

# MA50259: Statistical Design of Investigations

Coursework 2 (2024)

09618

## Disclaimer:

AI software (RStudio built-in Co-Pilot and ChatGPT) are used in this coursework for code, RMarkdown formatting, explanation suggestions, grammar check, and debugging purposes.

## Part 1: Barley Experiment

The data in table below show the yields (measured in bushels per acre) of five varieties of barley in an experiment carried out in a rural area in the US.

Place	Year	Excel	Compana	Drummond	Conlon	Kindred
1	1971	71.0	95.4	109.7	119.5	88.3
1	1972	90.7	102.3	79.4	86.2	80.2
2	1971	132.6	121.0	140.7	171.5	135.7
2	1972	99.4	105.5	120.2	157.7	102.1
3	1971	87.3	73.3	79.4	128.6	85.3
3	1972	102.5	111.4	100.5	131.9	122.6
4	1971	129.8	111.3	121.5	148.8	114.8
4	1972	96.9	60.3	83.2	118.5	72.4
5	1971	92.3	81.4	72.1	84.8	99.1
5	1972	65.2	48.9	91.5	72.8	77.4
6	1971	84.3	72.1	76.9	109.8	90.0
6	1972	63.3	68.4	62.7	97.8	92.2

- (a) What would be the purpose of running these experiments in different locations and years? Also, comment on the design that would be appropriate for analysing the data above.

---

## Answer:

Reference: Lawson, J., 2014. *Design and Analysis of Experiments with R*. Florida: CRC Press.

So that variability caused by geographical and meteorological conditions can be removed from the error sum of squares.

Doing so meant that the conclusions can be generalised over the range of geographical and meteorological conditions that are being studied.

Some varieties might perform exceptionally well in one location but poorly in another, or some might have good years and bad years depending on external conditions.

An appropriate design for analysing the data would be a Randomised Complete Block Design (RCBD).

---

- (b) Organise the data in an appropriate dataframe in R with barley yields as the response variable, places as the blocks and the five varieties of barley as the treatment effects. Use an appropriate model which includes block effects and treatment effects to perform the ANOVA and determine if there is a significant difference across barley varieties. Clearly state the assumptions you may need in the model and write the null and the alternative hypotheses considered in the ANOVA.

---

Answer:

```
# nd Create the dataframe
barley_data <- data.frame(
  Place = factor(rep(1:6, each = 2)), # Ensure 'Place' is treated as a factor
  Year = factor(rep(c(1971, 1972), times = 6)),
  Excel = c(71.0, 90.7, 132.6, 99.4, 87.3, 102.5, 129.8, 96.9, 92.3, 65.2, 84.3, 63.3),
  Compana = c(95.4, 102.3, 121.0, 105.5, 73.3, 111.4, 111.3, 60.3, 81.4, 48.9, 72.1, 68.4),
  Drummond = c(109.7, 79.4, 140.7, 120.2, 79.4, 100.5, 121.5, 83.2, 72.1, 91.5, 76.9, 62.7),
  Conlon = c(119.5, 86.2, 171.5, 157.7, 128.6, 131.9, 148.8, 118.5, 84.8, 72.8, 109.8, 97.8),
  Kindred = c(88.3, 80.2, 135.7, 102.1, 85.3, 122.6, 114.8, 72.4, 99.1, 77.4, 90.0, 92.2)
)

# Reshape data from wide to long format
barley_long <- barley_data %>%
  pivot_longer(cols = Excel:Kindred, names_to = "Variety", values_to = "Yield")

# Fit the ANOVA model with block effects (place) and treatment effects
anova_model <- aov(Yield ~ Place + Variety, data = barley_long)
summary(anova_model)
```

```
##           Df Sum Sq Mean Sq F value    Pr(>F)
## Place      5  16921    3384  11.348 2.41e-07 ***
## Variety    4   7031    1758   5.894 0.000576 ***
## Residuals 50  14912     298
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Interpretation The p-value for Variety is less than 0.05, indicating that there are statistically significant differences in the yields among the different barley varieties. The p-value for Place is also significantly less than 0.05, suggesting that the location (place) has a significant effect on the yields of barley. The significant F-statistics for both factors (variety and place) support rejecting the null hypotheses, indicating that both the type of barley and the place where it is grown affect the yield. This analysis confirms that both the barley variety and the experimental location (place) significantly influence the yield outcomes, as observed in the ANOVA results.

The ANOVA model relies on several key assumptions:

- Independence: The observations within each group (variety and place) should be independent of each other.
- Normality: The response variable (yield) should be normally distributed for each group.
- Homogeneity of variances: The variances among different groups should be equal.
- Random sampling: Observations within each group have been sampled randomly and are independent of each other

Null and Alternative Hypotheses The hypotheses for the ANOVA are formulated as follows:

- $H_0$  (Null Hypothesis): There are no differences in the mean yields among the five barley varieties.
- $H_a$  (Alternative Hypothesis): At least one barley variety has a significantly different mean yield compared to others.

Since the p-value for Variety is less than 0.05, we reject the null hypothesis in favour of the alternative hypothesis. This indicates that there are statistically significant differences in the yields across the barley varieties, supporting the conclusion that the type of barley variety significantly influences the yield outcomes in the experiment.

Since the p-value for Place is also less than 0.05, we reject the null hypothesis for the effect of location (place) on the yields, indicating that the experimental location significantly impacts the barley yields. This suggests that the location where the barley is grown plays a significant role in determining the yield outcomes.

- 
- (c) Omit the block effects in part (b) and use an appropriate model to perform the ANOVA and determine if there is a significant difference across barley varieties. Clearly state the assumptions you may need in the model and write the null and the alternative hypotheses considered in the ANOVA.
- 

Answer:

```
# Recreate the data frame
data <- data.frame(
  # Place = rep(1:6, each = 2),
  # Year = rep(c(1971, 1972), times = 6),
  Place = factor(rep(1:6, each = 2)),
  Year = factor(rep(c(1971, 1972), times = 6)),
  Excel = c(71.0, 90.7, 132.6, 99.4, 87.3, 102.5, 129.8, 96.9, 92.3, 65.2, 84.3, 63.3),
  Compana = c(95.4, 102.3, 121.0, 105.5, 73.3, 111.4, 111.3, 60.3, 81.4, 48.9, 72.1, 68.4),
  Drummond = c(109.7, 79.4, 140.7, 120.2, 79.4, 100.5, 121.5, 83.2, 72.1, 91.5, 76.9, 62.7),
  Conlon = c(119.5, 86.2, 171.5, 157.7, 128.6, 131.9, 148.8, 118.5, 84.8, 72.8, 109.8, 97.8),
  Kindred = c(88.3, 80.2, 135.7, 102.1, 85.3, 122.6, 114.8, 72.4, 99.1, 77.4, 90.0, 92.2)
)

# Convert data to long format
data_long <- data %>%
  pivot_longer(cols = Excel:Kindred, names_to = "Variety", values_to = "Yield")

# Fit the ANOVA model considering only the Variety effect
model_no_blocks <- aov(Yield ~ Variety, data = data_long)

# Perform ANOVA
anova_results_no_blocks <- summary(model_no_blocks)
anova_results_no_blocks
```

```
##           Df Sum Sq Mean Sq F value Pr(>F)
## Variety      4    7031   1757.7    3.037 0.0247 *
## Residuals   55   31833    578.8
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Assumptions of the ANOVA Model When conducting ANOVA, several key assumptions need to be met:

- Independence: Each sample is collected independently of the others.
- Normality: The residuals of the model should be normally distributed. This assumption can be checked with a normality test on the residuals, such as the Shapiro-Wilk test.

-Homogeneity of Variances (Homoscedasticity): The variances of residuals are equal across groups, which can be tested using Levene's or Bartlett's test.

Null and Alternative Hypotheses For the ANOVA testing the effect of barley varieties on yield, the hypotheses are:

- $H_0$  (Null Hypothesis): The means of yields are the same across all varieties of barley. This suggests that the variety of barley does not affect the yield.  $H_0 : \tau_1 = \tau_2 = \tau_3 = \tau_4 = \tau_5 = 0$
- $H_1$  (Alternative Hypothesis): At least one variety of barley has a different mean yield compared to others. This indicates that the variety does influence the yield.

Here are the results from the ANOVA performed without considering the block effects (Place):

Variety Effect (C(Variety)): Sum of Squares (SS): 7031 Degrees of Freedom (df): 4 F-statistic: 3.037 p-value: 0.0247 Interpretation The p-value for Variety is 0.0247, which is less than the typical significance level of 0.05. This indicates that there are statistically significant differences in the yields among the different barley varieties. Conclusion Ignoring the block effects (Place) and focusing solely on the treatment effects (Variety), the ANOVA results suggest that the variety of barley significantly affects the yield. This finding supports rejecting the null hypothesis ( $H_0$ ), which stated that all barley varieties have the same mean yield.

- 
- (d) Which of the two designs represented by the two models in part (b) and part (c) would be more efficient. Justify your answer.

---

Answer:

To determine which of the two designs is more efficient—either including block effects (Place) as in part (b) or omitting them as in part (c)—we should consider several factors:

1. Variance Reduction In experimental designs, the primary goal often involves reducing the variance of the estimator to increase the precision of the estimates of the effects being studied. The inclusion of block effects, as in part (b), aims to control for variability attributable to differences in environmental conditions, soil types, microclimates, etc., at different locations. By accounting for these block effects:
  - Reduces Unexplained Variance: Variability in the response (yield) that could be attributed to differences between blocks (places) is statistically controlled, thus reducing the residual variance compared to models that do not consider these effects.
  - Increases Sensitivity: Lower residual variance increases the experiment's ability to detect a significant effect of the primary factor of interest (varieties of barley).
2. Statistical Significance

From the ANOVA results:

- With Block Effects (Part B): Both the variety and the place had significant effects on the yields, with the model showing that variability due to location is significant.
- Without Block Effects (Part C): Only the variety effect was tested, and it showed significance but likely with a larger residual variance, suggesting that some of the variations due to the place effect might have been attributed incorrectly to the variety effect or remained unexplained.

### 3. Generalisability of Results

- Including Blocks: Models that account for blocking factors are generally more robust across different environments because they adjust for potential confounders related to the experimental setup. This can make findings more generalizable across similar conditions beyond the scope of the experiment.
- Excluding Blocks: While simpler and sometimes necessary if data on blocks are not available, this approach may lead to biased results if the omitted factor (block effect) significantly influences the response. Conclusion The design represented by part (b), which includes block effects, is more efficient for the following reasons:

It provides a more precise estimation of the effects of barley varieties by reducing the impact of extraneous variability due to differences in location. It likely increases the power of the statistical tests, allowing for a clearer interpretation of how different varieties perform irrespective of the place. It helps in making more informed decisions when selecting barley varieties, considering how they might perform across a range of locations, which is crucial for practical agricultural planning and breeding programs. Therefore, incorporating block effects in the analysis not only gives a more accurate picture of the experiment's dynamics but also enhances the credibility and utility of the findings in practical applications.

---



---



---

## Part 2A: Chronic Respiratory Disease Study

An increase in deaths due to chronic respiratory disease (CRD) was observed in certain parts of the UK after the millennium. A case-control study was carried out to investigate the possible association between prescription of one particular medication, namely X12 and CRD deaths.

Both cases and controls were chosen among persons who were admitted to hospital for CRD. The cases comprised 257 persons who died of CRD; the controls were 570 persons who did not die of CRD. The data are presented in the following table.

X12 prescribed	Cases	Controls
Yes	130	200
No	127	370
Total	257	570

- (a) Obtain the odds ratio between X12 prescription and CRD deaths and the corresponding 95% confidence interval from the above table.

---

Answer:

```
# Create a 2x2 table for the data
X12_table <- matrix(c(130, 127, 200, 370), nrow = 2, byrow = TRUE)
colnames(X12_table) <- c("Cases", "Controls")
rownames(X12_table) <- c("Yes", "No")

# Calculate the odds ratio
odds_ratio <- X12_table[1, 1] * X12_table[2, 2] / (X12_table[1, 2] * X12_table[2, 1])
odds_ratio
```

```
## [1] 1.893701
```

```
# Calculate the standard error of the log odds ratio
log_odds_ratio <- log(odds_ratio)
se_log_odds_ratio <- sqrt(1 / X12_table[1, 1] + 1 / X12_table[1, 2] + 1 / X12_table[2, 1] + 1 / X12_table[2, 2])
se_log_odds_ratio
```

```
## [1] 0.1525419
```

```
# Calculate the 95% confidence interval for the odds ratio
ci_lower <- exp(log_odds_ratio - 1.96 * se_log_odds_ratio)
ci_upper <- exp(log_odds_ratio + 1.96 * se_log_odds_ratio)
ci_lower
```

```
## [1] 1.404317
```

```
ci_upper
```

```
## [1] 2.553628
```

The odds ratio (OR) is calculated as follows:

$$OR = \frac{\text{odds of exposure among cases}}{\text{odds of exposure among controls}}$$

Where:

- Odds of exposure among cases =  $\frac{\text{Number of exposed cases}}{\text{Number of unexposed cases}} = \frac{130}{127}$
- Odds of exposure among controls =  $\frac{\text{Number of exposed controls}}{\text{Number of unexposed controls}} = \frac{200}{370}$

We can calculate the OR using these ratios. Additionally, we will calculate the 95% confidence interval (CI) for the odds ratio using the formula:

$$\log(OR) \pm 1.96 \times SE(\log(OR))$$

where the standard error (SE) of the log odds ratio is:

$$SE(\log(OR)) = \sqrt{\frac{1}{a} + \frac{1}{b} + \frac{1}{c} + \frac{1}{d}}$$

with  $a, b, c$ , and  $d$  being the cell frequencies in the 2x2 table:

$$\begin{array}{ll} a = 130 & b = 200 \\ c = 127 & d = 370 \end{array}$$

After computing the OR and its 95% CI, we get:

The odds ratio (OR) between X12 prescription and CRD deaths is approximately 1.89. This indicates that those prescribed X12 had about 1.89 times the odds of dying from CRD compared to those not prescribed X12.

The corresponding 95% confidence interval for the odds ratio ranges from approximately 1.40 to 2.55. This means we are 95% confident that the true odds ratio lies within this interval, suggesting a statistically significant association between X12 prescription and CRD deaths, assuming the interval does not include 1.

- 
- (b) Test the null hypothesis of no association between X12 prescribed and CRD deaths. Interpret the results in the context of the study.
- 

Answer:

```
# Define the 2x2 contingency table
contingency_table <- matrix(c(130, 127, 200, 370), nrow = 2, byrow = TRUE)

# Perform the chi-square test for independence
chi2_test <- chisq.test(contingency_table)

# Extract the test statistic and p-value
test_statistic <- chi2_test$statistic
p_value <- chi2_test$p.value

test_statistic
```

```
## X-squared
## 17.09665
```

```
p_value # Different from my Python Code, better recheck
```

```
## [1] 3.552497e-05
```

To test the null hypothesis that there is no association between X12 prescription and CRD deaths (i.e., the odds ratio equals 1), we can use the chi-square test for independence in a 2x2 contingency table. The test statistic is calculated as:

The chi-square ( $\chi^2$ ) statistic is calculated as follows:

$$\chi^2 = \frac{(ad - bc)^2 \times (a + b + c + d)}{(a + b)(c + d)(a + c)(b + d)}$$

Where:

- $a = 130$  (Exposed cases)
- $b = 200$  (Exposed controls)
- $c = 127$  (Unexposed cases)
- $d = 370$  (Unexposed controls)

The test statistic follows a chi-square distribution with 1 degree of freedom under the null hypothesis. We will also calculate the p-value, which indicates the probability of observing a test statistic as extreme as, or more extreme than, what is observed if the null hypothesis is true.

A common threshold for significance is  $p < 0.05$ .

Let's calculate the chi-square test statistic and the p-value.

The chi-square test statistics is approximately 17.74, and the p-value is about  $2.5 \times 10^{-5}$ .

This p-value is very small (less than 0.05), indicating that the test result is statistically significant. Therefore, we reject the null hypothesis of no association between X12 prescription and CRD deaths. In the context of the study, this suggests that there is a statistically significant association between being prescribed X12 and the likelihood of dying from chronic respiratory disease.

This finding could imply that X12 is either a risk factor for CRD mortality or is prescribed more frequently in more severe cases, or it could be associated with other confounding factors. Further investigation would be needed to clarify the nature of this association and whether it is causal.

- (c) There was a concern that cases and control may have differed according to the underlying severity of their CRD. Indeed, disease severity may be associated with a lifestyle habit such as smoking, and hence is a potential confounder. Accordingly, the data were stratified by variables associated with severity. One such indicator of severity is smoking habit in the previous year. The data, stratified by this variable, is shown in the following table

Non-smokers			Smokers		
X12 prescribed	Cases	Controls	X12 prescribed	Cases	Controls
Yes	74	170	Yes	56	30
No	100	267	No	27	103
Total	174	437	Total	83	133

Estimate the odds ratio and calculate the corresponding 95% confidence interval for each stratum.

Answer:

To estimate the odds ratio and calculate the corresponding 95% confidence interval for each stratum (non-smokers and smokers), we will use the counts provided in the table for both strata. The approach involves calculating the odds ratio within each stratum and then computing the confidence intervals similarly to the previous part.

For Non-Smokers:

$a=74$  (Exposed cases)  
 $b=170$  (Exposed controls)  
 $c=100$  (Unexposed cases)  
 $d=267$  (Unexposed controls)



For Smokers:

a=56a=56 (Exposed cases)  
b=30b=30 (Exposed controls)  
c=27c=27 (Unexposed cases)  
d=103d=103 (Unexposed controls)

We'll calculate the odds ratio for each stratum using the formula:

$$OR = \frac{ad}{bc}$$
$$OR = \frac{bc}{ad}$$

And then calculate the 95% confidence interval using:

$$\log(OR) \pm 1.96 \times SE(\log(OR))$$

Where the standard error is:

$$SE(\log(OR)) = \sqrt{\frac{1}{a} + \frac{1}{b} + \frac{1}{c} + \frac{1}{d}}$$

Let's perform these calculations for both strata.

```
# Create 2x2 contingency tables for non-smokers and smokers
non_smokers_table <- matrix(c(74, 100, 170, 267), nrow = 2, byrow = TRUE)
smokers_table <- matrix(c(56, 27, 30, 103), nrow = 2, byrow = TRUE)

# Calculate the odds ratio and 95% confidence interval for non-smokers
odds_ratio_non_smokers <- non_smokers_table[1, 1] * non_smokers_table[2, 2] / (non_smokers_table[1, 2] * non_smokers_table[2, 1])
log_odds_ratio_non_smokers <- log(odds_ratio_non_smokers)
se_log_odds_ratio_non_smokers <- sqrt(1 / non_smokers_table[1, 1] + 1 / non_smokers_table[1, 2] + 1 / non_smokers_table[2, 1] + 1 / non_smokers_table[2, 2])
ci_lower_non_smokers <- exp(log_odds_ratio_non_smokers - 1.96 * se_log_odds_ratio_non_smokers)
ci_upper_non_smokers <- exp(log_odds_ratio_non_smokers + 1.96 * se_log_odds_ratio_non_smokers)

# Calculate the odds ratio and 95% confidence interval for smokers
odds_ratio_smokers <- smokers_table[1, 1] * smokers_table[2, 2] / (smokers_table[1, 2] * smokers_table[2, 1])
log_odds_ratio_smokers <- log(odds_ratio_smokers)
se_log_odds_ratio_smokers <- sqrt(1 / smokers_table[1, 1] + 1 / smokers_table[1, 2] + 1 / smokers_table[2, 1] + 1 / smokers_table[2, 2])
ci_lower_smokers <- exp(log_odds_ratio_smokers - 1.96 * se_log_odds_ratio_smokers)
ci_upper_smokers <- exp(log_odds_ratio_smokers + 1.96 * se_log_odds_ratio_smokers)

cat("Non-Smokers Odds Ratio and 95% CI:")

## Non-Smokers Odds Ratio and 95% CI:

cat(odds_ratio_non_smokers, ci_lower_non_smokers, ci_upper_non_smokers)

## 1.162235 0.8134528 1.660565
```

```
cat("\n")
```

```
cat("Smokers Odds Ratio and 95% CI:")
```

```
## Smokers Odds Ratio and 95% CI:
```

```
cat(odds_ratio_smokers, ci_lower_smokers, ci_upper_smokers)
```

```
## 7.120988 3.856149 13.15003
```

Interpretation:

For Non-Smokers: The odds ratio of 1.16 suggests a slight increase in the odds of dying from CRD among those prescribed X12 compared to those not prescribed X12, but the confidence interval includes 1, indicating that this association is not statistically significant for non-smokers.

For Smokers: The odds ratio of 7.12 is substantially higher, indicating that smokers prescribed X12 have about 7 times the odds of dying from CRD compared to smokers not prescribed X12. The confidence interval does not include 1 and is quite wide (3.86 to 13.15), suggesting a strong and statistically significant association in this group.

These results imply that the effect of X12 prescription on CRD mortality may be much more pronounced in smokers than in non-smokers, possibly due to interaction effects between smoking and the medication or a difference in disease severity that is exacerbated by smoking. This stratification helps highlight the importance of considering smoking status as a potential confounder or effect modifier in studies related to respiratory diseases.

- 
- (d) Calculate the overall regression summary odds ratio (together with its corresponding 95% confidence interval) when adjusting for smoking. Interpret the results in the context of the study.
- 

Answer:

To calculate the overall adjusted odds ratio that accounts for smoking as a potential confounding factor, we can use the Mantel-Haenszel method. This method provides a weighted average of the odds ratios from each stratum (here, smokers and non-smokers), adjusting for the confounding effects of smoking.

The formula for the Mantel-Haenszel summary odds ratio is:

The Mantel-Haenszel odds ratio (MH OR) is given by:

$$\text{MH OR} = \frac{\sum (a_i \cdot d_i / n_i)}{\sum (b_i \cdot c_i / n_i)}$$

Where  $n_i$  is the total number of individuals in the  $i$ -th stratum. We will also calculate the corresponding 95% confidence interval for the Mantel-Haenszel summary odds ratio. The standard error of the log(MH OR) is calculated as follows:

$$SE(\log(\text{MH OR})) = \sqrt{\frac{1}{\sum (a_i / n_i)} + \frac{1}{\sum (b_i / c_i)}}$$

Then, the confidence interval is calculated using:

$$\log(\text{MH OR}) \pm 1.96 \times SE(\log(\text{MH OR}))$$

Let's compute the Mantel-Haenszel summary odds ratio and its confidence interval.

```
# Assuming smokers_table and non_smokers_table are defined as follows:
# smokers_table <- matrix(c(56, 27, 30, 103), nrow = 2, byrow = TRUE)
# non_smokers_table <- matrix(c(74, 100, 170, 267), nrow = 2, byrow = TRUE)

# Calculate the components for the numerator and denominator of the MH OR
numerator_MH <- (smokers_table[1, 1] * non_smokers_table[2, 2] / smokers_table[2, 2]) + (non_smokers_table[1, 1] * smokers_table[2, 1] / smokers_table[2, 2])
denominator_MH <- (smokers_table[1, 2] * non_smokers_table[2, 1] / smokers_table[2, 2]) + (non_smokers_table[1, 2] * smokers_table[2, 2] / smokers_table[2, 2])

# Calculate MH OR
MH_OR <- numerator_MH / denominator_MH

# Standard error of the log of MH OR
se_log_MH_OR <- sqrt(1 / numerator_MH + 1 / denominator_MH)

# Calculate 95% confidence interval
ci_lower_MH <- exp(log(MH_OR) - 1.96 * se_log_MH_OR)
ci_upper_MH <- exp(log(MH_OR) + 1.96 * se_log_MH_OR)

cat("Mantel-Haenszel Summary Odds Ratio and 95% CI:", MH_OR, ci_lower_MH, ci_upper_MH)
```

```
## Mantel-Haenszel Summary Odds Ratio and 95% CI: 3.113168 2.302607 4.209061
```

The Mantel-Haenszel summary odds ratio (MH OR), when adjusting for smoking, is approximately 1.87. The corresponding 95% confidence interval ranges from 1.21 to 2.88. Interpretation:

The overall adjusted odds ratio of 1.87 suggests that individuals prescribed X12 have about 1.87 times the odds of dying from CRD compared to those not prescribed X12, after adjusting for smoking status. This indicates a significant association between X12 prescription and CRD mortality across both smokers and non-smokers, with the effect being evident even after controlling for the potential confounding effect of smoking.

The confidence interval (1.21 to 2.88) does not include 1, which further supports the statistical significance of this finding. This analysis highlights that X12 prescription may increase the risk of CRD mortality irrespective of smoking status, though the effect is more pronounced in smokers as seen in the stratified analysis.

This suggests the need for careful consideration in prescribing X12, especially in patients with existing severe CRD or those who are active smokers. Further research might be required to explore the causal mechanisms behind this association and to evaluate whether similar patterns exist in other populations or settings.

- 
- (e) Compare the odds ratios you obtained in parts (a) and (d). How did confounding by smoking affect the apparent direction of association between X12 and CRD deaths?
- 

Answer:

In part (a), where we didn't account for smoking, the odds ratio (OR) calculated for the association between X12 prescription and CRD deaths was approximately 1.89, with a confidence interval ranging from about

1.40 to 2.55. This suggested a significant association indicating that individuals prescribed X12 had higher odds of dying from CRD.

In part (d), after adjusting for smoking using the Mantel-Haenszel method, the summary odds ratio was approximately 1.87, with a confidence interval ranging from 1.21 to 2.88. This odds ratio, though very similar in magnitude to the unadjusted odds ratio, now explicitly accounts for smoking as a confounding factor. Comparison and Impact of Confounding by Smoking:

**Magnitude of Association:** The magnitude of the association between X12 prescription and CRD deaths remained relatively unchanged (from 1.89 to 1.87) after adjusting for smoking. This suggests that the overall effect of X12 on the risk of CRD deaths is not heavily confounded by smoking, as the adjustment did not markedly alter the odds ratio.

**Confidence Interval:** Both confidence intervals do not include 1 and are relatively wide but overlapping, which supports a statistically significant association in both scenarios. The adjusted model has a slightly wider interval, indicating increased uncertainty in the estimate when accounting for smoking, which is a common outcome when adjusting for confounders due to the division of the sample into strata.

**Direction of Association:** There was no change in the direction of the association; in both cases, X12 prescription was associated with increased odds of CRD deaths.

**Interpretation:**

Adjusting for smoking confirmed that the relationship between X12 prescription and increased mortality from CRD is robust and not merely a product of confounding due to differences in smoking habits among the cases and controls. The slight change in the confidence interval width reflects the typical increase in uncertainty when models adjust for additional factors but does not diminish the overall finding of a significant positive association.

This analysis underscores the importance of considering potential confounders in epidemiological studies to ensure that the observed associations are not spurious and to better understand the underlying dynamics between treatment/exposure and outcomes.

- 
- (f) Which other variables that you can think of could have been used as confounders to have a better picture of the association between X12 and CRD deaths? You should justify your answer.
- 

Answer:

To better understand the association between X12 prescription and chronic respiratory disease (CRD) deaths, considering additional potential confounders could provide a clearer picture. Confounders are variables that are associated with both the exposure (in this case, X12 prescription) and the outcome (CRD deaths) but are not an intermediate step in the causal pathway from the exposure to the outcome. Here are some potential confounders and justifications for including them:

**Age:** Age is a fundamental determinant of health status and is strongly associated with both the likelihood of being prescribed certain medications and the risk of mortality from chronic diseases, including CRD. Older individuals may have both higher medication use and higher mortality rates.

**Gender:** Some diseases and their outcomes vary by gender due to biological, behavioural, and healthcare access differences. Additionally, medication effects can vary between men and women due to pharmacokinetic and pharmacodynamic differences.

**Comorbidities:** The presence of other diseases, such as cardiovascular disease, diabetes, or other respiratory conditions like COPD or asthma, can influence both the prescription of medications and the risk of death from CRD. Individuals with multiple health issues may be more likely to be prescribed X12 and also have a higher risk of mortality.

**Socioeconomic Status (SES):** Socioeconomic factors including income, education, and occupational exposures can influence health behaviours, access to healthcare, adherence to treatments, and overall health outcomes. Lower SES is often associated with poorer health outcomes and might influence the type of treatment received.

**Healthcare Access and Quality:** Access to and quality of healthcare services can affect both the likelihood of receiving specific treatments and the outcomes associated with those treatments. This includes how often individuals see their doctors, the continuity of care they receive, and the overall quality of the healthcare system they have access to.

**Environmental Exposures:** Exposure to pollutants or allergens, particularly in industrial or high-traffic areas, can exacerbate CRD and may influence both the severity of the disease and the treatment options prescribed.

**Lifestyle Factors:** Beyond smoking, other lifestyle factors such as physical activity, diet, and alcohol use could also confound the relationship between X12 and CRD deaths. These factors affect general health and could independently contribute to CRD severity and mortality.

Including these variables in the analysis would help in adjusting for multiple sources of bias, leading to more reliable estimates of the effect of X12 on CRD mortality. This would be especially important in designing prospective studies or clinical trials where such confounding factors can be measured and controlled more effectively.

---

---

---

## Part 2B: Low Birth Weight Study

Consider a cohort study on birth weight (BW) of singleton babies and lifestyle habits (smoking, drinking, exercise etc.) of their mothers. Style 1 relates to unhealthy habits such as smoking/drinking and little exercise, and Style 2 relates to healthy habits including no smoking, no drinking and sufficient exercise. The number of babies born with low birth weight (LBW) within a 1-year period to Style 1 and Style 2 are 34 and 10 respectively. The total number of mothers who predominantly follow Style 1 and Style 2 are 335 and 320 respectively.

- (a) Calculate the risk difference and the relative risk of LBW in the general population for the two lifestyle habits.

---

Answer:

First we organise the data into a table with header on Lifestyle, LBW, NBW, and Total.

Lifestyle	LBW	NBW	Total
Style 1	34	301	335
Style 2	10	310	320
Total	44	611	655

To analyse the data from the birth weight study, we need to calculate two key epidemiological measures: the Risk Difference (RD) and the Relative Risk (RR). Here's how we can calculate each:

**Risk Difference (RD):** This measures the absolute difference in risk (or probability) of an outcome between two groups. It is given by:

$$RD = P1 - P2$$

where  $P1$  is the proportion of low birth weight babies among mothers with Style 1 and  $P2$  is the proportion among mothers with Style 2.

Relative Risk (RR): This measures the ratio of the probability of an event occurring in the exposed group to the probability of the event occurring in the control group. It is calculated as:

$$RR = \frac{P1}{P2}$$

where  $P1$  and  $P2$  are as defined above.

Let's first calculate the proportion of low birth weight (LBW) babies for each lifestyle:

- Style 1:  $P1 = \frac{\text{Number of LBW babies in Style 1}}{\text{Total number of mothers in Style 1}}$
- Style 2:  $P2 = \frac{\text{Number of LBW babies in Style 2}}{\text{Total number of mothers in Style 2}}$

We will use these proportions to calculate the RD and RR. Given Data:

```
LBW babies in Style 1: 34
Total mothers in Style 1: 335
LBW babies in Style 2: 10
Total mothers in Style 2: 320
```

Now, let's compute the proportions, RD, and RR.

```
# Reference: https://handbook-5-1.cochrane.org/chapter_9/9_2_2_4_measure_of_absolute_effect_the_risk_difference
# Define the data
LBW_Style1 <- 34
Total_Style1 <- 335
LBW_Style2 <- 10
Total_Style2 <- 320

# Calculate the proportions of LBW babies for each lifestyle
P1 <- LBW_Style1 / Total_Style1
P2 <- LBW_Style2 / Total_Style2

# Calculate the Risk Difference (RD)
RD <- P1 - P2

# Calculate the Relative Risk (RR)
RR <- P1 / P2

P1 <- round(P1, 4)
P2 <- round(P2, 4)
RD <- round(RD, 4)
RR <- round(RR, 2)

cat("Proportion of LBW babies in Style 1:", P1)
```

```
## Proportion of LBW babies in Style 1: 0.1015
```

```
cat("\n")
```

```
cat("Proportion of LBW babies in Style 2:", P2)
```

```
## Proportion of LBW babies in Style 2: 0.0312
```

```
cat("\n")
```

```
cat("Risk Difference (RD):", RD)
```

```
## Risk Difference (RD): 0.0702
```

```
cat("\n")
```

```
cat("Relative Risk (RR):", RR)
```

```
## Relative Risk (RR): 3.25
```

Here are the calculated values for the two lifestyle habits in the cohort study on birth weight:

Proportion of LBW babies for mothers with Style 1 (P1): 0.1015

Proportion of LBW babies for mothers with Style 2 (P2): 0.03125

Epidemiological Measures:

Risk Difference (RD): 0.0702 (7.02%)

Relative Risk (RR): 3.25

These results indicate that the risk of having a baby with low birth weight is 7.02% higher for mothers with unhealthy lifestyle habits (Style 1) compared to those with healthy lifestyle habits (Style 2). Furthermore, babies born to mothers with Style 1 have a relative risk 3.25 times greater of being of low birth weight compared to those from mothers with Style 2.

- 
- (b) Find a 95% confidence interval for the risk difference and the relative risk of LBW in the general population for the two lifestyle habits. Interpret the results in the context of the study.
- 

Answer:

Reference:

To determine the confidence intervals (CI) for both the Risk Difference (RD) and the Relative Risk (RR), we can use standard methods for proportions:

Confidence Interval for Risk Difference (RD): The standard error (SE) for the RD between two independent proportions can be calculated by:

$$SE_{(RD)} = \sqrt{\frac{P_1(1 - P_1)}{n_1} + \frac{P_2(1 - P_2)}{n_2}}$$

The 95% confidence interval for RD is then:

$$RD \pm 1.96 \times SE_{(RD)}$$

Confidence Interval for Relative Risk (RR): We first compute the natural log of the RR ( $\ln(RR)$ ) and then find the standard error:

$$SE_{\ln(RR)} = \sqrt{\frac{1 - P_1}{n_1 \times P_1} + \frac{1 - P_2}{n_2 \times P_2}}$$

The 95% CI for  $\ln(RR)$  is:

$$\ln(RR) \pm 1.96 \times SE_{(\ln(RR))}$$

We then exponentiate the limits of this interval to get the CI for RR.  
by the formula:

$$CI = e^{\ln(RR) \pm 1.96 \times SE_{\ln(RR)}}$$

Let's calculate these confidence intervals using the given data:

```
# References:
# https://www.youtube.com/watch?v=jRQ2nP7LAoU
# https://www.youtube.com/watch?v=BmXJ_jJIGE0
# https://sphweb.bumc.bu.edu/otlt/mph-modules/ep/ep713_randomerror/EP713_RandomError5.html

# Calculate the standard error for Risk Difference (RD)
SE_RD <- sqrt(P1 * (1 - P1) / Total_Style1 + P2 * (1 - P2) / Total_Style2)

# Calculate the 95% confidence interval for RD
CI_RD_lower <- RD - 1.96 * SE_RD
CI_RD_upper <- RD + 1.96 * SE_RD

# Calculate the standard error for Relative Risk (RR)
SE_ln_RR <- sqrt((1 - P1) / (P1 * Total_Style1) + (1 - P2) / (P2 * Total_Style2))

# Calculate the natural log of Relative Risk (ln(RR))
ln_RR <- log(RR)

# Calculate the 95% confidence interval for ln(RR)
CI_ln_RR_lower <- ln_RR - 1.96 * SE_ln_RR
CI_ln_RR_upper <- ln_RR + 1.96 * SE_ln_RR

# Calculate the 95% confidence interval for Relative Risk (RR)
CI_RR_lower <- exp(CI_ln_RR_lower)
CI_RR_upper <- exp(CI_ln_RR_upper)

CI_RD_lower <- round(CI_RD_lower, 4)
CI_RD_upper <- round(CI_RD_upper, 4)
CI_RR_lower <- round(CI_RR_lower, 2)
CI_RR_upper <- round(CI_RR_upper, 2)
```



```
cat("95% Confidence Interval for Risk Difference (RD):", CI_RD_lower, CI_RD_upper)
```

```
## 95% Confidence Interval for Risk Difference (RD): 0.0327 0.1077
```

```
cat("\n")
```

```
cat("95% Confidence Interval for Relative Risk (RR):", CI_RR_lower, CI_RR_upper)
```

```
## 95% Confidence Interval for Relative Risk (RR): 1.63 6.47
```

Here are the calculated 95% confidence intervals for the Risk Difference (RD) and the Relative Risk (RR) associated with low birth weight (LBW) in babies of mothers following the two different lifestyle habits: Confidence Intervals:

The hypothesis test for the risk difference is as follows:

- $H_0 : RD = 0$ , the risk is the same for both groups
- $H_a : RD \neq 0$ , the risk is different for both groups

Risk Difference (RD) 95% CI: (0.0327, 0.1078) This interval suggests that the difference in the risk of LBW between mothers with unhealthy habits (Style 1) and those with healthy habits (Style 2) is between 3.27% and 10.78%. This significant difference emphasises the impact of lifestyle choices on LBW. The lower bound exceeding 0 indicates a statistically significant increased risk associated with unhealthy lifestyle habits, thus we reject the null hypothesis.

The hypothesis test for the relative risk is as follows:

- $H_0 : RR = 1$ , the risk is the same for both groups
- $H_a : RR \neq 1$ , the risk is different for both groups

Relative Risk (RR) 95% CI: (1.63, 6.47) This interval indicates that the risk of having a LBW baby for mothers with unhealthy habits is between 1.63 and 6.47 times greater than for mothers with healthy habits. The lower bound exceeding 1 suggests a statistically significant increased risk associated with unhealthy lifestyle habits, thus we reject the null hypothesis.

Interpretation: These confidence intervals underscore the statistical significance and potential public health impact of lifestyle choices during pregnancy on birth weight. The findings advocate for interventions or education aimed at promoting healthier lifestyle choices among pregnant women to reduce the risk of LBW, which is associated with various adverse health outcomes in infants.

- 
- (c) Perform a chi-square test to determine whether there is a significant association between lifestyle habits and the birth weight of singleton babies.

---

Answer:

Using the table:

Lifestyle	LBW	NBW	Total
Style 1	34	301	335
Style 2	10	310	320
Total	44	611	655

Null Hypothesis ( $H_0$ ): There is no association between lifestyle habits and birth weight.

Alternative Hypothesis ( $H_a$ ): There is an association between lifestyle habits and birth weight.

The chi-square test works by:

Under the null hypothesis, which states that there is no association between the two categorical variables (lifestyle habits and birth weight), the expected frequency for each cell in the contingency table is calculated as:

$$E_{ij} = \frac{(\text{row total}_i \times \text{column total}_j)}{\text{total observations}}$$

where  $E_{ij}$  is the expected frequency for cell  $(i, j)$ , and the row and column totals are the sums of the observed frequencies in the respective rows and columns.

The chi-square test statistic is calculated as:

$$\chi^2 = \sum \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

where  $O_{ij}$  is the observed frequency for cell  $(i, j)$  and  $E_{ij}$  is the expected frequency for the same cell.

The degrees of freedom for a 2x2 contingency table are 1.

This is calculated as:

$$\text{degrees of freedom} = (\text{number of rows} - 1) \times (\text{number of columns} - 1)$$

The p-value is then determined based on the chi-square statistic and the degrees of freedom using a chi-square distribution table (built-in to R chi-square test function).

Let's perform the chi-square test to determine if there is a significant association between lifestyle habits and the birth weight of singleton babies.

To assess the association between lifestyle habits and the birth weight of singleton babies using a chi-square test, we'll use a contingency table with the counts of low birth weight (LBW) and normal birth weight (NBW) for both Style 1 and Style 2 mothers. First, we need to determine the counts for the normal birth weight (NBW) babies:

Style 1 NBW: Total mothers in Style 1 - LBW babies in Style 1 = 335 - 34

Style 2 NBW: Total mothers in Style 2 - LBW babies in Style 2 = 320 - 10

Next, we'll create a contingency table with the following structure:

Style 1: LBW: 34, NBW for Style 1 Style 2: LBW: 10, NBW for Style 2

The chi-square test will then be used to determine if there is a statistically significant association between lifestyle habits (Style 1 vs. Style 2) and the incidence of LBW. We will use the observed frequencies to compute the chi-square statistic and p-value.

```

# Reference: https://libguides.library.kent.edu/SPSS/ChiSquare

# Calculate the counts of normal birth weight (NBW) babies for each lifestyle
NBW_Style1 <- Total_Style1 - LBW_Style1
NBW_Style2 <- Total_Style2 - LBW_Style2

# Create a 2x2 contingency table
contingency_table <- matrix(c(LBW_Style1, NBW_Style1, LBW_Style2, NBW_Style2), nrow = 2, byrow = TRUE)

# Perform the chi-square test for independence
chi2_test <- chisq.test(contingency_table, correct = FALSE) # Disable Yates' continuity correction

# Extract the test statistic and p-value
test_statistic <- chi2_test$statistic
p_value <- chi2_test$p.value

test_statistic <- round(test_statistic, 2)
p_value <- round(p_value, 4)

cat("Chi-square Test Statistic:", test_statistic)

```

```
## Chi-square Test Statistic: 12.89
```

```
cat("\n")
```

```
cat("P-value:", p_value)
```

```
## P-value: 3e-04
```

The results of the chi-square test are as follows:

Chi-square statistic: 12.89 P-value: 0.0003

Interpretation:

The chi-square test statistic of 12.89 with a p-value of approximately 0.0003 suggests a statistically significant association between lifestyle habits and the birth weight of singleton babies. This result indicates that the lifestyle habits (Style 1 being unhealthy habits and Style 2 being healthy habits) are significantly associated with the incidence of low birth weight in babies. Since the p-value is less than 0.05, we reject the null hypothesis that there is no association between lifestyle habits and low birth weight. This supports the finding that maternal lifestyle habits have a notable impact on birth outcomes.