

Categorical Data Analysis - Simple Logistic Regression

Alessio Crippa

Department of Medical Epidemiology and Biostatistics
Karolinska Institutet

Summary recap

Binary predictor

Continuous predictor

Categorical predictor

Summary recap

For a binary outcome Y , we considered two *absolute* measures (of disease occurrence):

- p the probability ($E[Y]$)
- the odds ($E[Y]/(1 - E[Y])$)

If we want to relate the outcome Y with a binary predictor X we considered two *relative* measures (of association):

- RR the relative risk (p_1/p_0)
- the OR odds ratio ($odds_1/odds_0$)

What were the corresponding measures in linear regression?

The linearity assumption between the probability $P(Y|X)$ and X is most often not appropriate.

We assume a logistic function to describe an S -shaped relation

$$\text{logistic}(x) = \frac{e^x}{1+e^x}$$

The empty model to estimate the odds or probability

$$P(Y) = \frac{\exp(\beta_0)}{1+\exp(\beta_0)}$$

$$\log(\text{odds}(Y)) = \beta_0$$

Binary predictor

Binary predictor

We want to relate the probability of the outcome Y to a binary predictor X , i.e. we want to test if the probability of the outcome is different in the two groups defined by X ($X = 1$ vs $X = 0$)

$$\text{logit}(P[Y|X]) = \log(\text{odds}(Y|X)) = \beta_0 + \beta_1 X$$

Or alternatively

$$P(Y|X) = \frac{\exp(\beta_0 + \beta_1 X)}{1 + \exp(\beta_0 + \beta_1 X)}$$

The model assumes that the **log odds** of the outcome *linearly* depends on X .

Or alternatively, that the probability of the outcome depends of X following a logistic function.

In general, all the modeling techniques apply to the linear model (log odds), while the results will be most often presented in terms of odds (or probability) and odds ratios (or risk ratio).

Interpretation

$$\log(\text{odds}(Y|X)) = \beta_0 + \beta_1 X$$

β_0 is the log odds of the outcome when $X = 0$.

What about β_1 ?

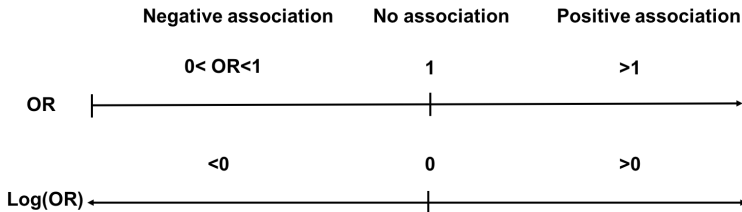
$$\begin{aligned} \beta_1 &= \log(\text{odds}(Y|X=1)) - \log(\text{odds}(Y|X=0)) = \\ &= \log\left(\frac{\text{odds}(Y=1|X=1)}{\text{odds}(Y=1|X=0)}\right) = \log(OR) \end{aligned}$$

β_1 is the the log odds ratio of the outcome comparing $X = 1$ vs. $X = 0$.

Interpretation of model coefficients is more informative on the exponential scale.

$\exp(\beta_0)$ is the odds of the outcome when $X = 0$.

$\exp(\beta_1)$ is the odds ratio of the outcome comparing $X = 1$ vs. $X = 0$.



Question: Is sex (female) a predictor of (the risk of) hyponatremia (nas135)?

Or, is the risk of hyponatremia different between men and women?

$$\log(\text{odds}(\text{nas135}|\text{female})) = \beta_0 + \beta_1 \text{female}$$

Estimation

We can use maximum likelihood to estimate the β coefficients.

$$Y \sim \text{Bernoulli}(p(x)), p(x) = P(Y|X = x).$$

The likelihood is defined as

$$\mathcal{L}(p|y, x) = \prod_{i=1}^n p(x_i)^{y_i} (1 - p(x_i))^{1-y_i}$$

or alternatively the log-likelihood

$$\ell(p|y, x) = \sum_{i=1}^n y_i \log(p(x_i)) + (1 - y_i) \log(1 - p(x_i))$$

$$p(x_i) = \frac{\exp(\beta_0 + \beta_1 x_i)}{1 + \exp(\beta_0 + \beta_1 x_i)}$$

so the log-likelihood becomes

$$\ell(\beta|y, x) = \sum_{i=1}^n y_i \log\left(\frac{\exp(\beta_0 + \beta_1 x_i)}{1 + \exp(\beta_0 + \beta_1 x_i)}\right) + (1 - y_i) \log\left(1 - \frac{\exp(\beta_0 + \beta_1 x_i)}{1 + \exp(\beta_0 + \beta_1 x_i)}\right)$$

The $\hat{\beta}_0$ and $\hat{\beta}_1$ which maximizes $\ell(\beta|y, x)$ can be obtained using iterative algorithms (there is no close formula).

Properties of MLE

The MLE \hat{B} has the following nice properties:

1. Consistency: the distributions of the estimators become more and more concentrated near the true value of the parameter being estimated.
2. Asymptotically normal: $\hat{B} \sim N(\beta, \text{Var}(\hat{B}))$
3. Asymptotic optimality: MLE has the smallest asymptotic variance and we say that the MLE is asymptotically efficient.
4. Invariance property: the maximum likelihood estimate of a function of the parameter being estimated ($\tau = g(\beta)$) is the function evaluated at the maximum likelihood estimate of the parameter ($\hat{\tau} = g(\hat{\beta})$)

Hypothesis testing

If sex is not associated with (risk of) hyponatremia, it means that $p_1 = p_0$, and thus $\text{odds}_1 = \text{odds}_0$. In this case, the OR will be 1.

The null hypothesis of no association between sex and (risk of) hyponatremia can be written as $H_0 : \beta_1 = 0$ (if $\beta_1 = \log(OR) = 0$, it means that $OR = 1$).

Based on large sample, the distribution of betas is approximately normal and Z test can be adopted.

$$Z = \frac{\hat{\beta}_1 - 0}{SE(\hat{\beta}_1)}$$

which, under the null hypothesis, follows a standard normal distribution.

Confidence intervals for the OR

95% confidence intervals are first defined on the log scale (i.e. for the $\log(OR)$)

$$\hat{\beta}_1 \pm 1.96 \cdot SE(\hat{\beta}_1)$$

and are exponentiated to obtain the corresponding confidence intervals for the OR (invariance property)

$$\exp\left(\hat{\beta}_1 \pm 1.96 \cdot SE(\hat{\beta}_1)\right)$$


```
mod <- glm(nas135 ~ female, data = marathon, family = "binomial")
summary(mod)
```

Call:

```
glm(formula = nas135 ~ female, family = "binomial", data = marathon)
```

Deviance Residuals:

Min	1Q	Median	3Q	Max
-0.7102	-0.7102	-0.4020	-0.4020	2.2608

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-2.4749	0.2082	-11.884	< 2e-16
femalefemale	1.2260	0.2795	4.386	0.0000116

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 371.60 on 487 degrees of freedom
Residual deviance: 351.93 on 486 degrees of freedom
AIC: 355.93

Number of Fisher Scoring iterations: 5

Interpretation

Alternative ways:

- the log odds of hyponatremia for male runners is -2.47 ;
- the odds of hyponatremia is 0.08 (8 cases for every 100 non-cases) among men;
- the log odds ratio of hyponatremia comparing female vs male runners is 1.23;
- the odds ratio of hyponatremia comparing female vs male runners is 3.41;
- the odds of hyponatremia among women is 3.41 times the odds for men;
- you multiply by 3.41 the odds of hyponatremia among male (0.08) to get the one among female (-1.24).

The 95% CI for the $\log(OR)$ is calculated as
 $1.23 \pm 1.96 \cdot 0.28 = (0.68, 1.77)$

The 95% CI for the OR is calculated as
 $\exp(0.68, 1.77) = (1.97, 5.87)$

The odds of hyponatremia among women was significantly higher than men ($OR = 3.41$).

We are 95% confident that the odds ratio relating sex (being woman compared to man) to hyponatremia is between 1.97 and 5.87.

The multiplicative model

The logistic model is a linear model in terms of the log odds

$$\log(\text{odds}(Y|X)) = \beta_0 + \beta_1$$

If we take the exponential, we can write the previous model as a multiplicative model in terms of the odds

$$\begin{aligned} \exp(\log(\text{odds}(Y|X))) &= \text{odds}(Y|X) = \exp(\beta_0 + \beta_1) = \\ \exp(\beta_0) \exp(\beta_1) &= \text{odds}(Y|X = 0) \cdot OR_{x=1 \text{ vs } x=0} \end{aligned}$$

What is the odds of hyponatremia among men? $\exp(\hat{\beta}_0) = 0.08$

What is the odds of hyponatremia among women?

$$\exp(\hat{\beta}_0) \cdot \exp(\hat{\beta}_1) = 0.08 \cdot 3.41 = 0.27$$

Linear function of regression coefficients

Question: What is the odds of hyponatremia among women?

Coefficients:

	b0	b1
	-2.475	1.226

Covariance matrix:

	b0	b1
b0	0.043	-0.043
b1	-0.043	0.078

Let's first calculate the 95% CI for the *log odds*

$$\widehat{\text{est}} = \log(\text{odds}(Y|X = 1)) = \hat{\beta}_0 + \hat{\beta}_1 = -2.47 + (1.23) = -1.249$$

$$\begin{aligned} \text{Var}(\widehat{\text{est}}) &= \text{Var}(\hat{\beta}_0) + \text{Var}(\hat{\beta}_1) + 2\text{Cov}(\hat{\beta}_0, \hat{\beta}_1) = \\ &0.043 + 0.078 + 2(-0.043) = 0.035 \end{aligned}$$

95% confidence interval for $\widehat{\text{est}}$ is

$$\widehat{\text{est}} \pm 1.96\sqrt{\text{Var}(\widehat{\text{est}})} = (-1.616, -0.882)$$

Now we can exponentiate both the point estimate and the confidence interval:

$$\exp(\hat{\beta}_0 + \hat{\beta}_1) = 0.29$$

95% confidence interval for odds($Y|X = 1$) is

$$\exp(-1.616, -0.882) = (0.199, 0.414)$$

The odds of hyponatremia is 0.29 (29 cases for every 100 non-cases) among women (95% CI: 0.199, 0.414).

Predicted probabilities

Once the β coefficients have been estimated it is possible to calculate the predicted probabilities of the outcome for any covariate values (covariate pattern).

What is the estimated probability of hyponatremia among male female?

$$P(Y = 1|X = 0) = \text{invlogit}(\beta_0) = \frac{\exp(-2.47)}{1+\exp(-2.47)} = \frac{0.08}{1+0.08} = 0.078$$

$$P(Y = 1|X = 1) = \text{invlogit}(\beta_0 + \beta_1) = \frac{\exp(-2.47+1.23)}{1+\exp(-2.47+1.23)} = \frac{0.29}{1+0.29} = 0.224$$

The 95% CI for the predicted probabilities can be obtained as the `invlogit` of the 95% CI for the corresponding *log odds*

95% CI for $\log(\text{odds}(Y|X = 0))$ ($\rightarrow \hat{\beta}_0$)

$$\hat{\beta}_0 \pm z_{0.975} \text{SE}(\hat{\beta}_0) = -2.47 \pm 1.96\sqrt{0.04} = (-2.86, -2.08)$$

95% CI for $P(Y|X = 0)$

$$\left(\frac{\exp(-2.86)}{1 + \exp(-2.86)}, \frac{\exp(-2.08)}{1 + \exp(-2.08)} \right) = (0.05, 0.11)$$

	b0	b1
b0	0.04	-0.04
b1	-0.04	0.08

95% CI for $\log(\text{odds}(Y|X = 1))$ ($\rightarrow \hat{\beta}_0 + \hat{\beta}_1$)

$$\text{Var}(\beta_0 + \beta_1) = \text{Var}(\hat{\beta}_0) + \text{Var}(\hat{\beta}_1) + 2\text{Cov}(\hat{\beta}_0, \hat{\beta}_1) =$$

$$0.04 + 0.08 + 2(-0.04) = 0.04$$

$$(\hat{\beta}_0 + \hat{\beta}_1) \pm z_{0.975}\text{SE}(\hat{\beta}_0 + \hat{\beta}_1) = -1.24 \pm 1.96\sqrt{0.04} = (-1.63, -0.85)$$

95% CI for $P(Y|X = 1)$

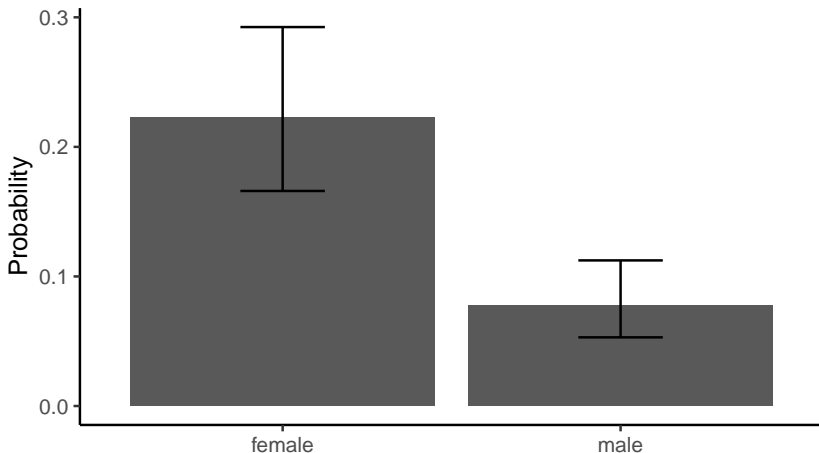
$$\left(\frac{\exp(-1.63)}{1+\exp(-1.63)}, \frac{\exp(-0.85)}{1+\exp(-0.85)} \right) = (0.16, 0.3)$$

Predicted response

X	odds	p	OR	RR
0	$\exp(\beta_0)$	$\text{invlogit}(\beta_0)$	ref	ref
1	$\exp(\beta_0 + \beta_1)$	$\text{invlogit}(\beta_0 + \beta_1)$	$\exp(\beta_1)$	p_1/p_0

X	odds	p	OR	RR
male	0.08	0.078	ref	ref
female	0.27	0.224	3.41	2.878

Graphical presentation of predicted probabilities



The logistic regression model estimated above with a binary covariate (female) is called saturated because the number of possible combination of covariate patterns (male, female) is equal to the number of parameters estimated.

<hr/>		
cases	non-cases	female
<hr/>		
25	297	male
37	129	female
<hr/>		

The consequence is that the fitted values from the saturated model will exactly fit the observed data. No model is more complicated than that.

cases	non-cases	female	p obs	p pred
25	297	male	0.0776398	0.0776398
37	129	female	0.2228916	0.2228916

We want know to related the probability of the outcome Y to a continuous predictor X .

Question: is weight change (`wtdiff`) a predictor of the (risk of) hyponatremia (`nas135`)?

Or alternatively, does the risk oh hyponatremia vary as a function of weight change? if yes, how does it vary?

$$\log(\text{odds}(\text{nas135}|\text{wtdiff})) = \beta_0 + \beta_1 \text{wtdiff}$$

NB: we assume that the log odds of hyponatremia *linearly* varies as a function of weight change.

More in general, the log odds of the outcome for any two values of X (x_1 vs x_2) are

$$\log(\text{odds}(Y|X = x_1)) = \beta_0 + \beta_1 x_1$$

$$\log(\text{odds}(Y|X = x_2)) = \beta_0 + \beta_1 x_2$$

$$\log(\text{odds}(Y|X = x_1)) - \log(\text{odds}(Y|X = x_2)) = \log(OR_{x_1 \text{ vs } x_2}) = \beta_1(x_1 - x_2)$$

is the log OR associated with a $(x_1 - x_2)$ change in X .

The odds ratio can be obtained exponentiating the log OR

$$OR = \exp(\beta_1(x_1 - x_2))$$

That is the OR comparing the sub-population having x_1 vs. x_2 of the quantitative covariate X .

Confidence intervals for OR

$$\widehat{\text{est}} = \log(OR) = \beta_1(x_1 - x_2)$$

$$\text{Var}(\widehat{\text{est}}) = \text{Var}(\beta_1(x_1 - x_2)) = (x_1 - x_2)^2 \text{Var}(\beta_1)$$

95% CI for **log** odds ratio

$$\widehat{\text{est}} \pm z_{0.975} \sqrt{\text{Var}(\widehat{\text{est}})}$$

95% CI for odds ratio (invariance property)

$$\exp\left(\widehat{\text{est}} \pm z_{0.975} \sqrt{\text{Var}(\widehat{\text{est}})}\right)$$

$$\widehat{\text{est}} = \hat{\beta}_1(x_1 - x_2) = 0.73(3) = 2.19$$

$$\text{Var}(\widehat{\text{est}}) = \text{Var}(3\hat{\beta}_1) = 9 \cdot \text{Var}(\hat{\beta}_1) = 9 \cdot 0.012 = 0.11$$

$$\text{SE}(\widehat{\text{est}}) = \sqrt{\text{Var}(\widehat{\text{est}})} = \sqrt{(0.11)} = 0.329$$

95% CI for **log** OR

$$2.19 \pm z_{0.975} \times 0.329 = (1.55, 2.83)$$

$$OR = \exp(3\hat{\beta}_1) = 8.94$$

95% CI for OR

$$\exp(1.55, 2.83) = (4.69, 17.03)$$

The odds of hyponatremia among those who increased 2 kg was 9 (95% CI 4.69, 17.03) times the odds for those runners who lost 1kg.

