

Applied Survival Analysis Using R

Chapter 6: Model Selection and Interpretation

Qi Guo

Department of Mathematical Sciences
The University of Texas at Dallas

April, 13 2019

- 1 Covariate Adjustment
- 2 Categorical and Continuous Covariates
- 3 Nested Models
- 4 The Akaike Information Criterion for Comparing Non-nested Models

Covariate Adjustment

- Most of our study is a randomized clinical trial, the focus will be on *comparing the effectiveness of different treatments*.
- A successful randomization procedure should ensure that *confounding covariates are balanced* between the treatments.
- **Goal:** Use methods to **sift through** a potentially large number of potential explanatory variables to find the important ones.
- **Example**(`coxph` model of the effect of treatment on survival unadjusted for the genetic mutation status of the patients)

```

1 > coxph(Surv(ttAll, status) ~ trt)
2      coef      exp(coef)      se(coef)      z      p
3 trt 0.464        1.59        0.117      3.96  7.6e-05
4
5 Likelihood ratio test=15.5 on 1 df, p=8.2e-05

```

Conclusion and stratify

- Conclusion:

- The estimate of the log hazard ratio treatment effect, $\hat{\beta}$, is 0.464, so higher hazards are associated with the treatment than with the control, unfortunate result.
- The value of $e^{\hat{\beta}} = 1.59$, suggesting (incorrectly, as we know) that the treatment is associated with a 59% additional risk of death over the control.

- Stratify on **genotype**:

```

1 > coxph(Surv(ttAll, status) ~ trt + strata(genotype))
2 coef      exp(coef)      se(coef)      z      p
3 trt -0.453      0.636      0.164     -2.76    0.0058
4
5 Likelihood ratio test=7.66 on 1 df, p=0.00566

```

- Conclusion: The coefficient is negative, indicating that, within each genotype, *the treatment is effective*.

Explicitly estimating the genetic effect

```

1 > coxph(Surv(ttAll, status) ~ trt + genotype)
2           coef      exp(coef)      se(coef)      z      p
3 trt          -0.453         0.636        0.164    -2.76  0.0058
4 genotypewt    -1.568         0.209        0.183   -8.59  0.0000
5 Likelihood ratio test=93.4 on 2 df, p=0

```

- Conclusion: The **wild type genotype** has lower hazard than the **reference (mutant) genotype**, and thus that the mutant genotype incurs additional risk of death.

Indicator or Dummy variable

- The previous sections considered a partial likelihood for comparing two groups, indexed by a covariate z . Since z can take the values 0 or 1 depending on which of two groups a subject belongs to, this covariate is called an *indicator or dummy variable*.
- When *categorical variables* with three or more variables, we will need multiple dummy variables. For example:
 - If a research question is how survival in *non-white* groups compares to survival in *whites*, one would select “white” as the reference variable. Since there are four levels, we need to create *three* dummy variables, say, z_2 , z_3 , and z_4 to represent “race”. Then for a *white* patient, all three would take the value zero. For an Asian person, we would have $z_2 = 1$, and $z_1 = z_3 = 0$.

k covariates model and enhance

- k covariates model:

$$\log(\psi_i) = z_{i1}\beta_1 + z_{i2}\beta_2 + \cdots + z_{ik}\beta_k. \quad (1)$$

- For each covariate, the parameter β_j is the log hazard ratio for the effect of that parameter on survival, adjusting for the other covariates.
- Matrix form: $\log(\psi_i) = z_i'\beta$ (for Patient i), where z_i' (the transpose of z_i) is a $1 \times k$ matrix (i.e. a row matrix) of covariates, and β is a $k \times 1$ matrix (i.e. a column matrix) of parameters.
- **Enhance** the model:
 - \otimes If a continuous variable is not linearly related to the log hazard, transform it using, eg: a logarithmic or square root function.
 - \otimes “discretize” a variable, eg: split the “age” into three pieces, “under 50” and “50–64”, and “65 and above” and entered into the model as a categorical variable.

Difference

- But it is incorporate for *interaction* terms.
- Difference with linear and logistic regression model:
 - Survival data can evolve over time, there is a possibility that some covariate values may also *change as time passes*. eg: Time-related variables like `age` must also be defined and fixed by taking their value at the beginning of the trial, even though patients will age as the trial progresses(Chapter 8).
 - There is *no intercept term* in proportional hazards models, if there were one, it would be absorbed into the baseline hazard(canceled out in `num` and `den`).

Example

Tidy the data

Suppose that we have two black patients, two white patients, and two patients of other races, with ages 48, 52, 87, 82, 67, and 53, respectively. We may enter these data values as follows:

```
1 > race <- factor(c("black", "black", "white", "white", "other", "other"))
2 age <- c(48, 52, 87, 82, 67, 53)
```

- Create matrix using “model.matrix” function (In my “R for data science presentation”)

```
1 > model.matrix(~ race + age)[,-1]
2      raceother      racewhite      age
3 1          0              0       48
4 2          0              0       52
5 3          0              1       87
6 4          0              1       82
7 5          1              0       67
8 6          1              0       53
```

Example

- If we need to use `whites` as the reference, we can change the race factor to have “whites” as the reference level

```

1 > race <- relevel(race, ref="white")
2 > model.matrix(~ race + age)[,-1]
3      raceblack      raceother      age
4 1           1             0       48
5 2           1             0       52
6 3           0             0       87
7 4           0             0       82
8 5           0             1       67
9 6           0             1       53

```

- Three covariates, say, z_1 , z_2 , and z_3 , the first two of which are dummy variables for black race and other race, and the third a continuous variable, `age`.
- For a black 48-year old person, the **log hazard ratio** is:

$$\log(\psi_1) = z_{11}\beta_1 + z_{12}\beta_2 + z_{13}\beta_3 = 1 \times \beta_1 + 0 \times \beta_2 + 48 \times \beta_3. \quad (2)$$

Example

- β_1 represents the log hazard ratio for blacks as compared to whites, and β_3 represents the change in log hazard ratio that would correspond to a one-year change in age.
- Add **Interaction**:

```

1 > model.matrix(~ race + age + race:age)[,-1]
2      raceblack  raceother  age  raceblack:age  raceother:age
3 1           1           0   48             48             0
4 2           1           0   52             52             0
5 3           0           0   87              0             0
6 4           0           0   82              0             0
7 5           0           1   67              0            67
8 6           0           1   53              0            53

```

Example

Simulation

Generate a small survival data set and show how models are incorporated into a survival problem

- Generate 60 ages between 40 and 80 at random and categorize, and make “white” the reference category:

```
1 > age <- runif(n=60, min=40, max=80)
2 > race <- factor(c(rep("white", 20), rep("black", 20),
3 rep("other", 20)))
4 > race <- relevel(race, ref="white")
```

- The variables are exponentially distributed with a particular rate parameter that depends on the covariates, specified the log rate parameter to have baseline -4.5, and “age” increase the log rate by 0.05 per year:

```
1 > log.rate.vec <- -4.5 + c(rep(0,20), rep(1,20), rep(2,20))
2 + age*0.05
```

Example

- No censoring

```
1 >tt <- rexp(n=60, rate=exp(log.rate.vec))
2 >status <- rep(1, 60)
```

- Fit a *Cox proportional hazards model*

```
1 > library(survival)
2 > result.cox <- coxph(Surv(tt, status) ~ race + age)
3 > summary(result.cox)
4 n= 60, number of events= 60
5      coef      exp(coef)    se(coef)      z      Pr(>|z|)
6 raceblack  1.15154      3.16305    0.36752    3.133    0.00173 **
7 raceother  2.49905     12.17087    0.42936    5.820    5.87e-09 ***
8 age        0.07798      1.08110    0.01448    5.385    7.24e-08 ***
9 ---
10 Signif. codes:  0 *** 0.001 ** 0.01 * 0.05 . 0.1 1
```

- Conclusion:

- The coefficient estimates, 1.15, 2.50, and 0.08, are close to the true values from the simulation, (1, 2, and 0.05).
- The blacks have 3.16 times the risk of death as do whites.

Nested Models

Nested Models

When comparing two models, the covariates of one model must be a **subset** of the covariates in the other.

- Example:

Model A: ageGroup4

Model B: employment

Model C: ageGroup4 + employment

Here, Model A is nested in Model C, and Model B is also nested in Model C, so these models can be compared using statistical tests

```
1 > levels(ageGroup4)
2 [1] "21-34" "35-49" "50-64" "65+"
3 > levels(employment) # "ft" refers to full-time, "pt" is part-time
4 [1] "ft" "other" "pt"
```

Coxph model

Model A:

```

1 > modelA.coxph <- coxph(Surv(ttr, relapse) ~ ageGroup4)
2 > modelA.coxph
3               coef      exp(coef)    se(coef)      z        p
4 ageGroup435-49  -0.453      0.636      0.164    -2.76    0.0058
5 ageGroup450-64  -1.568      0.209      0.183    -8.59    0.0000
6 ageGroup465+    -0.3173     0.728      0.444    -0.7153   0.470
7
8 Likelihood ratio test=12.2 on 3 df, p=0.00666

```

Model B:

```

1 > modelB.coxph <- coxph(Surv(ttr, relapse) ~ employment)
2 > modelB.coxph
3               coef      exp(coef)    se(coef)      z        p
4 employmentother -0.453      0.636      0.164    -2.76    0.0058
5 employmentpt    -1.568      0.209      0.183    -8.59    0.0000
6
7 Likelihood ratio test=2.06 on 2 df, p=0.357

```

Coxph model

Model C:

```

1 > modelC.coxph <- coxph(Surv(ttr, relapse) ~ ageGroup4 +
2 employment)
3 > modelC.coxph
4
5      coef      exp(coef)    se(coef)      z      p
6 ageGroup435-49 -0.130      0.878      0.321    -0.404  0.6900
7 ageGroup450-64 -1.024      0.359      0.359    -2.856  0.0043
8 ageGroup465+   -0.782      0.457      0.505    -1.551  0.1200
9 employmentother  0.526      1.692      0.275     1.913  0.0560
10 employmentpt   0.500      1.649      0.332     1.508  0.1300
11
12 Likelihood ratio test=16.8 on 5 df, p=0.00492

```

The log-likelihoods:

```

1 > logLik(modelA.coxph)
2 'log Lik.' -380.043 (df=3)
3 > logLik(modelB.coxph)
4 'log Lik.' -385.1232 (df=2)
5 > logLik(modelC.coxph)
6 'log Lik.' -377.7597 (df=5)

```


Compare

- Determining if “employment” belongs in the model by comparing Models A and C, the hypothesis test:
 - H_0 : The three coefficients for “employment” are zero.
 H_A : Not all zero.
- The likelihood ratio statistic is:

$$2(\ell(\hat{\beta}_{full}) - \ell(\hat{\beta}_{reduced})) = 2(-377.7597 + 380.043) = 4.567 \quad (3)$$

- Compare this to a chi-square distribution with $5 - 3 = 2$ degrees of freedom.

```
1 > pchisq(4.567, df=2, lower.tail=F)
2 [1] 0.1019268
```

- Conclusion: The effect of “employment” is not statistically significant when “ageGroup4” in the model.

ANOVA

- In our **STAT 6337:Advanced Statistic Method I** class, we use “anova” function, which is more direct.

```

1 > anova(modelA.coxph, modelC.coxph)
2 Analysis of Deviance Table
3 Cox model: response is Surv(ttr, relapse)
4 Model 1: ~ ageGroup4
5 Model 2: ~ ageGroup4 + employment
6      loglik      Chisq      Df      P(>|Chi|)
7 1 -380.04
8 2 -377.76      4.5666      2      0.1019

```

- The Akaike Information Criterion, or **AIC**:

$$AIC = -2 \cdot \ell(\hat{\beta}) + 2 \cdot k \quad (4)$$

- where $\ell(\hat{\beta})$ denotes the value of the partial log likelihood at the M.P.L.E. for a particular model, and k is the number of parameters in the model.
- The AIC balances two quantities which are properties of a model:
 - **Goodness of fit**: $-2 \cdot \ell(\hat{\beta})$, this quantity is smaller for models that fit the data well
 - The number of parameters is a measure of complexity.
- Conclusion: Smaller is better.

AIC

- Compute the AIC for model A:
 $AIC = 2 \times 380.043 + 2 \times 2 = 766.086$
- Use the “AIC” function:

```
1 > AIC(modelA.coxph)
2 [1] 766.086
3 > AIC(modelB.coxph)
4 [1] 774.2464
5 > AIC(modelC.coxph)
6 [1] 765.5194
```

- Conclusion: The best fitting model from among these three, using the AIC criterion, is Model C.

stepwise procedure

Using the AIC criterion:

```
1 > modelAll.coxph <- coxph(Surv(ttr, relapse) ~ grp + gender +
2 race + employment + yearsSmoking + levelSmoking +
3 ageGroup4 + priorAttempts + longestNoSmoke)
4 > result.step <- step(modelAll.coxph, scope=list(upper=~ grp +
5 gender + race + employment + yearsSmoking +
6 levelSmoking + ageGroup4 + priorAttempts + longestNoSmoke,
7 lower=~grp) )
```

```
1 Start: AIC=770.2 Surv(ttr, relapse) ~ grp + gender + race +
2 employment + yearsSmoking + levelSmoking + ageGroup4 +
3 priorAttempts + longestNoSmoke
4 - race                Df          AIC
5 - yearsSmoking        3        766.98
6 - gender              1        768.20
7 - priorAttempts       1        768.20
8 - levelSmoking        1        768.24
9 - longestNoSmoke      1        768.47
10 <none>                0        770.20
11 - employment         2        772.45
12 - ageGroup4          3        774.11
```

stepwise procedure

- Conclusion: The terms ordered from the one which, when deleted, yields the greatest AIC reduction (“race” in this case) to the smallest reduction (“ageGroup4”). Thus, “race” is deleted.
- Last step:

```

1 Step: AIC=758.42 Surv(ttr, relapse) ~ grp + employment
2 + ageGroup4
3
4                                     Df      AIC
5 <none>                                758.42
6 + longestNoSmoke                      1      759.10
7 - employment                          2      760.31
8 + yearsSmoking                        1      760.34
9 + gender                              1      760.39
10 + priorAttempts                      1      760.40
11 + levelSmoking                       1      760.41
12 + race                               3      761.53
13 - ageGroup4                          3      767.24

```

- The “+” sign shows the effect on AIC of adding certain terms. This table shows that no addition or subtraction of terms results in further reduction of the AIC.

BIC

- An alternative to the AIC is the “Bayesian Information Criterion”, sometimes called the “Schwartz criterion”, *BIC* is given by:

$$BIC = -2 \cdot \log(L) + k \cdot \log(n) \quad (5)$$

- BIC penalizes the number of parameters by a factor of $\log(n)$ rather than by a factor of 2 as in the AIC.
- The BIC in model selection will tend to result in models with **fewer** parameters as compared to AIC.