy =
13     sum(predicted_classes == data$status) / length(data$status)
14
15   return(accuracy)
16 }

```

```

1 cal_predict_value(fit1, heart_df_white)

```

```

1 ## [1] 0.8622314

```



```
1 | cal_predict_value(fit1, heart_df_other)
```

```
1 | ## [1] 0.8355263
```

According to the outcome here, we can find that the model's prediction accuracy is 0.836 to the white people and 0.836 to the other race people. There has a difference between two groups, but the degree is small which is just 3% and we can tolerate it. If we want to make a more accurate modification towards this shortcoming, we can just separate our original dataset and get two different model in different race condition.