

P8106 Midterm Project

Leyang Rui, Jinghan Zhao

2025-03-28

Load Data

```
load("data/dat1.RData")
train_data = dat1 |>
  janitor::clean_names() |>
  mutate(
    gender = as.factor(gender),
    diabetes = as.factor(diabetes),
    hypertension = as.factor(hypertension),
    race = fct_recode(race,
      White = "1",
      Asian = "2",
      Black = "3",
      Hispanic = "4"),
    gender = fct_recode(gender,
      Male = "1",
      Female = "0"),
    smoking = fct_recode(smoking,
      "Never smoked" = "0",
      "Former smoker" = "1",
      "Current smoker" = "2"))

load("data/dat2.RData")
test_data = dat2 |>
  janitor::clean_names() |>
  mutate(
    gender = as.factor(gender),
    diabetes = as.factor(diabetes),
    hypertension = as.factor(hypertension),
    race = fct_recode(race,
      White = "1",
      Asian = "2",
      Black = "3",
      Hispanic = "4"),
    gender = fct_recode(gender,
      Male = "1",
      Female = "0"),
    smoking = fct_recode(smoking,
      "Never smoked" = "0",
      "Former smoker" = "1",
```

```
)
    "Current smoker" = "2")
```

Modify Data

```
train_data1 =
  train_data %>%
  select(-id, -height, -weight, -hypertension)

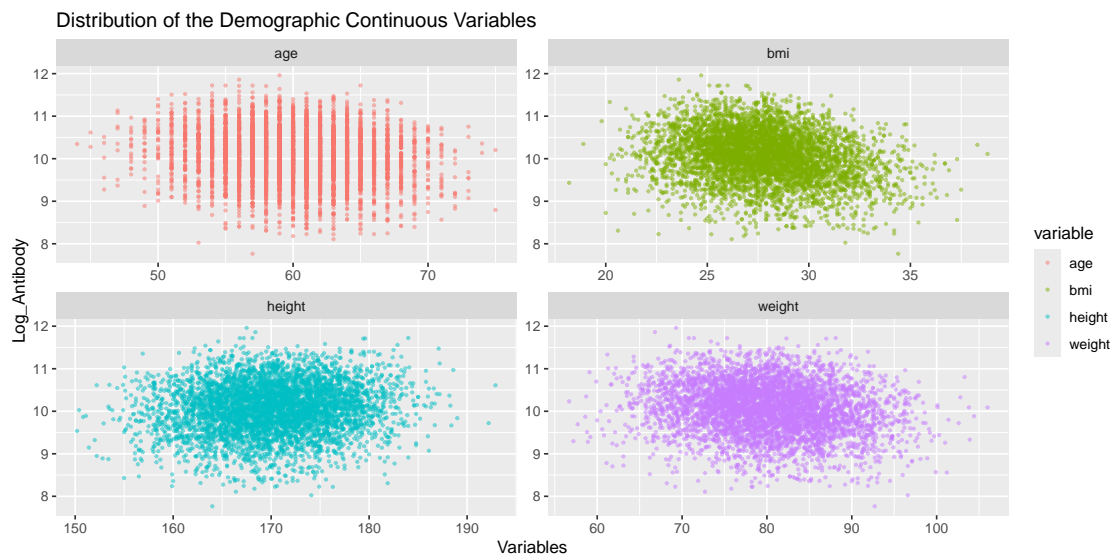
x_train = model.matrix(log_antibody ~ ., train_data1)[, -1]
colnames(x_train) = make.names(colnames(x_train), unique = TRUE)

y_train = train_data1[, "log_antibody"]
ctrl1 = trainControl(method = "cv", number = 10)
```

Descriptive Analysis

Numeric Variables

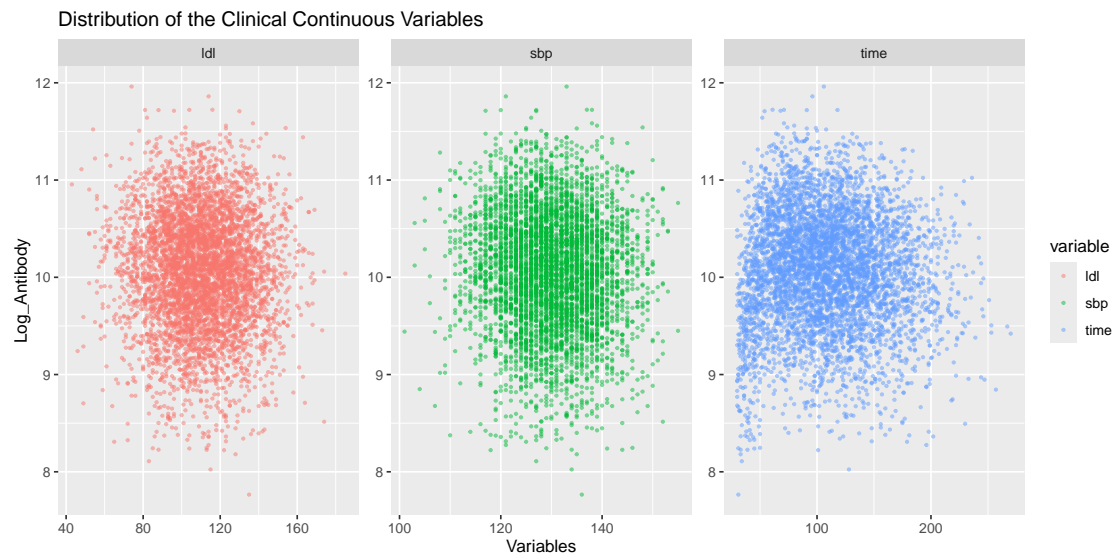
```
train_data |>
  pivot_longer(
    cols = c(age, height, weight, bmi),
    names_to = "variable",
    values_to = "value"
  ) |>
  ggplot(aes(x = value, y = log_antibody, color = variable)) +
  geom_point(alpha = 0.5, size = 0.6) +
  facet_wrap(variable ~ ., scales = "free") +
  labs(title = "Distribution of the Demographic Continuous Variables",
       x = "Variables",
       y = "Log_Antibody")
```



```

train_data |>
  pivot_longer(
    cols = c(sbp, ldl, time),
    names_to = "variable",
    values_to = "value"
  ) |>
  ggplot(aes(x = value, y = log_antibody, color = variable)) +
  geom_point(alpha = 0.5, size = 0.6) +
  facet_wrap(variable ~ ., scales = "free") +
  labs(title = "Distribution of the Clinical Continuous Variables",
       x = "Variables",
       y = "Log_Antibody")

```



```

train_data %>%
  pivot_longer(
    cols = c(age, height, weight, bmi, sbp, ldl, time, log_antibody),
    names_to = "variable_name",
    values_to = "value"
  ) %>%
  group_by(variable_name) %>%
  summarize(
    mean = mean(value),
    median = median(value),
    min = min(value),
    first_quantile = quantile(value, probs = 0.25),
    third_quantile = quantile(value, probs = 0.75),
    max = max(value)
  ) %>%
  ungroup() %>%
  arrange(desc(variable_name == "log_antibody"), variable_name) %>%
  knitr::kable(digits = 3, caption = "Descriptive Statistics")

```

Table 1: Descriptive Statistics

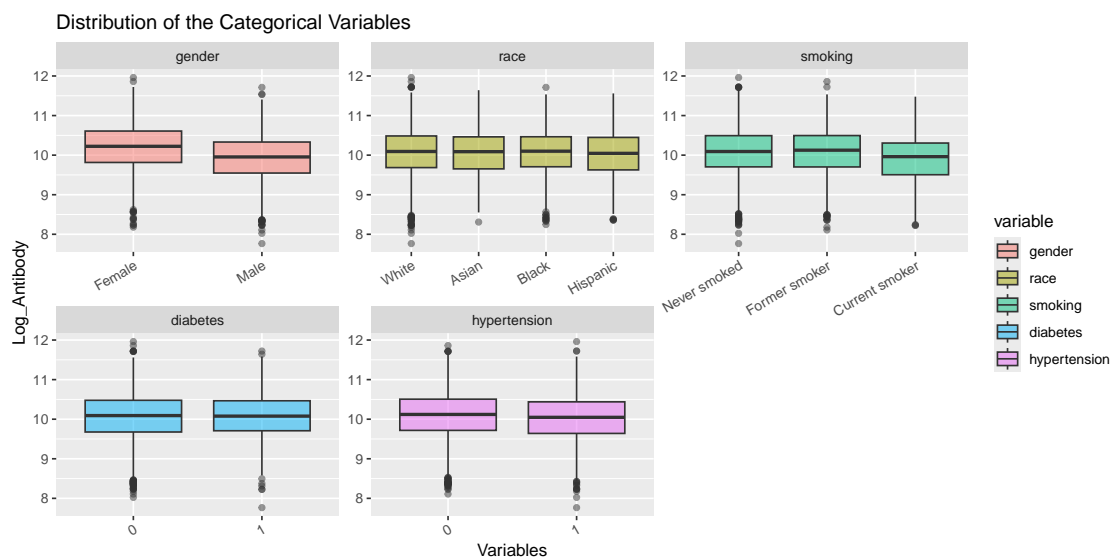
variable_name	mean	median	min	first_quantile	third_quantile	max
log_antibody	10.064	10.089	7.765	9.682	10.478	11.961
age	59.968	60.000	44.000	57.000	63.000	75.000
bmi	27.740	27.600	18.200	25.800	29.500	38.800
height	170.126	170.100	150.200	166.100	174.225	192.900
ldl	109.909	110.000	43.000	96.000	124.000	185.000
sbp	129.900	130.000	101.000	124.000	135.000	155.000
time	108.863	106.000	30.000	76.000	138.000	270.000
weight	80.109	80.100	56.700	75.400	84.900	106.000

Categorical Variables

```

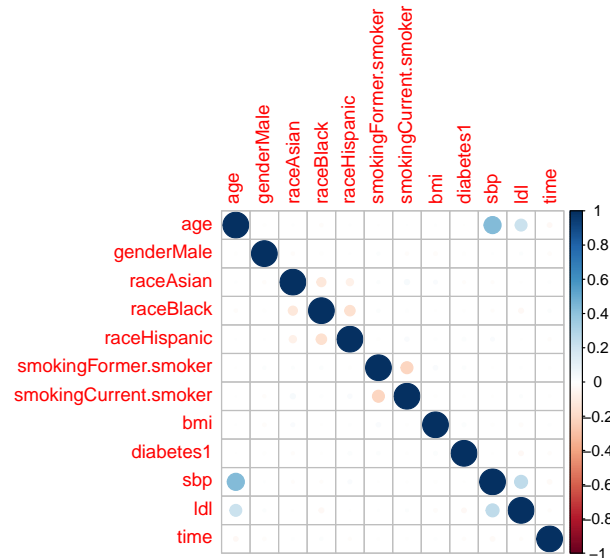
train_data |>
  pivot_longer(
    cols = c(gender, race, smoking, diabetes, hypertension),
    names_to = "variable",
    values_to = "value"
  ) |>
  mutate(
    variable = factor(variable, levels = c("gender", "race", "smoking", "diabetes", "hypertension"))
  ) |>
  ggplot(aes(x = value, y = log_antibody, fill = variable)) +
  geom_boxplot(alpha = 0.5) +
  facet_wrap(variable ~ ., scales = "free") +
  labs(title = "Distribution of the Categorical Variables",
       x = "Variables",
       y = "Log_Antibody") +
  theme(axis.text.x = element_text(angle = 30, vjust = 1, hjust = 1))

```



Correlation Plot

```
corrplot(cor(x_train), method = "circle", type = "full")
```



Regression

Elastic Net

PCR and PLS

GAM and MARS

```
set.seed(37)
```

```
gam.fit = train(x_train, y_train,  
               method = "gam",  
               trControl = ctrl1)
```

```
gam.fit$bestTune
```

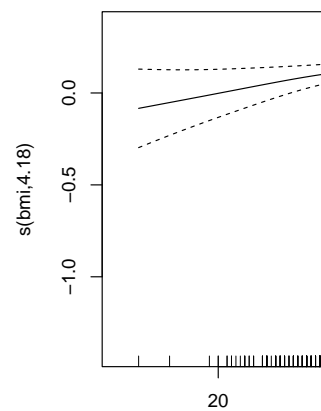
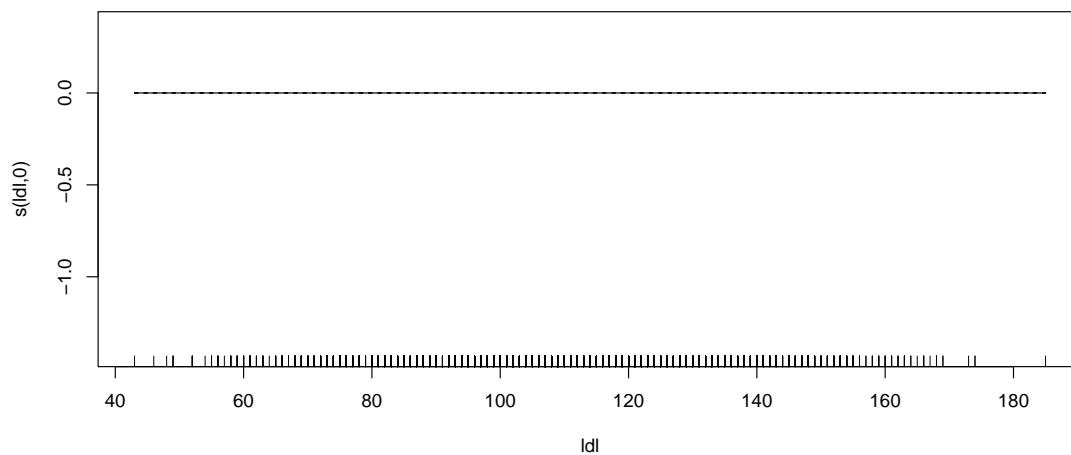
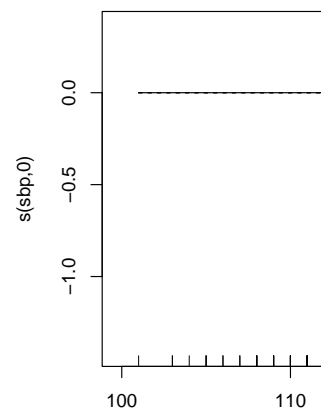
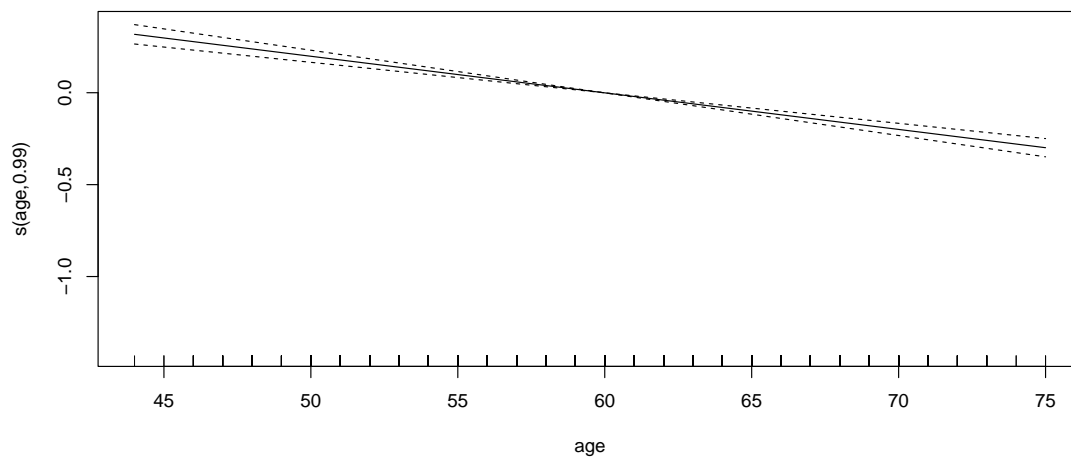
```
## select method  
## 2 TRUE GCV.Cp
```

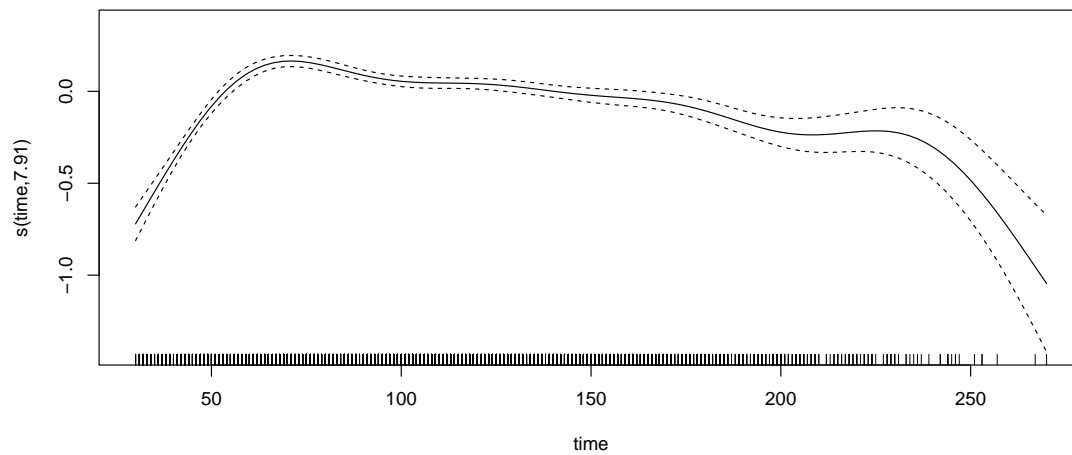
```
gam.fit$finalModel
```

```
##  
## Family: gaussian  
## Link function: identity  
##  
## Formula:
```

```
## .outcome ~ genderMale + raceAsian + raceBlack + raceHispanic +
##      smokingFormer.smoker + smokingCurrent.smoker + diabetes1 +
##      s(age) + s(sbp) + s(ldl) + s(bmi) + s(time)
##
## Estimated degrees of freedom:
## 0.992 0.000 0.000 4.179 7.915 total = 21.09
##
## GCV score: 0.2786375
```

```
plot(gam.fit$finalModel)
```



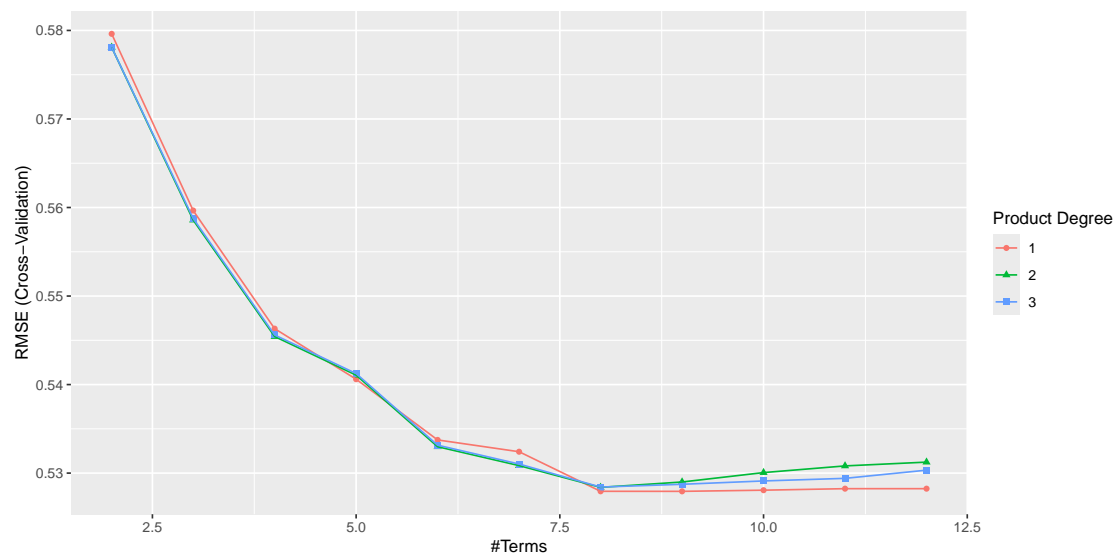


```
set.seed(37)

mars_grid = expand.grid(degree = 1:3,
                        nprune = 2:12)

mars.fit = train(x_train, y_train,
                 method = "earth",
                 tuneGrid = mars_grid,
                 trControl = ctrl1)

ggplot(mars.fit)
```



```
mars.fit$bestTune
```

```
##  nprune degree
##  8      9      1
```

```
coef(mars.fit$finalModel)
```

```
##          (Intercept)          h(27.8-bmi)          h(time-57)
##      10.847446930      -0.061997354      -0.002254182
##          h(57-time)          genderMale          h(age-59)
##      -0.033529326      -0.296290451      -0.022957648
##          h(59-age) smokingCurrent.smoker          h(bmi-23.7)
##      0.016138468      -0.205126851      -0.084380175
```

Regression Trees

```
set.seed(37)
```

```
tree_full = rpart(formula = log_antibody ~ ., data = train_data1,
                  control = rpart.control(cp = 0))
```

```
printcp(tree_full)
```

```
##
## Regression tree:
## rpart(formula = log_antibody ~ ., data = train_data1, control = rpart.control(cp = 0))
##
## Variables actually used in tree construction:
## [1] age      bmi      diabetes gender  ldl      race      sbp      smoking
## [9] time
##
## Root node error: 1778.6/5000 = 0.35572
##
## n= 5000
##
##      CP nsplit rel error  xerror    xstd
## 1  5.7918e-02      0  1.00000 1.00071 0.020042
## 2  2.7294e-02      1  0.94208 0.94320 0.018701
## 3  2.5172e-02      2  0.91479 0.92960 0.018467
## 4  2.0759e-02      3  0.88962 0.90421 0.018122
## 5  1.0560e-02      4  0.86886 0.88298 0.017677
## 6  7.9285e-03      5  0.85830 0.87618 0.017362
## 7  7.4831e-03      6  0.85037 0.87295 0.017328
## 8  6.8304e-03      7  0.84288 0.86572 0.017165
## 9  6.4129e-03      8  0.83605 0.86049 0.017049
## 10 6.4078e-03      9  0.82964 0.85726 0.016973
## 11 4.8526e-03     10  0.82323 0.84854 0.016789
## 12 3.2844e-03     11  0.81838 0.84659 0.016645
## 13 3.0338e-03     12  0.81510 0.84652 0.016679
## 14 2.9148e-03     13  0.81206 0.84669 0.016696
## 15 2.7907e-03     14  0.80915 0.84521 0.016657
## 16 2.7301e-03     15  0.80636 0.84505 0.016617
## 17 2.5602e-03     16  0.80363 0.84769 0.016691
## 18 2.3277e-03     17  0.80107 0.84831 0.016688
## 19 2.2239e-03     19  0.79641 0.84516 0.016615
## 20 2.1642e-03     21  0.79196 0.84473 0.016602
```


## 21	1.9279e-03	22	0.78980	0.84972	0.016801
## 22	1.7205e-03	23	0.78787	0.85048	0.016794
## 23	1.7191e-03	24	0.78615	0.85105	0.016785
## 24	1.7002e-03	26	0.78271	0.85150	0.016795
## 25	1.6647e-03	27	0.78101	0.85138	0.016788
## 26	1.6474e-03	28	0.77935	0.85082	0.016793
## 27	1.5736e-03	29	0.77770	0.84968	0.016751
## 28	1.4872e-03	30	0.77613	0.85101	0.016769
## 29	1.4587e-03	31	0.77464	0.85494	0.016959
## 30	1.3888e-03	32	0.77318	0.85549	0.016995
## 31	1.3702e-03	33	0.77179	0.85405	0.016944
## 32	1.3461e-03	34	0.77042	0.85372	0.016926
## 33	1.3077e-03	35	0.76908	0.85505	0.016914
## 34	1.2618e-03	36	0.76777	0.85756	0.016935
## 35	1.2616e-03	37	0.76651	0.86069	0.017033
## 36	1.2499e-03	40	0.76272	0.86146	0.017057
## 37	1.1923e-03	42	0.76022	0.86295	0.017042
## 38	1.1919e-03	43	0.75903	0.86387	0.017063
## 39	1.1819e-03	44	0.75784	0.86389	0.017062
## 40	1.1004e-03	46	0.75547	0.86856	0.017171
## 41	1.0945e-03	48	0.75327	0.87656	0.017327
## 42	1.0868e-03	55	0.74561	0.87721	0.017342
## 43	1.0720e-03	58	0.74235	0.87721	0.017364
## 44	1.0609e-03	59	0.74128	0.87786	0.017389
## 45	1.0586e-03	60	0.74022	0.87757	0.017392
## 46	1.0470e-03	61	0.73916	0.87900	0.017416
## 47	1.0469e-03	62	0.73811	0.87965	0.017412
## 48	1.0316e-03	64	0.73602	0.87792	0.017384
## 49	1.0221e-03	65	0.73499	0.88178	0.017468
## 50	1.0175e-03	70	0.72988	0.88240	0.017491
## 51	1.0085e-03	71	0.72886	0.88326	0.017571
## 52	9.9742e-04	72	0.72785	0.88569	0.017597
## 53	9.8051e-04	74	0.72586	0.88930	0.017684
## 54	9.6681e-04	75	0.72488	0.88971	0.017660
## 55	9.5372e-04	77	0.72294	0.89452	0.017809
## 56	9.5138e-04	78	0.72199	0.89620	0.017813
## 57	9.4890e-04	79	0.72104	0.89696	0.017833
## 58	9.4646e-04	80	0.72009	0.89913	0.017845
## 59	9.4501e-04	88	0.71238	0.89967	0.017867
## 60	9.2555e-04	91	0.70954	0.90291	0.017908
## 61	9.2335e-04	92	0.70862	0.90506	0.017938
## 62	9.2150e-04	93	0.70769	0.90437	0.017937
## 63	9.0635e-04	94	0.70677	0.90540	0.017963
## 64	9.0262e-04	95	0.70587	0.90601	0.017943
## 65	8.9875e-04	96	0.70496	0.90503	0.017918
## 66	8.9679e-04	97	0.70406	0.90480	0.017882
## 67	8.8749e-04	101	0.70048	0.90534	0.017892
## 68	8.7173e-04	103	0.69870	0.90590	0.017932
## 69	8.7111e-04	106	0.69609	0.90717	0.017949
## 70	8.6830e-04	107	0.69522	0.90724	0.017950
## 71	8.5924e-04	110	0.69261	0.90873	0.017987
## 72	8.5500e-04	117	0.68623	0.91160	0.018022
## 73	8.4650e-04	119	0.68452	0.91338	0.018107
## 74	8.3561e-04	120	0.68367	0.91450	0.018179

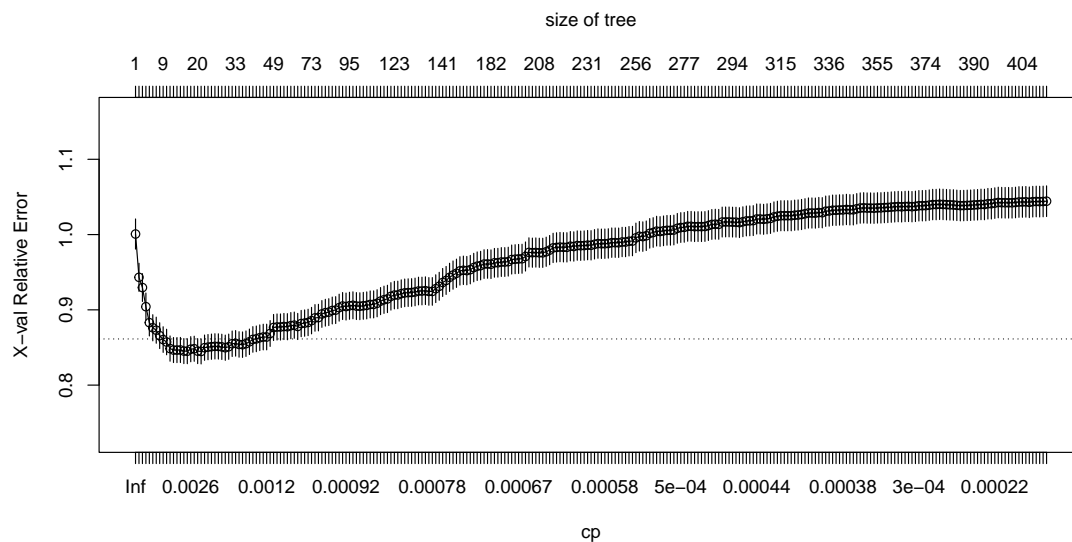
## 75	8.2232e-04	121	0.68284	0.91791	0.018240
## 76	8.2068e-04	122	0.68201	0.91957	0.018281
## 77	8.1855e-04	123	0.68119	0.91980	0.018283
## 78	8.0718e-04	125	0.67956	0.92151	0.018337
## 79	8.0119e-04	126	0.67875	0.92255	0.018378
## 80	7.9965e-04	127	0.67795	0.92287	0.018367
## 81	7.9120e-04	128	0.67715	0.92303	0.018369
## 82	7.8853e-04	129	0.67636	0.92377	0.018380
## 83	7.8826e-04	130	0.67557	0.92498	0.018401
## 84	7.8820e-04	131	0.67478	0.92498	0.018401
## 85	7.8228e-04	132	0.67399	0.92544	0.018405
## 86	7.8116e-04	133	0.67321	0.92447	0.018388
## 87	7.7954e-04	134	0.67243	0.92447	0.018388
## 88	7.7590e-04	138	0.66930	0.92811	0.018450
## 89	7.6555e-04	139	0.66852	0.93129	0.018568
## 90	7.6092e-04	140	0.66776	0.93636	0.018616
## 91	7.4559e-04	143	0.66548	0.93908	0.018600
## 92	7.4333e-04	147	0.66249	0.94314	0.018735
## 93	7.3053e-04	148	0.66175	0.94617	0.018778
## 94	7.2126e-04	152	0.65883	0.94858	0.018813
## 95	7.2064e-04	154	0.65739	0.95175	0.018872
## 96	7.1810e-04	158	0.65450	0.95253	0.018886
## 97	7.1557e-04	159	0.65378	0.95212	0.018870
## 98	7.1132e-04	160	0.65307	0.95340	0.018905
## 99	6.9696e-04	162	0.65165	0.95615	0.018958
## 100	6.9278e-04	164	0.65025	0.95796	0.018978
## 101	6.9053e-04	165	0.64956	0.95879	0.018976
## 102	6.8985e-04	172	0.64374	0.96080	0.019059
## 103	6.8634e-04	173	0.64305	0.96098	0.019044
## 104	6.8270e-04	181	0.63756	0.96060	0.019035
## 105	6.7836e-04	184	0.63538	0.96248	0.019103
## 106	6.7448e-04	187	0.63335	0.96258	0.019117
## 107	6.7441e-04	188	0.63267	0.96353	0.019157
## 108	6.7436e-04	189	0.63200	0.96353	0.019157
## 109	6.7214e-04	191	0.63065	0.96402	0.019172
## 110	6.6951e-04	193	0.62931	0.96682	0.019263
## 111	6.6879e-04	197	0.62663	0.96727	0.019289
## 112	6.6713e-04	199	0.62529	0.96763	0.019290
## 113	6.6364e-04	202	0.62329	0.96788	0.019293
## 114	6.4326e-04	203	0.62263	0.97087	0.019324
## 115	6.4242e-04	204	0.62198	0.97630	0.019393
## 116	6.4100e-04	205	0.62134	0.97594	0.019395
## 117	6.3949e-04	206	0.62070	0.97617	0.019397
## 118	6.3868e-04	207	0.62006	0.97580	0.019390
## 119	6.3755e-04	208	0.61942	0.97576	0.019387
## 120	6.2299e-04	209	0.61878	0.97718	0.019412
## 121	6.0990e-04	210	0.61816	0.97935	0.019445
## 122	6.0873e-04	216	0.61448	0.98228	0.019482
## 123	6.0510e-04	217	0.61387	0.98294	0.019500
## 124	6.0485e-04	219	0.61266	0.98338	0.019496
## 125	6.0210e-04	221	0.61145	0.98309	0.019502
## 126	5.9737e-04	223	0.61024	0.98303	0.019488
## 127	5.9728e-04	225	0.60905	0.98453	0.019499
## 128	5.9323e-04	226	0.60845	0.98417	0.019500

## 129	5.8916e-04	227	0.60786	0.98540	0.019505
## 130	5.8882e-04	228	0.60727	0.98518	0.019494
## 131	5.8663e-04	229	0.60668	0.98540	0.019495
## 132	5.8555e-04	230	0.60609	0.98597	0.019500
## 133	5.8536e-04	234	0.60375	0.98574	0.019512
## 134	5.7914e-04	235	0.60317	0.98708	0.019524
## 135	5.7768e-04	236	0.60259	0.98803	0.019548
## 136	5.7595e-04	237	0.60201	0.98779	0.019534
## 137	5.7527e-04	238	0.60143	0.98785	0.019535
## 138	5.7110e-04	239	0.60086	0.98820	0.019538
## 139	5.7084e-04	240	0.60029	0.98922	0.019543
## 140	5.7046e-04	242	0.59915	0.98922	0.019543
## 141	5.6757e-04	243	0.59858	0.98913	0.019530
## 142	5.6646e-04	244	0.59801	0.99014	0.019537
## 143	5.6561e-04	247	0.59631	0.99011	0.019533
## 144	5.6506e-04	248	0.59574	0.99114	0.019554
## 145	5.6502e-04	249	0.59518	0.99141	0.019554
## 146	5.4007e-04	255	0.59172	0.99561	0.019612
## 147	5.3640e-04	258	0.59010	0.99774	0.019605
## 148	5.3634e-04	262	0.58795	0.99729	0.019603
## 149	5.3391e-04	263	0.58742	0.99842	0.019619
## 150	5.2940e-04	264	0.58688	1.00161	0.019640
## 151	5.1708e-04	266	0.58582	1.00275	0.019702
## 152	5.1658e-04	267	0.58531	1.00456	0.019766
## 153	5.1652e-04	268	0.58479	1.00410	0.019740
## 154	5.1199e-04	270	0.58376	1.00496	0.019756
## 155	5.0908e-04	271	0.58324	1.00563	0.019757
## 156	5.0746e-04	272	0.58274	1.00563	0.019749
## 157	5.0600e-04	273	0.58223	1.00605	0.019747
## 158	5.0178e-04	274	0.58172	1.00829	0.019771
## 159	5.0060e-04	275	0.58122	1.00934	0.019800
## 160	4.9056e-04	276	0.58072	1.00960	0.019789
## 161	4.8924e-04	277	0.58023	1.01096	0.019815
## 162	4.8844e-04	278	0.57974	1.01096	0.019815
## 163	4.8612e-04	279	0.57925	1.01049	0.019813
## 164	4.8474e-04	280	0.57876	1.01062	0.019800
## 165	4.8429e-04	282	0.57780	1.01054	0.019799
## 166	4.8051e-04	283	0.57731	1.01072	0.019806
## 167	4.8046e-04	285	0.57635	1.01194	0.019810
## 168	4.7671e-04	286	0.57587	1.01345	0.019836
## 169	4.7617e-04	287	0.57539	1.01381	0.019837
## 170	4.7324e-04	288	0.57492	1.01360	0.019827
## 171	4.6544e-04	289	0.57444	1.01664	0.019849
## 172	4.6469e-04	291	0.57351	1.01677	0.019852
## 173	4.6438e-04	292	0.57305	1.01677	0.019852
## 174	4.6280e-04	293	0.57258	1.01656	0.019855
## 175	4.6091e-04	295	0.57166	1.01646	0.019859
## 176	4.5851e-04	296	0.57120	1.01586	0.019849
## 177	4.5040e-04	298	0.57028	1.01704	0.019863
## 178	4.4731e-04	299	0.56983	1.01807	0.019858
## 179	4.4672e-04	300	0.56938	1.01823	0.019868
## 180	4.4257e-04	302	0.56849	1.01857	0.019868
## 181	4.4176e-04	303	0.56805	1.02067	0.019966
## 182	4.3347e-04	305	0.56716	1.02068	0.019963

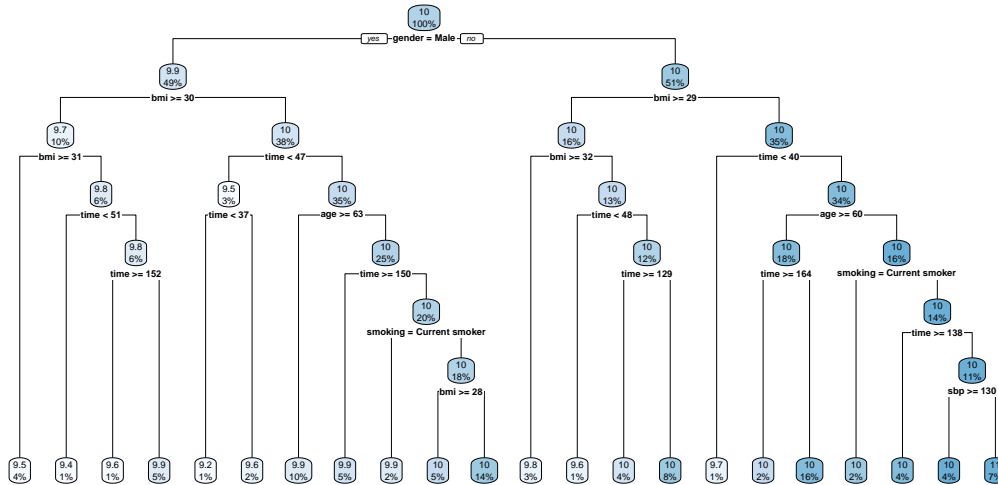
## 183	4.3156e-04	309	0.56543	1.02025	0.019958
## 184	4.3037e-04	310	0.56500	1.02089	0.019962
## 185	4.2651e-04	311	0.56457	1.02147	0.019978
## 186	4.1873e-04	312	0.56414	1.02340	0.019978
## 187	4.1464e-04	313	0.56372	1.02435	0.019991
## 188	4.1452e-04	314	0.56331	1.02502	0.020007
## 189	4.1444e-04	315	0.56289	1.02502	0.020007
## 190	4.1260e-04	316	0.56248	1.02502	0.020007
## 191	4.1181e-04	317	0.56207	1.02500	0.020006
## 192	4.0874e-04	318	0.56165	1.02572	0.020035
## 193	4.0617e-04	320	0.56084	1.02599	0.020041
## 194	4.0575e-04	321	0.56043	1.02665	0.020050
## 195	4.0037e-04	322	0.56002	1.02764	0.020046
## 196	4.0031e-04	323	0.55962	1.02864	0.020049
## 197	3.9993e-04	324	0.55922	1.02864	0.020049
## 198	3.9775e-04	326	0.55842	1.02853	0.020048
## 199	3.9763e-04	328	0.55763	1.02925	0.020098
## 200	3.9213e-04	332	0.55603	1.02902	0.020092
## 201	3.8337e-04	334	0.55524	1.03104	0.020142
## 202	3.8282e-04	335	0.55486	1.03135	0.020163
## 203	3.8077e-04	336	0.55448	1.03234	0.020175
## 204	3.7965e-04	337	0.55410	1.03234	0.020175
## 205	3.7959e-04	338	0.55372	1.03269	0.020175
## 206	3.7735e-04	339	0.55334	1.03277	0.020175
## 207	3.6711e-04	340	0.55296	1.03335	0.020180
## 208	3.6572e-04	341	0.55259	1.03290	0.020171
## 209	3.5205e-04	342	0.55223	1.03293	0.020164
## 210	3.5080e-04	346	0.55082	1.03427	0.020176
## 211	3.3596e-04	348	0.55012	1.03521	0.020159
## 212	3.2867e-04	349	0.54978	1.03551	0.020179
## 213	3.2801e-04	351	0.54912	1.03521	0.020176
## 214	3.2780e-04	352	0.54880	1.03521	0.020185
## 215	3.2374e-04	353	0.54847	1.03504	0.020182
## 216	3.2165e-04	354	0.54814	1.03543	0.020199
## 217	3.2122e-04	357	0.54716	1.03549	0.020182
## 218	3.2005e-04	358	0.54684	1.03591	0.020184
## 219	3.1900e-04	359	0.54652	1.03607	0.020184
## 220	3.1636e-04	360	0.54620	1.03649	0.020193
## 221	3.1401e-04	362	0.54557	1.03685	0.020195
## 222	3.0849e-04	363	0.54525	1.03721	0.020200
## 223	3.0776e-04	364	0.54495	1.03685	0.020196
## 224	3.0525e-04	366	0.54433	1.03731	0.020194
## 225	2.9975e-04	367	0.54402	1.03746	0.020197
## 226	2.9961e-04	368	0.54373	1.03704	0.020197
## 227	2.9629e-04	369	0.54343	1.03748	0.020212
## 228	2.9479e-04	370	0.54313	1.03826	0.020228
## 229	2.9436e-04	371	0.54283	1.03832	0.020226
## 230	2.9106e-04	373	0.54225	1.03845	0.020225
## 231	2.8704e-04	374	0.54195	1.03908	0.020238
## 232	2.7713e-04	376	0.54138	1.03987	0.020237
## 233	2.7325e-04	377	0.54110	1.03985	0.020256
## 234	2.6914e-04	378	0.54083	1.04026	0.020274
## 235	2.6745e-04	380	0.54029	1.04004	0.020265
## 236	2.6078e-04	381	0.54002	1.03973	0.020260

```
## 237 2.5964e-04 382 0.53976 1.03950 0.020250
## 238 2.5461e-04 383 0.53950 1.03933 0.020250
## 239 2.5262e-04 384 0.53925 1.03883 0.020245
## 240 2.5058e-04 385 0.53900 1.03842 0.020245
## 241 2.4999e-04 386 0.53875 1.03854 0.020247
## 242 2.4750e-04 387 0.53850 1.03880 0.020248
## 243 2.4474e-04 388 0.53825 1.03896 0.020256
## 244 2.4108e-04 389 0.53800 1.03937 0.020261
## 245 2.3052e-04 390 0.53776 1.03957 0.020266
## 246 2.2787e-04 391 0.53753 1.04004 0.020322
## 247 2.2524e-04 392 0.53730 1.03995 0.020312
## 248 2.2232e-04 393 0.53708 1.04096 0.020321
## 249 2.1999e-04 394 0.53686 1.04100 0.020320
## 250 2.1731e-04 395 0.53664 1.04179 0.020339
## 251 2.0458e-04 396 0.53642 1.04259 0.020351
## 252 2.0276e-04 397 0.53622 1.04237 0.020353
## 253 2.0102e-04 398 0.53601 1.04226 0.020350
## 254 2.0029e-04 399 0.53581 1.04226 0.020350
## 255 1.9140e-04 400 0.53561 1.04223 0.020343
## 256 1.8924e-04 401 0.53542 1.04261 0.020349
## 257 1.7999e-04 402 0.53523 1.04327 0.020357
## 258 1.7854e-04 403 0.53505 1.04332 0.020364
## 259 1.7478e-04 404 0.53487 1.04350 0.020363
## 260 1.5917e-04 405 0.53470 1.04288 0.020360
## 261 1.3108e-04 406 0.53454 1.04370 0.020376
## 262 1.1037e-04 407 0.53441 1.04395 0.020375
## 263 9.0997e-05 408 0.53430 1.04379 0.020374
## 264 5.9714e-05 409 0.53421 1.04396 0.020373
## 265 0.0000e+00 410 0.53415 1.04433 0.020380
```

```
cpTable = tree_full$cptable
plotcp(tree_full)
```



```
## Find the cp that yields the minimum cross-validation error
minErr = which.min(cpTable[,4])
tree_final = rpart::prune(tree_full, cp = cpTable[minErr, 1])
rpart.plot(tree_final)
```



```
summary(tree_final)
```

```
## Call:
## rpart(formula = log_antibody ~ ., data = train_data1, control = rpart.control(cp = 0))
## n = 5000
##
##          CP nsplit rel error   xerror   xstd
## 1  0.057918182      0 1.0000000 1.0007146 0.02004169
## 2  0.027293800      1 0.9420818 0.9432020 0.01870086
## 3  0.025172140      2 0.9147880 0.9295966 0.01846651
## 4  0.020759383      3 0.8896159 0.9042136 0.01812179
## 5  0.010560350      4 0.8688565 0.8829836 0.01767741
## 6  0.007928454      5 0.8582961 0.8761783 0.01736189
## 7  0.007483148      6 0.8503677 0.8729469 0.01732773
## 8  0.006830408      7 0.8428845 0.8657230 0.01716499
## 9  0.006412912      8 0.8360541 0.8604904 0.01704889
## 10 0.006407846      9 0.8296412 0.8572604 0.01697315
## 11 0.004852623     10 0.8232334 0.8485364 0.01678947
## 12 0.003284421     11 0.8183808 0.8465873 0.01664462
## 13 0.003033803     12 0.8150963 0.8465219 0.01667932
## 14 0.002914817     13 0.8120625 0.8466921 0.01669623
## 15 0.002790737     14 0.8091477 0.8452142 0.01665712
## 16 0.002730117     15 0.8063570 0.8450512 0.01661721
## 17 0.002560197     16 0.8036269 0.8476923 0.01669101
## 18 0.002327749     17 0.8010667 0.8483089 0.01668839
## 19 0.002223936     19 0.7964112 0.8451645 0.01661475
## 20 0.002164228     21 0.7919633 0.8447269 0.01660173
##
## Variable importance
```

```

##      bmi      time  gender      age smoking      sbp      ldl
##      31       28      27        7       3       3       1
##
## Node number 1: 5000 observations,      complexity param=0.05791818
## mean=10.06434, MSE=0.355722
## left son=2 (2427 obs) right son=3 (2573 obs)
## Primary splits:
##      gender splits as  RL,      improve=0.057918180, (0 missing)
##      bmi      < 29.65 to the right, improve=0.049512600, (0 missing)
##      time     < 46.5  to the left, improve=0.041303690, (0 missing)
##      age      < 60.5  to the right, improve=0.018144830, (0 missing)
##      smoking splits as  RRL,      improve=0.008813525, (0 missing)
## Surrogate splits:
##      sbp < 119.5 to the left, agree=0.518, adj=0.007, (0 split)
##      time < 37.5  to the left, agree=0.518, adj=0.006, (0 split)
##      ldl < 140.5 to the right, agree=0.517, adj=0.005, (0 split)
##      age < 66.5  to the right, agree=0.516, adj=0.003, (0 split)
##      bmi < 20.75 to the left, agree=0.515, adj=0.002, (0 split)
##
## Node number 2: 2427 observations,      complexity param=0.02517214
## mean=9.91655, MSE=0.3395789
## left son=4 (524 obs) right son=5 (1903 obs)
## Primary splits:
##      bmi      < 29.85 to the right, improve=0.054323820, (0 missing)
##      time     < 46.5  to the left, improve=0.049666810, (0 missing)
##      age      < 60.5  to the right, improve=0.017649280, (0 missing)
##      smoking splits as  RRL,      improve=0.009762305, (0 missing)
##      ldl      < 141.5 to the right, improve=0.004210017, (0 missing)
##
## Node number 3: 2573 observations,      complexity param=0.0272938
## mean=10.20375, MSE=0.3309125
## left son=6 (820 obs) right son=7 (1753 obs)
## Primary splits:
##      bmi      < 28.95 to the right, improve=0.057015340, (0 missing)
##      time     < 46.5  to the left, improve=0.036007110, (0 missing)
##      age      < 59.5  to the right, improve=0.023746020, (0 missing)
##      smoking splits as  RRL,      improve=0.010252650, (0 missing)
##      sbp      < 128.5 to the right, improve=0.006713466, (0 missing)
## Surrogate splits:
##      ldl < 164.5 to the right, agree=0.683, adj=0.005, (0 split)
##      sbp < 108.5 to the left, agree=0.682, adj=0.002, (0 split)
##
## Node number 4: 524 observations,      complexity param=0.004852623
## mean=9.657717, MSE=0.3166048
## left son=8 (207 obs) right son=9 (317 obs)
## Primary splits:
##      bmi      < 31.45 to the right, improve=0.05202457, (0 missing)
##      time     < 50.5  to the left, improve=0.04364456, (0 missing)
##      smoking splits as  RRL,      improve=0.01948327, (0 missing)
##      sbp      < 120.5 to the right, improve=0.01676732, (0 missing)
##      age      < 59.5  to the right, improve=0.01274723, (0 missing)
## Surrogate splits:
##      ldl < 152   to the right, agree=0.615, adj=0.024, (0 split)
##      time < 31.5  to the left, agree=0.609, adj=0.010, (0 split)

```

```

##      sbp < 109   to the left,  agree=0.607, adj=0.005, (0 split)
##
## Node number 5: 1903 observations,      complexity param=0.02075938
##   mean=9.987821, MSE=0.3223782
##   left son=10 (137 obs) right son=11 (1766 obs)
##   Primary splits:
##     time < 46.5  to the left,  improve=0.060185330, (0 missing)
##     age  < 62.5  to the right, improve=0.023530880, (0 missing)
##     ldl  < 141.5 to the right, improve=0.008688514, (0 missing)
##     bmi  < 28.35 to the right, improve=0.008031965, (0 missing)
##     smoking splits as RRL,      improve=0.007860098, (0 missing)
##
## Node number 6: 820 observations,      complexity param=0.006830408
##   mean=10.00291, MSE=0.3091452
##   left son=12 (162 obs) right son=13 (658 obs)
##   Primary splits:
##     bmi < 32.25 to the right, improve=0.047923770, (0 missing)
##     time < 47.5  to the left, improve=0.047697360, (0 missing)
##     age < 60.5  to the right, improve=0.021396790, (0 missing)
##     ldl < 73.5  to the right, improve=0.011921340, (0 missing)
##     sbp < 137.5 to the right, improve=0.008408437, (0 missing)
##
## Node number 7: 1753 observations,      complexity param=0.01056035
##   mean=10.29769, MSE=0.313402
##   left son=14 (53 obs) right son=15 (1700 obs)
##   Primary splits:
##     time < 39.5  to the left, improve=0.034188120, (0 missing)
##     age  < 59.5  to the right, improve=0.028638110, (0 missing)
##     smoking splits as RRL,      improve=0.017123260, (0 missing)
##     sbp  < 128.5 to the right, improve=0.009433160, (0 missing)
##     bmi  < 26.75 to the right, improve=0.005238462, (0 missing)
##
## Node number 8: 207 observations
##   mean=9.498897, MSE=0.2980599
##
## Node number 9: 317 observations,      complexity param=0.002223936
##   mean=9.761427, MSE=0.3014877
##   left son=18 (27 obs) right son=19 (290 obs)
##   Primary splits:
##     time < 51    to the left, improve=0.03944451, (0 missing)
##     smoking splits as RRL,      improve=0.03307206, (0 missing)
##     age  < 58.5  to the right, improve=0.02350311, (0 missing)
##     sbp  < 121.5 to the right, improve=0.02005357, (0 missing)
##     ldl  < 80.5  to the right, improve=0.01272572, (0 missing)
##
## Node number 10: 137 observations,      complexity param=0.003033803
##   mean=9.487714, MSE=0.3536406
##   left son=20 (46 obs) right son=21 (91 obs)
##   Primary splits:
##     time < 36.5  to the left, improve=0.11137440, (0 missing)
##     sbp  < 118.5 to the right, improve=0.05219794, (0 missing)
##     bmi  < 26.65 to the right, improve=0.04629031, (0 missing)
##     age  < 65.5  to the right, improve=0.04068284, (0 missing)
##     ldl  < 88.5  to the left, improve=0.02321797, (0 missing)

```



```

##
## Node number 11: 1766 observations,      complexity param=0.007483148
##   mean=10.02662, MSE=0.2990454
##   left son=22 (508 obs) right son=23 (1258 obs)
##   Primary splits:
##     age      < 62.5  to the right, improve=0.025202130, (0 missing)
##     time     < 97.5  to the right, improve=0.015378130, (0 missing)
##     ldl      < 142.5 to the right, improve=0.012513270, (0 missing)
##     smoking splits as RRL,      improve=0.009228690, (0 missing)
##     bmi      < 28.35 to the right, improve=0.008363651, (0 missing)
##   Surrogate splits:
##     sbp < 140.5 to the right, agree=0.737, adj=0.085, (0 split)
##     ldl < 159.5 to the right, agree=0.715, adj=0.008, (0 split)
##
## Node number 12: 162 observations
##   mean=9.757604, MSE=0.2946382
##
## Node number 13: 658 observations,      complexity param=0.006407846
##   mean=10.06331, MSE=0.2942538
##   left son=26 (51 obs) right son=27 (607 obs)
##   Primary splits:
##     time < 47.5  to the left,  improve=0.058863320, (0 missing)
##     age  < 61.5  to the right, improve=0.018077840, (0 missing)
##     bmi  < 31.35 to the right, improve=0.011334980, (0 missing)
##     ldl  < 73.5  to the right, improve=0.009520452, (0 missing)
##     sbp  < 137.5 to the right, improve=0.009235878, (0 missing)
##
## Node number 14: 53 observations
##   mean=9.71145, MSE=0.4536552
##
## Node number 15: 1700 observations,      complexity param=0.007928454
##   mean=10.31597, MSE=0.2979808
##   left son=30 (914 obs) right son=31 (786 obs)
##   Primary splits:
##     age      < 59.5  to the right, improve=0.027837610, (0 missing)
##     smoking splits as RRL,      improve=0.016679170, (0 missing)
##     time     < 159.5 to the right, improve=0.014116890, (0 missing)
##     sbp      < 128.5 to the right, improve=0.009671418, (0 missing)
##     bmi      < 26.75 to the right, improve=0.004398652, (0 missing)
##   Surrogate splits:
##     sbp < 126.5 to the right, agree=0.645, adj=0.232, (0 split)
##     ldl < 91.5  to the right, agree=0.569, adj=0.069, (0 split)
##     time < 176.5 to the left,  agree=0.545, adj=0.017, (0 split)
##     bmi  < 21.35 to the right, agree=0.540, adj=0.005, (0 split)
##
## Node number 18: 27 observations
##   mean=9.404035, MSE=0.2771699
##
## Node number 19: 290 observations,      complexity param=0.002223936
##   mean=9.794701, MSE=0.2907525
##   left son=38 (62 obs) right son=39 (228 obs)
##   Primary splits:
##     time      < 151.5 to the right, improve=0.04911456, (0 missing)
##     smoking splits as RRL,      improve=0.04214043, (0 missing)

```

```

##      sbp      < 121.5 to the right, improve=0.02250072, (0 missing)
##      age      < 58.5  to the right, improve=0.02190715, (0 missing)
##      ldl      < 80.5  to the right, improve=0.01055996, (0 missing)
##
## Node number 20: 46 observations
##   mean=9.208578, MSE=0.2594872
##
## Node number 21: 91 observations
##   mean=9.628815, MSE=0.3419384
##
## Node number 22: 508 observations
##   mean=9.890003, MSE=0.3020677
##
## Node number 23: 1258 observations,   complexity param=0.006412912
##   mean=10.08178, MSE=0.287245
##   left son=46 (250 obs) right son=47 (1008 obs)
##   Primary splits:
##     time      < 149.5 to the right, improve=0.031564790, (0 missing)
##     age       < 54.5  to the right, improve=0.011762700, (0 missing)
##     ldl       < 126.5 to the right, improve=0.011752420, (0 missing)
##     smoking splits as RRL,         improve=0.010707140, (0 missing)
##     sbp       < 109.5 to the left,  improve=0.007957896, (0 missing)
##   Surrogate splits:
##     ldl < 162.5 to the right, agree=0.802, adj=0.004, (0 split)
##
## Node number 26: 51 observations
##   mean=9.609269, MSE=0.2289111
##
## Node number 27: 607 observations,   complexity param=0.002730117
##   mean=10.10146, MSE=0.2809679
##   left son=54 (225 obs) right son=55 (382 obs)
##   Primary splits:
##     time < 128.5 to the right, improve=0.02847191, (0 missing)
##     age  < 61.5  to the right, improve=0.02293396, (0 missing)
##     bmi  < 31.35 to the right, improve=0.01993356, (0 missing)
##     ldl  < 111.5 to the right, improve=0.01092987, (0 missing)
##     sbp  < 137.5 to the right, improve=0.01007648, (0 missing)
##   Surrogate splits:
##     sbp < 114.5 to the left, agree=0.636, adj=0.018, (0 split)
##     age < 52.5  to the left, agree=0.631, adj=0.004, (0 split)
##     bmi < 32.15 to the right, agree=0.631, adj=0.004, (0 split)
##
## Node number 30: 914 observations,   complexity param=0.002914817
##   mean=10.23151, MSE=0.2886223
##   left son=60 (101 obs) right son=61 (813 obs)
##   Primary splits:
##     time      < 163.5 to the right, improve=0.019652420, (0 missing)
##     smoking splits as RRL,         improve=0.012542120, (0 missing)
##     bmi       < 24.85 to the right, improve=0.010077130, (0 missing)
##     age       < 70.5  to the right, improve=0.004980359, (0 missing)
##     ldl       < 79.5  to the right, improve=0.004331232, (0 missing)
##
## Node number 31: 786 observations,   complexity param=0.002790737
##   mean=10.41418, MSE=0.2909223

```

```

## left son=62 (78 obs) right son=63 (708 obs)
## Primary splits:
## smoking splits as RRL, improve=0.021707010, (0 missing)
## time < 184.5 to the right, improve=0.020284770, (0 missing)
## sbp < 129.5 to the right, improve=0.013958710, (0 missing)
## race splits as LRL, improve=0.008200332, (0 missing)
## age < 54.5 to the right, improve=0.007226599, (0 missing)
## Surrogate splits:
## sbp < 149 to the right, agree=0.903, adj=0.026, (0 split)
##
## Node number 38: 62 observations
## mean=9.565542, MSE=0.2656172
##
## Node number 39: 228 observations
## mean=9.857017, MSE=0.2794241
##
## Node number 46: 250 observations
## mean=9.890584, MSE=0.2961252
##
## Node number 47: 1008 observations, complexity param=0.003284421
## mean=10.12921, MSE=0.273727
## left son=94 (87 obs) right son=95 (921 obs)
## Primary splits:
## smoking splits as RRL, improve=0.021171970, (0 missing)
## bmi < 28.25 to the right, improve=0.016615490, (0 missing)
## ldl < 126.5 to the right, improve=0.013732310, (0 missing)
## age < 52.5 to the right, improve=0.009102790, (0 missing)
## sbp < 109.5 to the left, improve=0.008120837, (0 missing)
##
## Node number 54: 225 observations
## mean=9.984915, MSE=0.2376411
##
## Node number 55: 382 observations
## mean=10.1701, MSE=0.293776
##
## Node number 60: 101 observations
## mean=10.01783, MSE=0.30527
##
## Node number 61: 813 observations
## mean=10.25805, MSE=0.2801773
##
## Node number 62: 78 observations
## mean=10.17476, MSE=0.3236525
##
## Node number 63: 708 observations, complexity param=0.002327749
## mean=10.44056, MSE=0.2803057
## left son=126 (182 obs) right son=127 (526 obs)
## Primary splits:
## time < 137.5 to the right, improve=0.019352530, (0 missing)
## sbp < 133.5 to the right, improve=0.014196050, (0 missing)
## age < 55.5 to the right, improve=0.009180245, (0 missing)
## race splits as LRL, improve=0.008380573, (0 missing)
## bmi < 28.45 to the right, improve=0.007482769, (0 missing)
## Surrogate splits:

```

```

##      ldl < 154.5 to the right, agree=0.744, adj=0.005, (0 split)
##
## Node number 94: 87 observations
##   mean=9.881514, MSE=0.2388818
##
## Node number 95: 921 observations,      complexity param=0.002560197
##   mean=10.1526, MSE=0.2706758
##   left son=190 (236 obs) right son=191 (685 obs)
##   Primary splits:
##     bmi < 28.25 to the right, improve=0.018266070, (0 missing)
##     ldl < 127.5 to the right, improve=0.013821740, (0 missing)
##     age < 55.5 to the right, improve=0.012300610, (0 missing)
##     sbp < 109.5 to the left, improve=0.011041330, (0 missing)
##     time < 95.5 to the right, improve=0.007309987, (0 missing)
##   Surrogate splits:
##     ldl < 164 to the right, agree=0.746, adj=0.008, (0 split)
##     age < 46.5 to the left, agree=0.745, adj=0.004, (0 split)
##
## Node number 126: 182 observations
##   mean=10.31535, MSE=0.2119669
##
## Node number 127: 526 observations,      complexity param=0.002327749
##   mean=10.48388, MSE=0.2966499
##   left son=254 (191 obs) right son=255 (335 obs)
##   Primary splits:
##     sbp < 129.5 to the right, improve=0.028452590, (0 missing)
##     time < 50.5 to the left, improve=0.013522720, (0 missing)
##     race splits as LRL, improve=0.012365230, (0 missing)
##     bmi < 28.35 to the right, improve=0.009747016, (0 missing)
##     age < 54.5 to the right, improve=0.007912211, (0 missing)
##   Surrogate splits:
##     ldl < 145.5 to the right, agree=0.646, adj=0.026, (0 split)
##     bmi < 22.7 to the left, agree=0.644, adj=0.021, (0 split)
##     time < 42.5 to the left, agree=0.641, adj=0.010, (0 split)
##
## Node number 190: 236 observations
##   mean=10.03281, MSE=0.2619733
##
## Node number 191: 685 observations
##   mean=10.19387, MSE=0.2670264
##
## Node number 254: 191 observations
##   mean=10.36221, MSE=0.2972116
##
## Node number 255: 335 observations
##   mean=10.55325, MSE=0.2830768

```