

Self-assessment test of prerequisite knowledge for Biostatistics III in R

Mark Clements, Karolinska Institutet

2017-10-31

- Participants in the course Biostatistics III are expected to have prerequisite knowledge equivalent to the learning outcomes of the courses Biostatistics I and Biostatistics II. In particular, participants should be comfortable interpreting the output from logistic regression models and we expect course participants to understand:
 1. how to interpret regression coefficients after fitting a logistic regression
 2. assessing confounding in a modelling framework
 3. assessing effect modification (interactions) in a modelling framework
 4. how to conduct a formal hypothesis tests (Wald and likelihood ratio tests) in a modelling framework
- This document contains a self-assessment test of the key concepts you are expected to understand prior to the course. Brief answers are provided at the end of this document. If you have difficulty with any questions we recommend you consult previous course notes and/or course texts book or consult a colleague.
- The questions are typical of exam questions from earlier biostatistics courses and the marks (in brackets) reflect the level of difficulty. If you attempt the test under examination conditions (i.e., without referring to the answers) we would recommend:
 1. if you score 70% or more then you possess the required prerequisite knowledge;
 2. if you score 40%-70% you should brush up on the areas where you lost marks;
 3. if you score less than 40% you should, at a minimum, undertake an extensive review of central concepts in statistical modelling and possibly consider studying intermediate-level courses (e.g., Biostatistics II) before taking Biostatistics III.
- Questions about this test should be addressed to Mark Clements (<mailto:mark.clements@ki.se>) via e-mail.

All questions are based on data from a cohort study designed to study risk factors for incidence of coronary heart disease (CHD). We will study three exposures of interest, body mass index (BMI), job type (3 categories) and energy intake (classified as high or low and where high is considered exposed). The R output shown on this page is not central to the question but is shown for completeness. The output below shows how a variable for BMI has been created and how job type and energy intake are coded. The data are available on the web (see the use statement below) so it is possible for you to reproduce all analyses shown in this document. There is also a do file available (<http://biostat3.net/download/self-assessment.R>).

We have analysed the data using logistic regression, which is not completely appropriate given that these data are from a cohort study where individuals were at risk for different amounts of time. For the purpose of this exercise you should interpret the results from the models as if logistic regression was appropriate. During Biostatistics III we will reanalyse these data using more appropriate methods (e.g., Cox regression and Poisson regression).

```

> library(foreign)
## function to calculate odds ratios (exponential form)
> eform <- function (object) {
  exp(cbind("exp(coef)" = coef(object), confint.default(object)))
}
> diet <- read.dta("http://biostat3.net/download/diet.dta")
## Generate a variable containing BMI
> diet <- transform(diet, bmi=weight/(height/100)^2)
> summary(diet$bmi)
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.   NA's
 15.88  21.59  24.11  24.12  26.50  33.29    5
> table(diet$job)

  driver conductor    bank
    102      84     151
> table(diet$hieng)

 low high
155  182

```

We now estimate a logistic regression model where the outcome is CHD (0 = No CHD 1 = CHD) and the exposures are coded as described above.

```

## Model 1
> model1 <- glm(chd ~ hieng + job + bmi, data=diet, family=binomial)
> summary(model1)

```

Call:

```
glm(formula = chd ~ hieng + job + bmi, family = binomial, data = diet)
```

Deviance Residuals:

Min	1Q	Median	3Q	Max
-0.8889	-0.6028	-0.4794	-0.4068	2.3551

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-3.62891	1.32946	-2.730	0.00634 **
hienghigh	-0.78827	0.33700	-2.339	0.01933 *
jobconductor	0.58399	0.44335	1.317	0.18777
jobbank	0.15623	0.39868	0.392	0.69515
bmi	0.07945	0.05225	1.521	0.12834

Signif. codes: 0 *** 0.001 ** 0.01 * 0.05 . 0.1 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 263.47 on 331 degrees of freedom
Residual deviance: 255.69 on 327 degrees of freedom
(5 observations deleted due to missingness)
AIC: 265.69

Number of Fisher Scoring iterations: 5

```

> eform(model1)
              exp(coef)      2.5 %    97.5 %
(Intercept) 0.02654522 0.001960413 0.3594389
hienghigh   0.45463158 0.234856575 0.8800685
jobconductor 1.79317539 0.752036418 4.2756945
jobbank     1.16909719 0.535168740 2.5539388
bmi         1.08269320 0.977310034 1.1994398

```

1. (1 mark) Interpret the estimated odds ratio for BMI, including a comment on statistical significance.

2. (2 marks) Is it possible to ascertain, using the output on this page, whether the effect of high energy intake is modified by BMI? If so, comment on whether the effect of high energy intake is modified by BMI. If not, describe how you could study this.

3. (1 mark) Both P-values for the parameters representing the effect of occupation are greater than 0.1. Does this mean that there is no evidence of a statistically significant overall association between occupation and CHD risk? If not, how could you test whether there is an association between occupation and CHD risk?

4. (1 marks) What is the estimated odds ratio for individuals working as bankers compared to conductors?

5. (1 mark) Individuals with a high energy intake (≥ 2750 kcals/day) appear to have a statistically significant lower risk of CHD compared to individuals with a low energy intake (< 2750 kcals/day). Should we recommend individuals with a low energy intake to eat more as a means of reducing CHD risk?

We now fit another model (labelled model 2).

```

## Model 2
> model2 <- glm(chd ~ hieng + bmi, data=diet, family=binomial)
> summary(model2)

```

```
Call:
glm(formula = chd ~ hieng + bmi, family = binomial, data = diet)
```

```
Deviance Residuals:
```

Min	1Q	Median	3Q	Max
-0.7739	-0.5865	-0.4827	-0.4229	2.2956

```
Coefficients:
```

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-2.98008	1.21973	-2.443	0.0146 *
hienghigh	-0.75899	0.33403	-2.272	0.0231 *
bmi	0.06159	0.05035	1.223	0.2213

```
---
```

```
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
(Dispersion parameter for binomial family taken to be 1)
```

```
Null deviance: 263.47 on 331 degrees of freedom
Residual deviance: 257.56 on 329 degrees of freedom
(5 observations deleted due to missingness)
AIC: 263.56
```

```
Number of Fisher Scoring iterations: 4
```

```
> eform(model2)
```

	exp(coef)	2.5 %	97.5 %
(Intercept)	0.05078865	0.00465085	0.5546271
hienghigh	0.46813905	0.24324747	0.9009515
bmi	1.06352566	0.96357659	1.1738422

- (1 marks) Based on model 2, among individuals with a BMI of 24, what is the estimated odds ratio for individuals with a high energy compared to those with a low energy intake? You do not have to comment on statistical significance.
- (2 marks) Based on model 2, what is the estimated odds ratio for individuals with a BMI of 30 compared to individuals with a BMI of 25? Is the difference statistically significant?
- (2 marks) Is it possible to ascertain, using the output from models 1 and/or 2, whether the effect of high energy intake is modified by job type? If so, comment on whether the effect of high energy intake is modified by job type. If not, describe how you could study this.

9. (2 marks) Is it possible to ascertain, using the output from models 1 and/or 2, whether the effect of high energy intake is confounded by job type? If so, comment on whether the effect of high energy intake is confounded by job type. If not, describe how you could study this.
10. (3 marks) Based on models 1 and/or 2, is there any evidence that job type is associated with CHD incidence? Conduct a formal hypothesis test using output from models 1 and/or 2. You should state the null hypothesis, alternative hypothesis, value of a test statistic, assumed distribution of the test statistic under the null hypothesis, the name of the statistical test you are using, and a comment on statistical significance.

We now reuse model 1, but report the estimated coefficients rather than the estimated odds ratios. We will label this Model 3 even though it is technically the same model as Model 1 but with estimates presented on a different scale.

```
## Model 3
> model3 <- model1
> summary(model3)
```

Call:

```
glm(formula = chd ~ hieng + job + bmi, family = binomial, data = diet)
```

Deviance Residuals:

Min	1Q	Median	3Q	Max
-0.8889	-0.6028	-0.4794	-0.4068	2.3551

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-3.62891	1.32946	-2.730	0.00634 **
hienghigh	-0.78827	0.33700	-2.339	0.01933 *
jobconductor	0.58399	0.44335	1.317	0.18777
jobbank	0.15623	0.39868	0.392	0.69515
bmi	0.07945	0.05225	1.521	0.12834

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 263.47 on 331 degrees of freedom
 Residual deviance: 255.69 on 327 degrees of freedom
 (5 observations deleted due to missingness)
 AIC: 265.69

Number of Fisher Scoring iterations: 5

11. (2 marks) What is the standard error and 95 percent confidence interval of the estimate for hieng? That is, what are the numbers indicated by X, Y and Z? You may make use of output from models 12 in your answer.

12. (1 mark) What is the interpretation of the intercept (i.e., the coefficient labelled (Intercept))?

```
> model4 <- glm(chd ~ hieng*job + bmi, data=diet, family=binomial)
> summary(model4)
```

Call:

```
glm(formula = chd ~ hieng * job + bmi, family = binomial, data = diet)
```

Deviance Residuals:

Min	1Q	Median	3Q	Max
-0.8669	-0.6019	-0.4865	-0.3956	2.3974

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-3.51644	1.36411	-2.578	0.00994 **
hienghigh	-0.96949	0.65116	-1.489	0.13652
jobconductor	0.46260	0.57680	0.802	0.42255
jobbank	0.07198	0.51302	0.140	0.88842
bmi	0.07768	0.05252	1.479	0.13914
hienghigh:jobconductor	0.29458	0.88619	0.332	0.73958
hienghigh:jobbank	0.21684	0.82026	0.264	0.79151

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 263.47 on 331 degrees of freedom
Residual deviance: 255.58 on 325 degrees of freedom
(5 observations deleted due to missingness)
AIC: 269.58

Number of Fisher Scoring iterations: 5

```
> eform(model4)
```

	exp(coef)	2.5 %	97.5 %
(Intercept)	0.02970487	0.002049705	0.4304908
hienghigh	0.37927458	0.105848046	1.3590162
jobconductor	1.58819688	0.512777999	4.9190280
jobbank	1.07463303	0.393164419	2.9372855
bmi	1.08078028	0.975054654	1.1979698
hienghigh:jobconductor	1.34256502	0.236379979	7.6253532
hienghigh:jobbank	1.24214121	0.248863632	6.1998404

13. (2 marks) What is the OR of high energy intake compared to low for the 3 different job types?

14. (3 marks) Based on models 3 and/or 4, is there any evidence of statistically significant effect modification? Conduct a formal hypothesis test using output from models 3 and/or 4. You should state the null hypothesis, alternative hypothesis, value of a test statistic, assumed distribution of the test statistic under the null hypothesis, the name of the statistical test you are using, and a comment on statistical significance.

Solutions

1. After adjusting for total energy intake (in two categories) and job type (in three categories) it is estimated that the odds of CHD incidence increases by a factor of 1.08 (and every 1 unit increase in BMI).
2. No it is not possible. That would require evaluating if there is an interaction between energy intake and BMI. This can be done by refitting the model with the relevant interaction term and subsequently performing a likelihood ratio test or a Wald test for that effect. Our decision on whether or not effect modification exists should then be based on the size and statistical significance of the interaction effect as well as knowledge of the underlying biology/physiology.
3. The p-values for the parameters representing the effect of occupation represent the pairwise comparison and we should not make a conclusion based on those tests alone. In order to test for a global (overall effect) occupation on CHD risk we could conduct a joint test of the two parameters representing occupation, e.g., a likelihood ratio test or a Wald test (see question 10).
4. The OR is given by $\frac{1.169}{1.793} = 0.652$
5. No, we should always be wary of interpreting associations as causal effects. In this specific case we would expect the association to be confounded by, for example, level of physical activity.
6. OR = 0.468. The OR is assumed to be the same within any level of BMI since the model does not account for possible effect modification.
7. OR = $(1.064)^5 = 1.364$. The effect is not statistically significant (the scale that is used, i.e. a one unit increase or a five unit increase does not affect the significance).
8. No it is not possible to assess effect modification based on the results from model 1 and/or 2. In order to do so we would need to include an interaction term between high energy intake and attained age.
9. There is no formal test for testing for confounding. If the effect of high energy was confounded by job type we would expect to see a substantial difference in the OR representing the effect of energy intake if we include job type in the model compared to when it is left out. The OR for energy intake goes from 0.468 to 0.455 so there is no convincing evidence of confounding by job type.
10. We can perform a likelihood ratio test by testing the null hypotheses that the 2 parameters representing the effect of job type are 0 against the alternative hypothesis that at least one of parameters is non-zero. That is, we test whether the likelihood for the more elaborate model is statistically greater than the likelihood for the reduced model. The test statistic is: $D : 2(\ln L_{(submodel)} - \ln L_{(fullmodel)}) = 2(128.78 - 127.84724) = 1.86552$ Under the null hypothesis, the test statistic follows a χ^2 distribution with 2 df (the difference in the number of parameters between the two models). The critical value of a χ^2 with 2 degrees of freedom is 5.99 at the 5% significance level. Since our test statistic is less than the critical value we conclude that there is no evidence that job type is statistically significant.

In R

```
> anova(model2, model1, test="Chisq")  
Analysis of Deviance Table
```

```

Model 1: chd ~ hieng + bmi
Model 2: chd ~ hieng + job + bmi
  Resid. Df Resid. Dev Df Deviance Pr(>Chi)
1         329      257.56
2         327      255.69  2    1.8655  0.3935

```

11. We can retrieve the standard error from model 1 by taking the log of the confidence limits, i.e. $Y = -1.44817$ and $Z = -0.127833$. The standard error is thus given by $\frac{0.788(1.44817)}{1.96} = 0.336821$ by re-organizing the formula for how to calculate .e.g. the lower confidence limit and solving for the standard error.
12. The constant represents the $\log(\text{odds})$ for an individual where all covariates are at their reference level (i.e., for a driver with low energy intake and $\text{BMI} = 0$). The constant does not always make any sense in practice (as in this case). We can nevertheless calculate $\exp(3.63) = 0.027$. This is the estimated odds of CHD for a driver with low energy intake and BMI of zero. The estimated odds of CHD for a driver with low energy intake and BMI of 25 is given by $\exp(3.63 + 25 \times 0.079) = 0.19$.
13. For drivers the $\text{OR} = 0.379$, for conductors the $\text{OR} = 0.379 \times 1.343 = 0.509$ and for bankers the $\text{OR} = 0.379 \times 1.242 = 0.471$
14. Use a likelihood ratio test as in Question 10. The test statistic is 0.12 which follows a χ^2 distribution with 2 df. We conclude that there is no evidence of a statistically significant interaction.