

# Assignment II

## Logistic regression

Silpa Soni Nallacheruvu (19980824-5287) Hernan Aldana (20000526-4999)

2024-11-28

---

### Exercise 2:1

This report analyses the data on periodontitis from a group of adult patients at a large dental clinic. The goal is to understand the influence of the parameter estimates between logistic regression models of the probability for periodontitis in the population as a function of dental floss use and the probability for using dental floss as a function of periodontitis status.

Table 1: Regular Use of Dental Floss and Periodontitis

	Periodontitis	No Periodontitis	Total
Used Dental Floss	22	75	148
Not Used Dental Floss	265	97	413
Total	170	340	510

### Question 1: Logistic regression model of the probability for periodontitis

The logistic regression model of the probability for periodontitis in the population as a function of dental floss use, say Model A is defined as:

$$\text{logit}(p_x) = \log\left(\frac{p_x}{1-p_x}\right) = \beta_0 + \beta_1 x$$

where:

- $x = 1$  if dental floss is regularly used,  $x = 0$  otherwise.
- $p_x = P(\text{periodontitis} | x)$ , the probability of periodontitis given  $x$ .
- $\beta_0$ : The log-odds of periodontitis when  $x = 0$  (no floss use).
- $\beta_1$ : The change in log-odds of periodontitis when  $x$  changes from 0 to 1 (effect of using floss).

#### Approach:

Fit the logistic regression model A to the data and interpret the parameter estimates  $\beta_0$  and  $\beta_1$ .

#### Estimate Parameters:

#### Fit the logistic regression model:

The summary of the coefficients of the logistic regression model A is presented below:

Table 2: Summary of Logistic Regression Model A Coefficients

	Term	Estimate	Std. Error	z value	P-value
(Intercept)	(Intercept)	-0.5825176	0.1026175	-5.676593	0.0000000
floss	floss	-0.6439281	0.2632835	-2.445759	0.0144548

The final logistic regression equation:  $\text{logit}(p_x) = -0.582 - 0.644x$

#### Interpret Parameters:

- 1.  $\beta_0 = -0.582$ :
  - The log-odds of periodontitis for individuals who do not use dental floss is approximately -0.582.
  - The corresponding probability of periodontitis is:  $p_0 = \frac{e^{-0.582}}{1+e^{-0.582}} \approx 0.358$
- 2.  $\beta_1 = -0.644$ :
  - The log-odds of periodontitis decreases by 0.644 when individuals use dental floss regularly.
  - This corresponds to an odds ratio of: Odds Ratio =  $e^{\beta_1} = e^{-0.644} \approx 0.525$

Individuals who use dental floss regularly have approximately 52.5% lower odds of developing periodontitis compared to those who do not.

#### Conclusion:

- The logistic regression model A provides interpretable estimates of the effect of dental floss use on periodontitis.
- Regular dental floss use is associated with a significant reduction in the odds of periodontitis.
- The logistic regression model A is appropriate because the data is binary (periodontitis: yes/no) and the explanatory variable (dental floss use) is categorical.

### Question 2: Logistic regression model of the probability for using dental floss

The logistic regression model of the probability for using dental floss as a function of periodontitis status, say Model B is defined as:

$$\text{logit}(p_y) = \log\left(\frac{p_y}{1-p_y}\right) = \gamma_0 + \gamma_1 y$$

where:

- $y = 1$  if periodontitis is present,  $y = 0$  otherwise.
- $p_y = P(\text{using dental floss} \mid y)$ , the probability of using dental floss given  $y$ .
- $\gamma_0$ : The log-odds of using dental floss when  $y = 0$  (i.e., when periodontitis is absent).
- $\gamma_1$ : The change in log-odds of using dental floss associated with the presence of periodontitis ( $y = 1$ ).

#### Approach:

Fit the logistic regression model B to the data and interpret the parameter estimates  $\gamma_0$  and  $\gamma_1$

#### Estimate Parameters:

##### Fit the logistic regression model:

The summary of the coefficients of the logistic regression model B is presented below:

Table 3: Summary of Logistic Regression Model B Coefficients

	Term	Estimate	Std. Error	z value	P-value
(Intercept)	(Intercept)	-1.2622417	0.1307934	-9.650652	0.0000000
periodontitis	periodontitis	-0.6439281	0.2632835	-2.445759	0.0144548

The final logistic regression equation:  $\text{logit}(p_y) = -1.262 - 0.644y$

#### Interpret Parameters:

- 1.  $\gamma_0 = -1.262$ :
  - The log-odds of using dental floss when periodontitis is absent is approximately -1.262
  - The corresponding probability of periodontitis is:  $p_0 = \frac{e^{-1.262}}{1+e^{-1.262}} \approx 0.2205$
  - This means approximately 22% of people without periodontitis use dental floss.
- 2.  $\gamma_1 = -0.644$ :
  - The log-odds of using dental floss decreases by 0.644 when periodontitis is present.
  - This corresponds to an odds ratio of: Odds Ratio =  $e^{\gamma_1} = e^{-0.644} \approx 0.525$

The odds of individuals with periodontitis using dental floss are approximately 52.4% lower compared to not using dental floss.

#### Conclusion:

- The logistic regression model provides interpretable estimates of the effect of dental floss use on periodontitis.
- The presence of periodontitis is associated with a significant reduction in the odds of regular dental floss.
- The logistic regression model B is appropriate because the data is binary (Dental Floss use: yes/no) and the explanatory variable (presence of periodontitis) is categorical.

#### Question 3:

##### Approach :

##### 1. Expected Estimates of the parameters of Model A:

###### 1.1. Intercept ( $\beta_0$ ):

- $\beta_0$  represents the log-odds of periodontitis when  $x = 0$  (no dental floss use):  $\beta_0 = \log\left(\frac{p_0}{1-p_0}\right)$  where:  
 $p_0 = \frac{\text{Periodontitis (No Floss)}}{\text{Total (No Floss)}} = \frac{148}{413} \approx 0.3585$
- Hence:  $\beta_0 \approx \log\left(\frac{0.3585}{1-0.3585}\right) = \log(0.558) \approx -0.582$

###### 1.2. Slope ( $\beta_1$ ):

- $\beta_1$  represents the change in log-odds when  $x = 1$  (dental floss is used):  $\beta_1 = \log\left(\frac{p_1}{1-p_1}\right) - \log\left(\frac{p_0}{1-p_0}\right)$   
 where:  $p_1 = \frac{\text{Periodontitis (Floss)}}{\text{Total (Floss)}} = \frac{22}{97} \approx 0.2268$
- Hence:  $\beta_1 \approx \log\left(\frac{0.2268}{1-0.2268}\right) - \log\left(\frac{0.3585}{1-0.3585}\right) \approx \log(0.293) - \log(0.558) = -1.226 + 0.582 = -0.644$

##### 2. Expected Estimates of the parameters of Model B:

###### 2.1. Intercept ( $\gamma_0$ ):

- $\gamma_0$  represents the log-odds of using dental floss when  $y = 0$  (periodontitis is absent):  $\gamma_0 = \log\left(\frac{p_0}{1-p_0}\right)$   
where:  $p_0 = \frac{\text{Using Floss (No Periodontitis)}}{\text{Total (No Periodontitis)}} = \frac{75}{340} \approx 0.2205$
- Hence:  $\gamma_0 \approx \log\left(\frac{0.2205}{1-0.2205}\right) = \log(0.282) \approx -1.262$

2.2. Slope ( $\gamma_1$ ):

- $\gamma_1$  represents the change in log-odds of using dental floss when  $y = 1$  (periodontitis is present):  
 $\gamma_1 = \log\left(\frac{p_1}{1-p_1}\right) - \log\left(\frac{p_0}{1-p_0}\right)$  where:  $p_1 = \frac{\text{Using Floss (Periodontitis)}}{\text{Total (Periodontitis)}} = \frac{22}{170} \approx 0.1294$
- Hence:  $\gamma_1 \approx \log\left(\frac{0.1294}{1-0.1294}\right) - \log\left(\frac{0.2205}{1-0.2205}\right) \approx \log(0.149) - \log(0.282) = -1.906 + 1.262 = -0.644$

### Observation :

- The expected estimates of the parameters of Model A and Model B are consistent with the actual estimates obtained from the logistic regression models.
  - The intercepts  $\beta_1$  and  $\gamma_1$  are the same, indicating that the log-odds of periodontitis and dental floss use are similar when the other variable is absent.
-

## Exercise 2:2

### Question 1:

For this question we estimate for each dose separately the risk, odds, and log-odds of developing a tumor. Also plot the risk (probability) of developing a tumor as a function of the log dose. Calculate the log-odds of developing a tumor and plot it as a function of the log dose.

Table 4: Original Data Table

log.dose.	-7.6	-6.22	-4.6	-3	-1.39	0.92
Tumor	1.0	2.00	4.0	9	12.00	32.00
No.tumor	17.0	17.00	24.0	23	16.00	8.00
Total	18.0	19.00	28.0	32	28.00	40.00

### Approach :

To answer these questions we have to calculate the following estimates for each dose level of the table:

1. The risk of developing a tumor is calculated as:

$$Risk = \frac{NumberofTumor}{TotalObservations}$$

2. The odds of developing a tumor is calculated as:

$$Odds = \frac{Risk}{1-Risk}$$

3. The log-odds of developing a tumor is calculated as:

$$\log(Odds) = \log\left(\frac{Risk}{1-Risk}\right)$$

After we calculate these estimates for each dose level, we plot the risk (probability) of developing a tumor as a function of the log dose. We also calculate the log-odds of developing a tumor and plot it as a function of the log dose.

### Results

Table 5: Original Data Table with Risk, Odds, and Log-Odds Calculations

log.dose.	-7.6000	-6.2200	-4.6000	-3.0000	-1.3900	0.9200
Tumor	1.0000	2.0000	4.0000	9.0000	12.0000	32.0000
No.tumor	17.0000	17.0000	24.0000	23.0000	16.0000	8.0000
Total	18.0000	19.0000	28.0000	32.0000	28.0000	40.0000
Risk	0.0556	0.1053	0.1429	0.2812	0.4286	0.8000
Odds	0.0588	0.1176	0.1667	0.3913	0.7500	4.0000
Log.Odds	-2.8332	-2.1401	-1.7918	-0.9383	-0.2877	1.3863

#### 1. Data

- As we can observe from the updated table and from the given data, at lower doses of BaP, few mice develop tumors, while most do not. As the dose increases, the risk of developing a tumor also increases.

#### 2. Risk

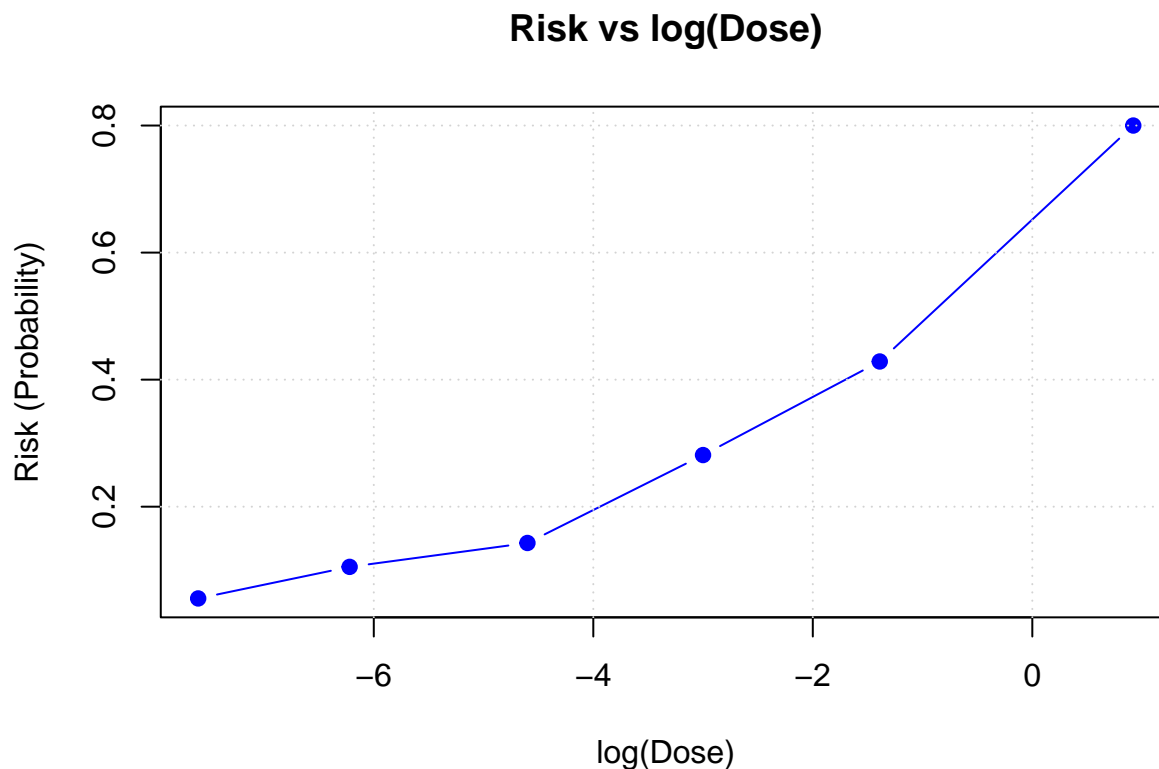
- Risk represents the probability of developing a tumor, and for the given data it increases with  $\log(dose)$ . As for the calculated estimates, the risk is very low(0.0556) for the lowest dose of BaP(-7.60), indicating a small proportion of tumors. At the highest dose (0.92) the risk is much higher(0.8), showing a significant increase in the development of tumors.

### 3. Odds

- Odds, which represent the ratio of tumor probability to no-tumor probability, increases exponentially with dose. The odds are very low(0.0588) for the lowest dose of BaP(-7.60), indicating a very low probability of developing a tumor. This means that tumors are 0.0588 as likely as no tumors. At the highest dose (0.92) the odds are much higher(4), showing a significant increase in the development of tumors. Meaning that tumors are four times more likely than no tumors at the given dose.

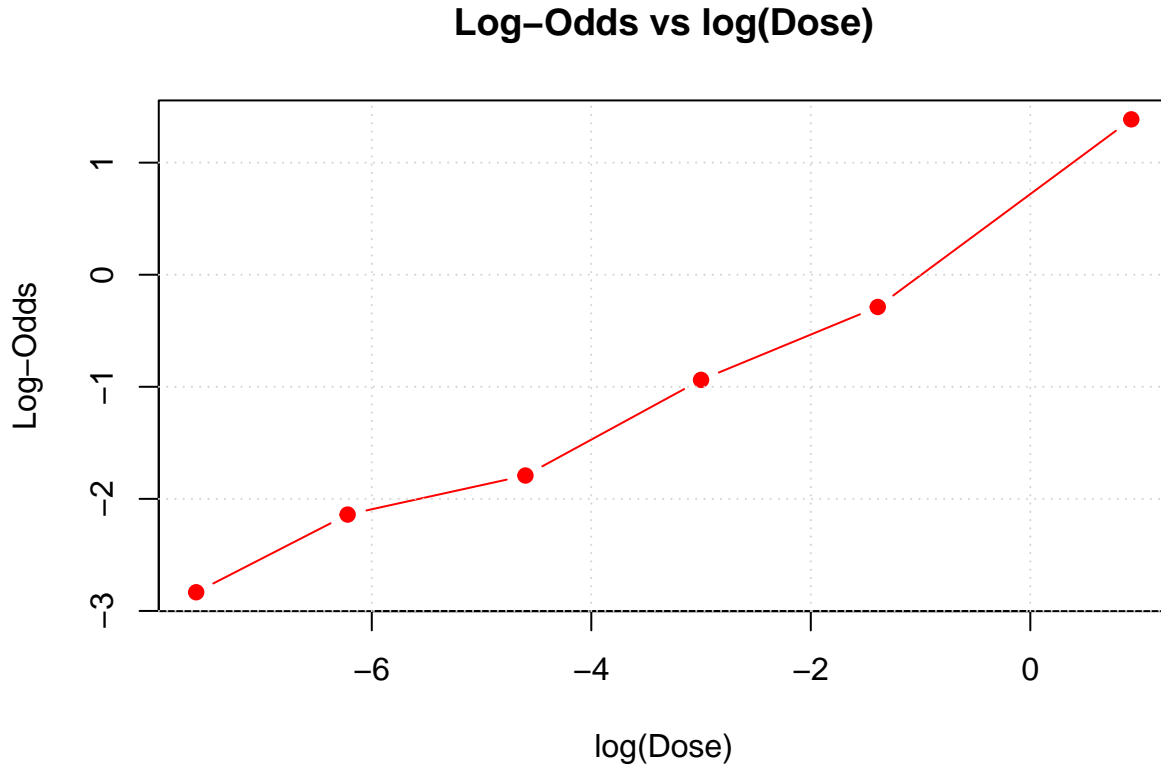
### 4. Log-Odds

- Log-Odds increase nearly linearly with  $\log(\text{dose})$ , ranging from -2.8332 at the lowest dose to 1.3863 at the highest. The linearity in log-odds is imperative, as it supports the use of a logistic regression.



### 1. Risk vs $\log(\text{dose})$

- The graph confirms a strong dose response relationship. As  $\log(\text{dose})$  increases, the probability of developing a tumor increases rapidly.



#### 2. Log-Odds vs log(dose)

- The graph shows a nearly linear relationship with log(dose) confirming that the data aligns well with the assumptions of logistic regression.

### Conclusion

The analysis confirms a strong dose response relationship between the dose of BaP and the probability of developing a tumor. As the dose or log(dose) increases, the risk and odds of tumor development also increase as shown in the “Risk vs. log(dose)” plot. The near linear relationship in the “Log-Odds vs. log(dose)” plot supports the use of logistic regression to model the data.

### Question 2

Fit a logistic regression model to the data and interpret the parameter estimates, particularly the slope parameter.

#### Approach :

To answer this question, we fit a logistic regression model where the probability of developing a tumor( $P$ ) is modeled as:

$$\text{logit}(P) = \log\left(\frac{P}{1-P}\right) = \beta_0 + \beta_1 \times \log(\text{dose})$$

where:

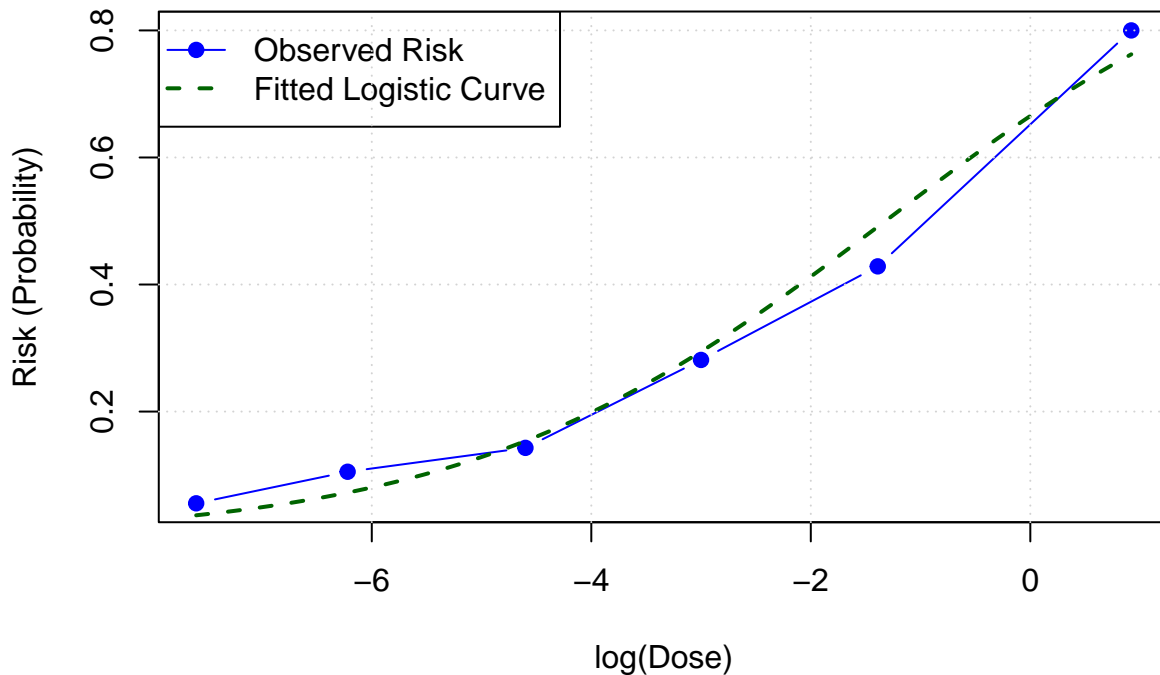
- $\beta_0$  is the intercept, which represents the log-odds of developing a tumor when the dose is zero.
- $\beta_1$  is the slope parameter, which represents the change in log-odds of developing a tumor for a one-unit increase in the log dose.

We will use the logistic regression model to estimate this parameters.

Results :

```
##
## Call:
## glm(formula = cbind(tumor, no_tumor) ~ log_dose, family = binomial)
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  0.68702    0.25755   2.668  0.00764 **
## log_dose      0.52037    0.08502   6.120  9.35e-10 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##    Null deviance: 56.5321  on 5  degrees of freedom
## Residual deviance:  1.2416  on 4  degrees of freedom
## AIC: 24.024
##
## Number of Fisher Scoring iterations: 4
```

### Risk vs log(Dose) with Logistic Fit



Intercept  $\beta_0$

According to the summary of the model, the intercept  $\beta_0$  is 0.687. This represents the log-odds of developing a tumor when the dose is zero. Its corresponding odds are  $e^{0.687} = 1.99$ . This means that the odds of developing a tumor are nearly twice as high as the odds of not developing a tumor when the log(dose) is zero.

Translating this to probabilities, the probability of developing a tumor is  $P = \frac{e^{0.687}}{1 + e^{0.687}} = 0.665$ . This means that the probability of developing a tumor is 66.5% when the dose is zero.

2. Slope  $\beta_1$

The slope parameter  $\beta_1$  is 0.52037. This quantifies the change in the log-odds of developing a tumor for a



one-unit increase in  $\log(\text{dose})$ . A positive slope indicates that higher doses increase the odds of developing a tumor. the odds ratio is  $e^{0.52037} \approx 1.683$

This means that for each one-unit increase in  $\log(\text{dose})$ , the odds of developing a tumor increases by approximately 68.3%

### Conclusion :

In conclusion, the logistic regression model reveals a clear relationship between  $\log(\text{dose})$  and the probability of developing a tumor. The intercept  $\beta_0$  represents the log-odds of developing a tumor when the  $\log(\text{dose})$  is zero, while the slope  $\beta_1$  quantifies the change in log-odds for a one-unit increase in  $\log(\text{dose})$ . The positive slope indicates that higher doses increase the odds of developing a tumor. The fitted logistic curve aligns well with the observed risk, confirming the appropriateness of the logistic regression model.

### Question 3

Find the covariance matrix for the estimates and assess if they are correlated. Find a 95% confidence interval for the parameters. Find a 95% confidence interval for the tumor risk at dose 0.25( $\log(\text{dose}) = -1.39$ )

### Approach :

1. Covariance matrix for the estimates.
  - The covariance matrix provides the variances and covariances of the parameter estimates  $\beta_0, \beta_1$ . The variances are the diagonal elements of the matrix, while the covariances are the off-diagonal elements. The correlation between the estimates can be calculated from their covariance:

$$\text{Corr}(\beta_0, \beta_1) = \frac{\text{Cov}(\beta_0, \beta_1)}{\sqrt{\text{Var}(\beta_0) \times \text{Var}(\beta_1)}}$$

2. 95% confidence interval for the parameters.
  - The 95% confidence interval for the parameters can be calculated using the standard errors of the estimates. The confidence interval is given by:

$$\text{CI} = \hat{\beta} \pm Z \cdot \text{SE}(\hat{\beta})$$

where  $\hat{\beta}$  is the estimate,  $Z$  is the critical value for a 95% confidence interval, and  $\text{SE}(\hat{\beta})$  is the standard error of the estimate.

3. 95% confidence interval for the tumor risk at dose 0.25( $\log(\text{dose}) = -1.39$ )
  - Using the logistic regression equation:

$$\text{logit}(P) = \beta_0 + \beta_1 \times \log(\text{dose})$$

we calculate the predicted probability  $P$  at  $\log(\text{dose}) = -1.39$

- To find the confidence interval for  $P$ , we:
  - Calculate the standard error of  $\text{logit}(P)$ .
  - Use it to derive the confidence interval and transform the interval back to the probability scale.

### Results :

```
^
##               (Intercept)      log_dose
## (Intercept)  0.06633368  0.014358467
## log_dose     0.01435847  0.007229089
## [1] 0.655691
```

```
## Waiting for profiling to be done...
```

```
##           2.5 %    97.5 %  
## (Intercept) 0.1980903 1.2132416  
## log_dose    0.3638310 0.6989154
```

```
## $predicted_risk  
## (Intercept)  
##    0.4909301  
##  
## $confidence_interval  
## (Intercept) (Intercept)  
##    0.3940875    0.5884584
```

## 1. Covariance Matrix and Correlation

### I. Covariance Matriz

- The covariance matrix for the estimates  $\beta_0$  and  $\beta_1$  is:
  - Variance of  $\beta_0$ : 0.06633.
  - Variance of  $\beta_1$ : 0.01436.
  - Covariance between  $\beta_0$  and  $\beta_1$ : 0.00723

The correlation between  $\beta_0$  and  $\beta_1$  is -0.98, indicating a strong negative correlation between the intercept and slope parameters.

### II. Correlation:

- The correlation between  $\beta_0$  and  $\beta_1$  is 0.655 approximately, indicating a moderate positive correlation between the intercept and slope parameters.

## 2. Confidence Intervals for Parameters

- Confidence Interval for  $\beta_0$ :
  - The 95% confidence interval for  $\beta_0$  is approximately [0.198, 1.213].
  - This means that the log-odds of developing a tumor at  $\log(\text{dose})=0$  is likely within this range with 95% confidence.
- Confidence Interval for  $\beta_1$ :
  - The 95% confidence interval for  $\beta_1$  is approximately [0.363, 0.699].
  - This means that each one unit increase in  $\log(\text{dose})$  increases the log-odds of developing a tumor by approximately 0.363 to 0.699 with 95% confidence.

## 3. Predict Risk and Confidence Interval for P at a dose of 0.25 or $\log(\text{dose}) = -1.39$

- The predicted risk of developing a tumor at  $\log(\text{dose}) = -1.39$  is  $P = 0.409$  or 40.9%
- The 95% confidence interval for the predicted risk is approximately [0.394, 0.588].

This means that at a dose of 0.25, the probability of developing a tumor is estimated to be between 39.4% and 58.8% with 95% confidence.

**Conclusion :**

**Question 4**

**Approach :**

**Results :**

**conclusion :**