

EXERCISE 2

$$y_i = \text{diabetes}_i = \begin{cases} 1 & \text{if woman } i \text{ has diabetes} \\ 0 & \text{if woman } i \text{ does not have diabetes} \end{cases}$$

$$i = 1, \dots, n \quad \text{with} \quad n = 724$$

$$\text{Moreover} \quad \sum_{i=1}^n y_i = 249 \quad \sum_{i=1}^n (1 - y_i) = 475$$

- a) • $Y_i \sim \text{Bernoulli}(\pi_i) \quad i = 1, \dots, n \quad \text{independent} \quad (\text{distribution})$
- $\eta_i = \beta_1 + \beta_2 x_{i2} + \beta_3 x_{i3} + \beta_4 x_{i4} + \beta_5 x_{i5} + \beta_6 x_{i6}$
- with $x_{i2} = \text{pregnant}_i \quad (\text{dummy}) \quad (\text{linear predictor})$
- $x_{i3} = \text{glucose}_i$
- $x_{i4} = \text{pressure}_i$
- $x_{i5} = \text{BMI}_i$
- $x_{i6} = \text{age}_i$
- $g(\pi_i) = \text{logit}(\pi_i) = \log \frac{\pi_i}{1 - \pi_i} = \eta_i \quad (\text{link function})$

b) • Interpretation of the "age" coefficient (β_6)

We study the odds of having diabetes for two women i and k

with $\underset{x_{k6}}{\text{age}_k} = \underset{x_{i6}}{\text{age}_i} + 1$ and all other covariates equal.

$$\text{odds woman } i: \frac{\pi_i}{1 - \pi_i} = \frac{e^{\eta_i^T \beta}}{1 + e^{\eta_i^T \beta}} \cdot (1 + e^{\eta_i^T \beta}) = e^{\eta_i^T \beta}$$

$$= \exp\{\beta_1 + \beta_2 x_{i2} + \beta_3 x_{i3} + \beta_4 x_{i4} + \beta_5 x_{i5} + \beta_6 x_{i6}\}$$

$$\text{odds woman } k: \frac{\pi_k}{1 - \pi_k} = e^{\eta_k^T \beta} = \exp\{\beta_1 + \beta_2 \underset{x_{i2}}{x_{k2}} + \beta_3 \underset{x_{i3}}{x_{k3}} + \beta_4 \underset{x_{i4}}{x_{k4}} + \beta_5 \underset{x_{i5}}{x_{k5}} + \beta_6 \underset{x_{i6}+1}{x_{k6}}\}$$

if we look at the ODDS RATIO

$$\frac{\frac{\pi_k}{1 - \pi_k}}{\frac{\pi_i}{1 - \pi_i}} = e^{\beta_6} \Rightarrow \frac{\pi_k}{1 - \pi_k} = e^{\beta_6} \cdot \frac{\pi_i}{1 - \pi_i} \quad \text{and} \quad e^{\hat{\beta}_6} = 1.033 > 1$$

odds with age increased by 1

The odds of having diabetes increase by a multiplicative factor 1.033 for every additional year, keeping all other covariates fixed.

• Interpretation of the "pregnant" coefficient (β_2)

We study the odds of two women i and k , where woman i had no pregnancies, and woman k had at least one pregnancy, and all other covariates are equal.

$$\text{odds woman } i: \frac{\pi_i}{1 - \pi_i} = \frac{e^{\eta_i^T \beta}}{1 + e^{\eta_i^T \beta}} \cdot (1 + e^{\eta_i^T \beta}) = e^{\eta_i^T \beta}$$

$$= \exp\{\beta_1 + \beta_2 \cancel{x_{i2}} + \beta_3 x_{i3} + \beta_4 x_{i4} + \beta_5 x_{i5} + \beta_6 x_{i6}\}$$

$$\text{odds woman } k: \frac{\pi_k}{1 - \pi_k} = e^{\eta_k^T \beta} = \exp\{\beta_1 + \beta_2 \underset{1}{x_{k2}} + \beta_3 \underset{x_{i3}}{x_{k3}} + \beta_4 \underset{x_{i4}}{x_{k4}} + \beta_5 \underset{x_{i5}}{x_{k5}} + \beta_6 \underset{x_{i6}}{x_{k6}}\}$$

if we look at the ODDS RATIO

$$\frac{\frac{\pi_k}{1 - \pi_k}}{\frac{\pi_i}{1 - \pi_i}} = e^{\beta_2} \Rightarrow \frac{\pi_k}{1 - \pi_k} = e^{\beta_2} \cdot \frac{\pi_i}{1 - \pi_i} \quad \text{and} \quad e^{\hat{\beta}_2} = 1.279 > 1$$

odds with at least 1 pregnancy

The odds of having diabetes increase by a multiplicative factor 1.279 when we consider a woman who had at least one pregnancy, compared to a woman with no pregnancies and all other covariates equal.

c) If we look at the p-value in the Table, it is 0.40.

It is the p-value α^{obs} of the test for testing

$$\begin{cases} H_0: \beta_2 = 0 \\ H_1: \beta_2 \neq 0 \end{cases}$$

Hence $0.40 = P_{H_0}(|Z| > |z^{\text{obs}}|)$ probability of observing "more extreme" values.

Since $\alpha^{\text{obs}} > \alpha$ for all usual α (0.01, 0.05, 0.1), we do not reject H_0 .

That is, β_2 is not significant and we can remove the corresponding covariate.

d) We want to compare model A and model B.

First of all, notice that model B is nested w.r.t. model A and can be obtained by constraining $\beta_2 = \beta_3 = 0$.

Hence, we want to test

$$\begin{cases} H_0: \beta_2 = \beta_3 = 0 \\ H_1: \exists r \in \{2, 3\}: \beta_r \neq 0 \end{cases}$$

We use the likelihood ratio test

$$W = 2(\hat{\ell}(\text{model A}) - \hat{\ell}(\text{model B})) \stackrel{H_0}{\sim} \chi^2_{p-p_0} \quad \text{Here } p = 6$$

$$p_0 = 4$$

max of the loglikelihood in the FULL model

max of the loglikelihood in the RESTRICTED model

observed value of the test

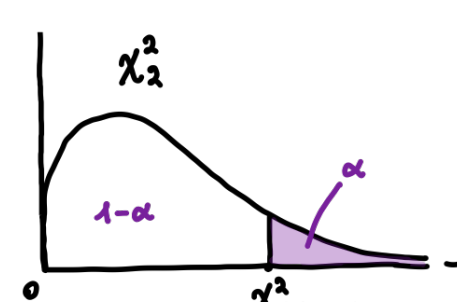
$$\begin{aligned} w^{\text{obs}} &= 2(\hat{\ell}(\text{model A}) - \hat{\ell}(\text{model B})) \\ &= 2(\hat{\ell}(\text{model A}) - \hat{\ell}(\text{model B}) + \tilde{\ell}(\text{saturated}) - \tilde{\ell}(\text{saturated})) \\ &= 2[(\tilde{\ell}(\text{saturated}) - \hat{\ell}(\text{model B})) - (\tilde{\ell}(\text{saturated}) - \hat{\ell}(\text{model A}))] \\ &= D(\text{model B}) - D(\text{model A}) \\ &= 696.15 - 694.45 = 1.70 \end{aligned}$$

Reject region is for large values of the test using a fixed significance level α

$$R = (\chi^2_{2; 1-\alpha}; +\infty)$$

$$\text{if } \alpha = 0.05 \quad R = (\chi^2_{2; 0.95}; +\infty) = (5.99; +\infty)$$

$w^{\text{obs}} \notin R$: do not reject H_0
| prefer model B



e) NULL model

is the model with only the intercept: it assumes that none of the covariates has an effect.

$$\begin{cases} Y_i \sim \text{Bern}(\pi_i) \\ \eta_i = \beta_1 \\ \text{logit } \pi_i = \eta_i \end{cases}$$

π_i does not depend on i : it is constant

$$\pi_i = \pi = \frac{e^{\beta_1}}{1 + e^{\beta_1}} \quad \text{for all } i = 1, \dots, n$$

In the null model, the MLE of π is $\hat{\pi} = \bar{y} = \frac{249}{724} = 0.344$

$$\text{Thus, } \hat{\beta}_1 = \log \frac{\bar{y}}{1 - \bar{y}} = -0.646$$