

Applied Data Science

DATA ANALYSIS ROADMAP

Shen, Yuchen

Zhejiang University

DATA ANALYSIS ROADMAP

YUCHEN SHEN

Table of contents

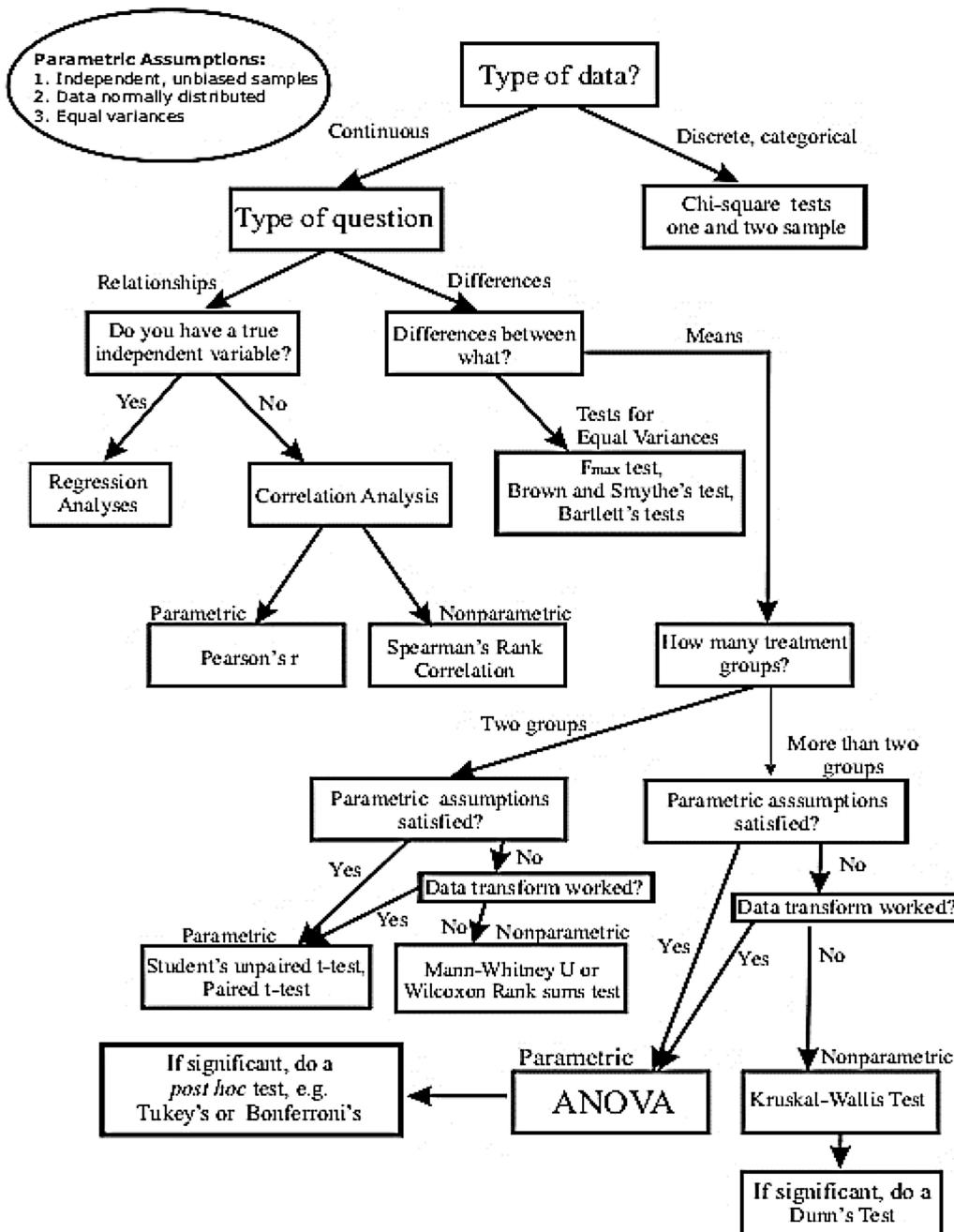
OVERVIEW	3
GENERAL REMINDERS	3
BEFORE-CODING-PREPARATIONS	4
WINDOWS OPERATIONS	4
DATA IMPORT, EXPLORATION & CLEANING	5
IMPORT THE DATA	5
EXPLORE THE DATA (MUST)	5
DATA CLEANING	5
VISUALIZATION	8
T-TEST	9
RUN A T-TEST	9
ANOVA	12
OVERVIEW	12
1-WAY ANOVA BY SIMULATION	12
1-WAY ANOVA BY STATISTICAL TESTS	12
2-WAY ANOVA BY STATISTICAL TESTS	13
POWER ANALYSIS	15
OVERVIEW	15
BALANCE α AND β LEVEL	15
T-TEST: CALCULATE POWER (SIMULATION)	15
T-TEST: CALCULATE POWER (BUILT-IN FUNCTION)	15
T-TEST: CALCULATE THE MINIMUM SAMPLE SIZE FOR AN AIMED POWER (SIMULATION)	15
T-TEST: PARAMETER SETTINGS (BOTH)	16
ANOVA: CALCULATE POWER (SIMULATION)	16
ANOVA: CALCULATE THE MINIMUM SAMPLE SIZE FOR AN AIMED POWER (SIMULATION)	16
CHI-SQUARED TEST	18
OVERVIEW	18
GOODNESS-OF-FIT TEST	18
TEST FOR HOMOGENEITY	19
TEST FOR INDEPENDENCE	19
3-WAY CHI-SQUARED TEST	19
BOOTSTRAPPING	20
OVERVIEW	20
JUSTIFY STATISTICAL TEST CHOICE	20
STEPS	20
EXAMINE ORIGINAL DISTRIBUTION	20
BOOTSTRAP FOR HYPOTHESIS TESTING	20
BOOTSTRAP FOR CONFIDENCE INTERVAL	21
CORRELATION AND REGRESSION	24
REGRESSION (ASSUMPTIONS ARE MET)	24
REGRESSION (ASSUMPTIONS NOT MET)	25
OPTIMIZE A MULTI-VARIABLE MODEL	25
TIME SERIES	26
CORE	26
CONSTRUCT TIME SERIES OBJECT	26
VISUALIZE TIME SERIES	26
DATA TRANSFORMATION (OPTIONAL)	26
CHECK ASSUMPTIONS	26

TABLE OF CONTENTS

NAVIGATE THE TIME SERIES COMPONENT	26
PREDICTIONS	27
IMPROVE MODEL AND PREDICTION	28
<u>ORDINARY DIFFERENTIAL EQUATIONS MODELING</u>	29
CONCEPT OVERVIEW	29
TUMOR GROWTH	29
BACTERIA GROWTH	30
FULL SIR MODEL - NON MARKOVIAN, CONTINUOUS	30
<u>MATRIX-BASED MODELING</u>	32
CONCEPT OVERVIEW	32
SIR EARLY-STAGE MODEL	32
FULL SIR - MARKOVIAN, DISCRETE	33
VESSEL BLOOD FLOW MODEL	35
<u>CONDITIONAL PROBABILITY</u>	38
BASICS	38
CONDITIONAL PROBABILITY	39
<u>BAYESIANISM</u>	40
BAYES' THEOREM	40
BAYESIAN LOGIC COMPONENTS	40
BAYESIAN FACTOR	40
A DICE GAME	41
<u>R AS A CALCULATOR</u>	42
BASIC ARITHMETIC OPERATIONS	42
MATHEMATICAL FUNCTIONS	42
VECTORIZED OPERATIONS AND VECTOR FUNCTIONS	42
MATRIX OPERATIONS	42
SOLVING EQUATIONS	43
<u>PROBABILITY PRACTICE PROBLEMS</u>	44
1 IVF TEST	44
2 VARICOCELE SURGERY	45
<u>CRITICAL THINKING & FUTURE DIRECTIONS</u>	47
<u>CONCEPT TABLE</u>	48
DISTRIBUTION TERMS	48
SAMPLING TERMS	48
GENERAL STATISTICAL TERMS	48
HYPOTHESIS TESTING TERMS	49
CORRELATION TERMS	50

Overview

Flow Chart for Selecting Commonly Used Statistical Tests



General Reminders

1. It is of extreme importance to notice *if the data is paired*, e.g., one before drug injection one after drug injection. If so, only the *difference* between two time points are meant to be used. So you should calculate it correspondingly.
2. One hack of re-structuring data during data cleaning is to let the independent variables be column names, while dependent variables be the filling values.
3. All the plots must be colored.
4. Just don't choose Kruskal-Wallis test when running ANOVA.
5. ANOVA plots = boxplot/violin plot (preferred, go visit www.r-graph-gallery.com)

Data import, exploration & cleaning

Import the Data

```

1 #Note: set ````{r eval=T, echo=F}
2 data <- read.table(file = "", header = , sep = "\t") # .tsv + .txt
3 data <- read_csv(file = "", col_names = , sep = "") # .csv

```

If built-in datasets are to be used, such as `data("cars")`:

```

1 data("cars")
2 Cars <- cars # next, use `Cars`"

```

Explore the data (must)

Explore the data to see if needs data cleaning.

```

1 str(data)
2 head(data) # must be written in the report!!!
3 table(data$col) # check the distribution of a column

```

Self-exploration (set `echo = F`)

```

1 dfSummary(data,graph.col=TRUE,style="grid") %>% stview()

```

Data cleaning

- Explanation
- Cleaning for each block
- Final re-examination

1 Remove unnecessary columns

Remove the column if that index is not necessary.

```

1 data <- data[,-x]

```

2 Check NAs (must)

First, to find out if any missing value exists.

```

1 anyNA(data)

```

Next, to clarify where exactly are they.

```

1 which(is.na(data), arr.ind = TRUE) # To see the index of NAs
2 data[apply(data, 1, function(row) { # To see the rows with NAs
3   any(is.na(row))
4 }), ]

```

Why NAs are associated with [experimental design] is unknown and worth discussion. But since we did not obtain sufficient evidence from the question to retain these abnormal values, we can just drop them.

```

1 data.noNA = data[complete.cases(data), ] # One way
2 df_after <- df_before %>% drop_na() # Another way

```

3 Check duplications (must)

First, to see if any duplicated entries exist and where are they.

DATA IMPORT, EXPLORATION & CLEANING

```
1 data[duplicated(data) | duplicated(data, fromLast = TRUE), ]  
2 which(duplicated(data) | duplicated(data, fromLast = TRUE))
```

Since duplicated entries are evident mistakes, they must be removed from the table.

```
1 data <- data[!duplicated(data), ]
```

4 Check Typos (must)

```
1 table(data$col)
```

That appears to be a typo, so next just zoom in and confirm.

```
1 data %>%  
2 filter(col == "A weird value")
```

A clear typo confirmed. Need to correct it.

```
1 data.noTypo <- data  
2 data.noTypo[data.noTypo$X== row_number, "A weird value"] <- "Corrected value"  
3 table(data.noTypo$col)
```

5 Change long/wide data format

Long format					Wide format				
ID	Treatment	Measurement	Glucose	Comment	ID	Treatment	Comment	Glucose_before	Glucose_after
1	Vehicle	Glucose_before	11.51		1	Vehicle		11.51	12.24
2	Vehicle	Glucose_before	10.95	Died	2	Vehicle	Died	10.95	NA
3	Vehicle	Glucose_before	9.54		3	Vehicle		9.54	7.98

```
1 # Option 1  
2 df_wide <- df_long %>%  
3 spread(key = Measurement, value = Glucose)  
  
4 # Option 2  
5 df_wide <- df_long %>%  
6 pivot_wider(  
7   names_from = Measurement, # Column name  
8   values_from = Glucose # Values of column  
9 )
```

Wide format					Long format				
ID	Treatment	Comment	Glucose_before	Glucose_after	ID	Treatment	Measurement	Glucose	Comment
1	Vehicle		11.51	12.24	1	Vehicle	Glucose_before	11.51	
2	Vehicle	Died	10.95	NA	2	Vehicle	Glucose_before	10.95	Died
3	Vehicle		9.54	7.98	3	Vehicle	Glucose_before	9.54	

```
1 # Option 1  
2 df_long <- df_wide %>%  
3 gather(key = "Measurement", value = "Glucose",  
4         Glucose_before, Glucose_after)  
  
5 # Option 2  
6 df_long <- df_wide %>%
```

DATA IMPORT, EXPLORATION & CLEANING

```
7   pivot_longer(  
8     cols = starts_with("Glucose"),  
9     names_to = "Measurement",  
10    values_to = "Glucose"  
11  )
```

6 Rename

Set column name

```
1 colnames(df)[3] <- "new_column_name" # option 1  
  
2 df <- df %>% # option 2  
3   rename(NewName1 = OldName1,  
4         NewName2 = OldName2)
```

Set row name

```
1 increasing_rownames <- as.character(1:nrow(df1))  
2 rownames(df1) <- increasing_rownames  
3 # or  
4 rownames(df)[3] <- "new_row_name"
```

7 Convert data type

Especially check independent variable must be factor.

```
1 str(data)
```

The [data's column] is [old data type], should be convert to [new data type].

```
1 data$col1 <- as.factor(data$col1)  
  
2 # Organize more than 1 columns  
3 data <- data %>%  
4   mutate(col2 = factor(col1, levels = c(0.5, 1, 2), ordered = T),  
5         col1 = as.factor(col1)) %>%  
6   relocate(col1, col2)
```

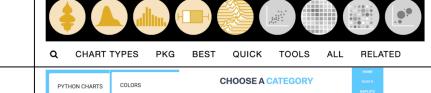
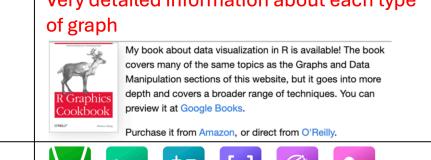
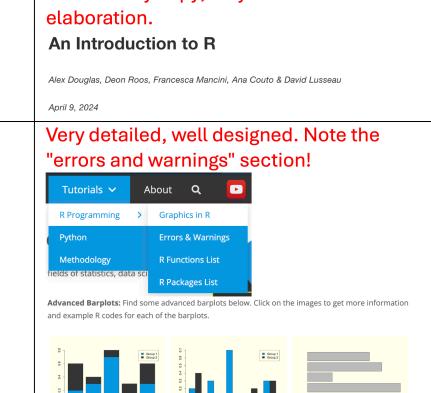
8 Check outliers

```
1 plot(density(data$col),  
2       main = "Check outliers for [data]'s [col]",  
3       lwd = 3)  
4 lines(density(data$col),  
5        col = "coral2",  
6        lwd = 3)
```

Now data is clean

Visualization

Variable type	Suggested plots
Categorical X – With unpaired continuous y – With paired continuous y – With categorical y	See below Boxplot, Whisker plot (+SD/SE/CI/IQR), Strip chart Boxplot (differences), Strip chart with linked points, Histogram of differences Bar chart, Table (with percentages), Pie chart
Continuous X	2D Scatter plot, Trend line

Stat + Plot	https://www.statology.org/tutorials/	Different types of statistical operations and plottings: explanations with their implications in R.
R Operations	https://www.statology.org/r-guides/	R operation walk through
Stat	https://statisticsbyjim.com/	
Plot	https://r-graph-gallery.com/	
Plot	https://www.r-charts.com	
R Markdown	https://bookdown.org/yihui/bookdown/	
R Operations	https://www.cookbook-r.com/	
R Operations + Stat	https://rc2e.com/	VERY VERY VERY DETAILED VERY DETAILED FOR MARKDOWN, STATS, OPERATIONS... R Cookbook, 2nd Edition James (JD) Long Paul Teator 2019-09-26 
Plot	https://r-graphics.org/	Very detailed information about each type of graph 
Math	https://www.desmos.com	
R Operation + R Markdown	https://intro2r.com/	Can directly copy; very detailed markdown elaboration. An Introduction to R Alex Douglas, Deon Roos, Francesca Mancini, Ana Couto & David Lusseau April 9, 2024
Plot R Package R Function	https://www.statisticsglobe.com	Very detailed, well designed. Note the "errors and warnings" section! 
R Packages	https://www.tidyverse.org/	

T-test

Run a t-test

Formulate hypotheses

- **Reminder:**
 - Hypothesis depends on tail and if they are paired!
 - Hypothesis depends on the question!
 - Same batch of cell cultures are paired two-samples. Not for mice.
- **One-tailed**
 - H_0 : A is not greater (or smaller) than B
 - H_A : A is greater (or smaller) than B
- **Two-tailed**
 - H_0 : The difference between A and B is 0
 - H_A : The difference between A and B is **not** 0
- **Paired two-samples (= one sample)**
 - H_0 : The difference between both conditions is 0
 - H_A : The difference between both conditions is **not** 0

State assumptions

The below are for **choosing**:

1. The test units are chosen randomly;
2. The dependent variable is continuous; (3. The mean and standard error are independent.)
3. For **unpaired t-test**, samples must be independent;
4. For **paired t-test**, each pair of measurements must be independent;
5. For **unpaired t-test**, normality must be satisfied.
6. For **paired t-test**, the differences between different conditions must be normally distributed.
7. (For **two-samples** t-test - Since two samples are involved, the homogeneity of variances should be checked.) - **Don't write it.**

Transformation (optional)

If one sample is transformed, then the reference should also be transformed.

Log transformation

```
1 # + 1 if data includes 0
2 transformed_data <- log(data)
```

Square transformation

```
1 # + 1 if data includes 0
2 transformed_data <- sqrt(data)
```

Check assumptions

- **Randomness:**
 - `sampled_data <- sample(data, size = 10, replace = F)`
 - or just explained by words in the question.
- **Continuous dependent variable:**
 - `is.numeric(data$col)`
- **Independency:**
 - Text-explained
- **Normality:**
 - **Reminders**
 - * H_0 : the data is normally distributed; H_A : the data is not normally distributed
 - * If 2 sample, check for EACH, unless they are PAIRED!
 - * Sometimes, need to declare a LACK OF OBSERVATIONS: ...thus, the histogram (and any diagnostic test as well) is difficult to judge. But it is difficult to reject H_0 of normality either. Thus, I assume that the assumption of normality of [] is satisfied.
 - Methods
 - * **Q-Q plot** `qqnorm(data$col)` with Q-Q line `qqline(data$col, col = "red", lwd = 3)`
 - The data points are evenly distributed along the `qqline`. So normality is satisfied. We choose t-test instead of non-parametric tests.
 - * **Histogram:** `hist(data$col, breaks = 10, main = "", xlab = "", ylab = "")`
 - The data distribute normally. So normality is satisfied. We choose t-test instead of non-parametric tests.
 - * **Shapiro-Wilk test:** `test_res <- shapiro.test(data$col)`

T-TEST

- $p > 0.05$: the null hypothesis of normality is not rejected (ideal).

- **Variance Homogeneity (if 2-sample)**

- F-test: `var.test(data$col1, data$col2)` or `numeric ~ factor, data`
 - * $p > 0.05$: student's t-test
 - * $p \leq 0.05$: Welch's t-test
 - p-value = [], So we use Student's t-test instead of Welch's t-test.

- **Mean and standard error independency**

- Each sample is independent and the sampling process is random, so this assumption is automatically satisfied.

Select the t-test

Questions (ask one by one)	Answer
1. Is it between 1 sample & value or 2 samples?	1 sample & a value → 1-sample test 2 samples → 2-samples test
2. If are 2-samples, are the variances equal?	Yes → Student's t-test No → Welch's t-test
3. Is the comparison direction stated?	Yes → 1-tailed No → 2-tailed
4. If are 2-samples, across different time points/conditions, is the object the same?	Yes → paired
5. Failed to satisfy normality?	No → unpaired Yes → Non-parametric test
6. For non-parametric test:	No → Student's t-test Paired/One-sample → Wilcoxon signed-rank test Unpaired → Mann-Whitney U test

- We got [two] groups.
- They are [not] paired.
- The direction of comparison is [not] stated.

So we use []-tail []-samples []paired [Student's] t-test/non-parametric test.

Run a t-test

- **Student's t-test 1**

```
1 sd(x)
2 sd(y)
```

```
1 t.test(x,
2     y = NULL,
3     alternative = c("two.sided", "less", "greater"),
4     mu = 0,
5     paired = F,
6     var.equal = F
7 )
```

- **Student's t-test 2**

```
1 sd(x)
2 sd(y)
```

```
1 t.test(DN ~ Condition, apoptosis, paired = T)
```

Compare whether the mean of DN is significantly different under two different Conditions.

- **Non-parametric test:** same as t-test other than the function names.

- Wilcoxon signed-rank test
 - * `wilcox.test(paired = TRUE)`
- Mann-Whitney U test
 - * `wilcox.test(paired = FALSE)`

Report the findings

- **Any unpaired test**

- There was a [increase/decrease] of alcohol consumed in the week after the end of semester ($M = 8.7$, $SD = 3.1$) compared to the week before the end of semester ($M = 3.2$, $SD = 1.5$), $t(52) = 4.8$, $p < 0.001$.
- But that was [not] significant. So we can/cannot reject the null hypothesis.

T-TEST

- Note: $t() = dof$
- **Paired test**
 - The average difference between the treatment condition is 7.82% ($P_{(t)} = 0.1$). Thus, we cannot reject H_0 .

ANOVA

Overview

- **Purpose:** compare means of more than two groups
 - more than 2 populations
 - * life expectancy of China, Canada and Japan
 - more than 1 predictor
 - * effect of diet (1_{st} predictor) and exercise (2_{nd} predictor) on weight
 - exactly 2 populations
 - * one-way ANOVA = independent two-sample t-test
 - * $F = t^2$
- **F-statistic** = $\frac{VarianceBetweenGroups}{VarianceWithinGroups}$
 - the higher, the more likely to reject H_0
- **Types (focused)**
 - 1-way ANOVA: 1 factor
 - * e.g. effect of 3 doses of drug A on heart rate
 - 2-way ANOVA: 2 factors
 - * e.g. effect of age and sex on salary
 - * Interactions (3 p-values):
 - Does the first factor influence the outcome?
 - Does the second factor influence the outcome?
 - Is there an interaction between both factors?

1-way ANOVA by simulation

See practical notes for details.

1-way ANOVA by statistical tests

Visualization

```

1 ggplot(data, aes(x = GroupLabelColumn, y = ValueColumn)) +
2   geom_boxplot() +
3   geom_jitter(width = 0.1, alpha = 0.5) +
4   labs(title = "", x = "", y = "") +
5   theme_minimal()

```

Justify statistical choice

- We want to compare the means of more than 2 groups
- There is 1 factor to compare

So we should try to use a 1-way ANOVA (parametric, use mean) if the data fit the requirements or we can run Kruskal-Wallis test (rank, use median).

Construct ANOVA model

```

1 model <- aov(Value ~ GroupingFactor, data = data)

```

Test assumptions

We still need to test assumptions first!

- **Independent random sampling:** from the given description of the experiment, we believe this assumption is met.
- **Normality of residuals:**
 - Q-Q plot: `plot(model, 2)`. Dots should be aligned along diagonal
 - Histogram: `hist(resid(model), main = 'residuals')`. Normality should be directly seen
 - Shapiro-Wilk test: `shapiro.test(resid(model))`. If p-value > 0.05, the null hypothesis of normality is not rejected (ideal).
- **Equality of variance:**
 - Residuals vs fitted plot: `plot(model, 1)`. Dots should be randomly scattered, similar heights of 'columns' should be seen.

Formulate hypotheses

- H_0 : there is no difference in means between `[grouping factor name]` groups
- H_A : at least one `[grouping factor name]` group mean is different from the others

ANOVA

Perform the actual ANOVA

```
1 summary(ANOVA)
```

Post-hoc test (*required even if not asked*)

```
1 TukeyHSD(model)
```

Visualize the post-hoc results

```
1 TukeyHSD(model) %>%
2   as.data.frame() %>%
3   ggplot(aes(x = diff, y = reorder(comparison, diff))) +
4   geom_point() +
5   geom_errorbarh(aes(xmin = lwr, xmax = upr), height = 0.2) +
6   labs(title = "Tukey's HSD Post-hoc Test",
7        x = "Difference in Means",
8        y = "Comparison") +
9   theme_minimal()
```

Present and interpret the results

Subjects treated with [intervention name or dose] showed a [increase/decrease] in [response variable] compared to before treatment, with a mean change of [effect size in units]. However, this difference [did/did not] reach statistical significance (red{p = [p-value]}). There was [no effect / minimal effect / mixed effect] observed at the [lower/higher/alternate] dose compared with the control group.

Despite the lack of statistical significance, the effect size measures indicate a [small/moderate/large] effect:

- Eta squared (η^2) = [value], suggesting that [X]% of the variability is explained by the treatment (small: 0.01, medium: 0.06, large: 0.14);
- Cohen's f = [value], which corresponds to a [small/medium/large] effect according to conventional thresholds (small: 0.1, medium: 0.25, large: 0.4).

2-way ANOVA by statistical tests

Question: explore the effects of attendance and previous grades on course performance

Visualization

```
1 ggplot(data, aes(x = Attendance, y = Performance, color = PreviousGrades)) +
2   geom_boxplot() +
3   geom_jitter(width = 0.1, alpha = 0.5) +
4   labs(title = "Course Performance by Attendance and Previous Grades",
5        x = "Attendance",
6        y = "Performance") +
7   theme_minimal()
```

Justify statistical choice

- There is [] groups to compare
- We are looking at the effects of more than 1 predictors

So we should try to use a 2-way ANOVA if the data fit the requirements. I see no reason to run ANOVA without interactions, so I would like to have a model like Course performance = Attendance + PreviousGrades + interaction.

Construct ANOVA model

```
1 model <- aov(weight_gain ~ genotype * diet, data = mouse) # 2 factors, with interaction
2 model <- aov(weight_gain ~ genotype + diet, data = mouse) # 2 factors, without interaction
```

Test assumptions

We still need to test assumptions first!

- Independent random sampling: from the given description of the experiment, we believe this assumption is met.
- Normality of residuals:
 - Q-Q plot: `plot(model, 2)`. Dots should be aligned along diagonal

ANOVA

- Histogram: `hist(resid(model), main = 'residuals')`. Normality should be directly seen
- Shapiro-Wilk test: `shapiro.test(resid(model))`. If p-value > 0.05, the null hypothesis of normality is not rejected (ideal).
- **Equality of variance:**
 - Residuals vs fitted plot: `plot(model, 1)`. Dots should be randomly scattered, similar heights of ‘columns’ should be seen.
- **Equality of group size:**
 - The group size can be noticed in the data diagnosis step.

So we use *parametric* ANOVA.

Formulate hypotheses

- Hypotheses set 1 for *main effect*:
 - H_0 : the mean of [y] of [x1] and [x2] is the same
 - H_A : the mean of [y] of [x1] and [x2] is not the same
- Hypotheses set 2 for *interactions*:
 - H_0 : there is no interaction between [x1] and [x2]
 - H_A : there is an interaction between [x1] and [x2]

Perform the actual ANOVA

```
1 summary(model)
```

Post-hoc test (*required even if not asked*)

```
1 TukeyHSD(model)
```

Visualize the post-hoc results

```
1 TukeyHSD(model) %>%
2   as.data.frame() %>%
3   ggplot(aes(x = diff, y = reorder(comparison, diff))) +
4   geom_point() +
5   geom_errorbarh(aes(xmin = lwr, xmax = upr), height = 0.2) +
6   labs(title = "Tukey's HSD Post-hoc Test",
7        x = "Difference in Means",
8        y = "Comparison") +
9   theme_minimal()
```

Present and interpret the results

- Overall statement: almost all terms show significant differences. This indicates that the [independent factors] have distinct effects on [response variable];
- Main effect A: treatment 1 is generally better than treatment 2 ($\sim \text{diff}$, $p < 0.05$);
- Main effect B: treatment I is generally better than treatment II ($\sim \text{diff}$, $p < 0.05$);
- Interaction: the interaction between Factor A and Factor B is also statistically detectable ($p = \text{interaction p-value}$); Or: the interaction term was not significant ($p = \text{interaction p-value}$), indicating that the effect of Factor A does not depend on the level of Factor B.
- Conclusion: altogether, these results suggest that [interpretation: e.g., the new formulation can substitute the natural source only under certain conditions].

Power Analysis

Overview

- **Definition:** statistical power is the likelihood that a study will detect an effect when THERE IS an effect there to be detected
- **Scenario:** Lack of statistical significance does not prove that there is no difference, instead, it may be a consequence of low power
- **Power requirements**
- **Effect factors**
 - Effect size $\Delta\mu = \mu_1 - \mu_0$: Power $\propto \frac{|\mu_1 - \mu_0|}{\sigma}$
 - Standard deviation σ : higher $\sigma \Rightarrow$ wider distributions \Rightarrow more overlap \Rightarrow lower power
 - **Sample size** n : $SE = \frac{\sigma}{\sqrt{n}}$ \Rightarrow larger $n \Rightarrow$ smaller SE \Rightarrow higher power
 - Desired significance level α : higher $\alpha \Rightarrow$ lower $\beta \Rightarrow$ higher power

Balance α and β level

- **General rules**
 - A Type I error is considered worse so α is rarely > 0.05
 - If we can tolerate a 5% type I error, we can tolerate a 20% type II error
 - Some say $TypeIError = TypeIIError$
- **Clinical settings**
 - Large clinical trials use 0.9 or 0.95 (90-95% power, $\beta = 0.1 - 0.05$)
 - Clinical trials phase III: $min = 80$
- **Animal studies**
 - Animal studies usually use 0.8 (80% power, $\beta = 0.2$)
- **Omics studies**
 - Omics studies: aim for high power, because you want to minimize Type II error

T-test: calculate power (simulation)

- **Given information**
 - Sample size (paired): $n = 5$
 - Mean difference: $\Delta = 7.82$
 - Standard deviation: $\sigma = 8.20$
 - Significance level: $\alpha = 0.05$

```

1 ps <- replicate(1e5, t.test(rnorm(5, 7.82, 8.2), mu = 0)$p.value)
2 power <- length(which(ps < 0.05)) / 1e5
3 power

```

E.g., if the drug is indeed effective as showed on the animal model, then what is the probability that they do not see a significant effect of the drug (p-value cutoff = 0.05)?

```

1 p.10 <- replicate(1e5, t.test(rnorm(10, 117, sd = 30), mu = 130, alternative = "less")$p.value)
2 length(which(p.10 <= 0.05)) / length(p.10) * 100

```

T-test: calculate power (built-in function)

- **Given information:** same as above

```

1 power.t.test(n = 5, delta = 7.80, sd = 8.2, sig.level = 0.05, type = "paired")

```

T-test: calculate the minimum sample size for an aimed power (simulation)

- **Goal:** what n gives us 80% power?
- **Interpretation:** $power = 0.8$, so $\beta = 0.2$
- **Given information:** same as above

```

1 power.t.test(power = 0.8, delta = 7.82, sd = 8.2, sig.level = 0.05, type = "two.sample")
2 # Output: n = 18.27 (per group)
3 # next, clarify which one is best: 18 or 19
4 power.t.test(delta = 7.82, sd = 8.2, sig.level = 0.05, n = 18, type = "two.sample")
5 power.t.test(delta = 7.82, sd = 8.2, sig.level = 0.05, n = 19, type = "two.sample")

```

POWER ANALYSIS

T-test: parameter settings (both)

- Consider:

- Study design (e.g., paired or unpaired ...)
- Effective or not in reality
- See the effect or not

Simulation

A college is providing a free bottle of milk per day to the students. Based on the history of a nearby country, giving a bottle of milk to the college students could improve the average height from 175cm to 178cm. Now, you go to the college and measure the heights of 10 male students, and perform a t-test to see if the students are indeed higher than the average of the country. Simulate the data you get from the students, and record the p-value in the t.test, which is the number?

```
1 p <- replicate(1e4, t.test(rnorm(10, 178, sd = 10), mu = 175,
2                     alternative = "greater")$p.value) # Note `alternative = "greater"`
3 length(which(p <= 0.05)) / length(p) * 100 # Note `p <= 0.05`
```

If the drug is indeed effective as showed on the animal model, then what is the probability that they do not see a significant effect of the drug (p-value cutoff = 0.05)? Please perform a simulation as you did before to give an answer.

```
1 p <- replicate(1e5, t.test(rnorm(10, 117, sd = 30), mu = 130,
2                     alternative = "less")$p.value) # Note `alternative = "less"`
3 length(which(p <= 0.05)) / length(p) * 100 # Note `p <= 0.05`
```

Built-in functions

If the company truly believes the effect of the drug and want to be sure that they will not largely miss the effect in the trial (type II error rate < 0.2), then how many volunteers do they need to recruit?

```
1 power.t.test(power = 0.8, delta = 13, sd = 30, sig.level = 0.05,
2               type = "two.sample", alternative = "one.sided") # Note: use `type = "one.sided"`
```

If the company changes their strategy, asking all the volunteers to take the pills and measuring their weights before and afterward, how many volunteers do they need?

```
1 power.t.test(power = 0.8, delta = 13, sd = 30, sig.level = 0.05,
2               type = "paired", alternative = "one.sided") # Note: use `type = "paired"`
```

ANOVA: calculate power (simulation)

```
1 sample_sig_interactions <- sapply(seq(0, 100, 5), function(sample_size) {
2   mean(replicate(100, {
3     sample_mouse <- mouse[sample(nrow(mouse), sample_size), ]
4     sample_model <- aov(weight_gain ~ genotype * diet, data = sample_mouse)
5     sample_interaction_pvalue <- summary(sample_model)[[1]][3, 5]
6     sample_interaction_pvalue <= 0.05 # see the effect
7   })) * 100
8 })
9 # Get sample size closest to 80% power
10 sample_sizes <- seq(0, 100, 5)
11 min_sample_size <- sample_sizes[which.min(abs(sample_sig_interactions - 80))]
12 print(min_sample_size)
```

ANOVA: calculate the minimum sample size for an aimed power (simulation)

```
1 power_at_60 <- mean(replicate(100, {
2   sample_mouse <- mouse[sample(nrow(mouse), 60), ]
3   sample_model <- aov(weight_gain ~ genotype * diet, data = sample_mouse)
4   sample_interaction_pvalue <- summary(sample_model)[[1]][3, 5]
5   sample_interaction_pvalue <= 0.05
```

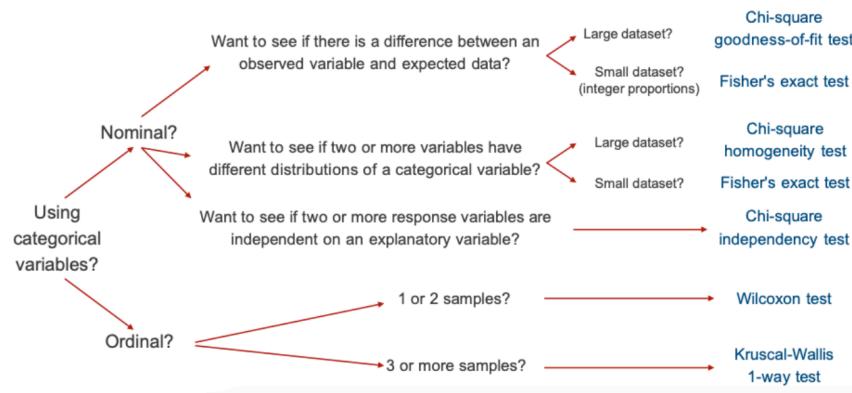
POWER ANALYSIS

```
6    })) * 100  
7    print(power_at_60)
```

For an integrated method: see week 2 practical notes.

Chi-Squared Test

Overview



Test Type	#	Purpose	Scenario
Goodness-of-Fit	1	Does the observed distribution match an expected one?	- Is there a difference between the season preferences? - Is their estimation correct?
Test of Independence	2	Are two categorical variables dependent (aka affecting each other)?	- Does gene X affect mice lifespan?
Test of Homogeneity	2	Are distributions of a categorical variable the same across groups?	- Is there a difference between the distribution of allergic reactions in different seasons?
3-way	3	3 categorical variables	- Are there relationships between gene X, sex and lifespan?
Fisher's Exact Test	x	For small sample sizes's test for independence	

Goodness-of-Fit Test

From dataframe to table

```
1 table(genotype$sex, genotype$genotype) # sex: row name
```

Visualize

```
1 t <- as.data.frame(table(genotype$sex, genotype$genotype))
2 pie(
3   t$Freq[t$Var1 == "male"],
4   labels = t$Var2[t$Var1 == "male"],
5   main   = "")
```

What would you expect?

Make a markdown table (table 1) and describe with a few sentences.

Choose and justify the appropriate statistical test

- It is an expected vs observed question.

So we use Chi-squared goodness-of-fit test.

Test Assumptions

- The variables must be categorical. *Fits, due to our examination.*
- Observations must be independent. *Fits, can assume from the task.*
- Cells in the contingency table are mutually exclusive. *Fits, due to our examination.*
- The expected value of cells should be 5 or greater in at least 80% of cells. *Fits, see table 1.*

Thus, we can run Chi-squared goodness-of-fit test safely.

Form hypotheses

- H_0 : The data follow the expected distribution (table 1).
- H_a : The data does not follow the expected distribution (table 1).

CHI-SQUARED TEST

Perform Chi-squared test

Note: values of the expected dataframe should be provided in *frequency*!

```
1 Seasons <- c(40, 30, 18, 28)
2 chisq.test(Seasons, p = rep(0.25, 4), correct = F)
```

Present and discuss the results

Clearly, [independent variable] affects the [dependent variable]. Observed objects are [percent]% ($[p >/<]$) of the expected number. Moreover, ... (if any)

Calculate p-value (not required)

```
1 pchisq(X2, dof, lower.tail = FALSE)
```

Test for homogeneity

Form hypotheses

- H_0 : the distribution of allergic reactions (a variable) is the same for people who prefer different seasons (another variable).
- H_a : the distribution of allergic reactions (a variable) is not the same for people who prefer different seasons (another variable).

Perform the test: Chi-squared or Fisher's exact test

```
1 # Two_categories: a dataframe with the column being seasons and the row being allergic reactions
2 chisq.test(Two_categories)
3 fisher.test(Two_categories)
```

Test for independence

Form hypotheses

- H_0 : mice survival (a variable) is independent of gene X (another variable).
- H_a : mice survival (a variable) is dependent of gene X (another variable).

Perform the test: Chi-squared or Fisher's exact test

```
1 # Two_categories: a dataframe with the column being gene X and the row being mice survival
2 chisq.test(Two_categories)
3 fisher.test(Two_categories)
```

3-way Chi-squared test

dof

- $dof = (\#1st\ variable)(\#2nd\ variable)(\#3rd\ variable) - 1 - (\#1st\ variable - 1) - (\#2nd\ variable - 1) - (\#3rd\ variable - 1)$

Form hypotheses

- H_0 : there is no inter-dependency between gene X, sex and lifespan
- H_a : there is an inter-dependency between gene X, sex and lifespan

Bootstrapping

Overview

- When to use:
 - Don't know population distribution
 - The *complex distribution* (e.g., *non-normal distribution*) limit the use of classical statistical methods (t-test, ANOVA...)
 - The population size is extremely small
 - Data lacks independence
 - Options of categorical data are not mutually exclusive
- What for?
 - For hypothesis testing
 - for a particular pair
 - simulate, with my own efforts, how distribution would look like if H_0 was true
 - then examine if the observed distribution is extreme (compare p and α)
 - For calculating confidence interval
 - not for a particular pair, multiple categories exist
 - after repeating, rank acquired values
 - get 2.5% as lower bound and 97.5% as upper bound for 95% Confidence Interval

Justify statistical test choice

- [Characteristics of this study] + [So A test and B tests are not appropriate].
- Bootstrapping is the only possible alternative to take account of these complications.

Steps

- Make a bootstrapped sample
 - Note: sampling from original data *with replacement*
- Calculate something (mean, median, sd...)
- Record the calculation
- Repeat step 1 ~ 3 for a bunch of times
 - Note: the number of bootstrapping replication will NOT change p-value systematically, but WILL narrow the p-value variation, that is, the more likely you will get a stable and accurate p-value

Examine original distribution

I first assess whether the data is normally distributed which would allow me to perform parametric tests such as a t-test.

```
1 shapiro.test(data$value)
2 hist(data$value)
```

Since [], the Shapiro-Wilks test shows that the data is normally distributed, however, the histogram disagrees. Probably the lack of significance in the Shapiro-Wilks test is due to a lack of statistical power rather than the data being truly normally distributed. To be conservative, I choose Bootstrapping.

Bootstrap for hypothesis testing

The question: is group 1 and group 2 significantly different?

Formulate hypotheses

- H_0 : there is no difference between the median of group 1 and group 2
- H_A : there is a difference between the median of group 1 and group 2

Visualization with box plot

```
1 boxplot(
2   # data$value~data$category,
3   # ylim = c(0, 10),
4   las = 2,
5   # ylab, xlab, main
6 )
```

The actual condition:

```
1 group1 <- median(subset(data, data$col == "Group1")) # might change for mean or sd with `sd()`-
2 group2 <- median(subset(data, data$col == "Group2"))
3 median_diff <- group1 - group2
4 median_diff
```

BOOTSTRAPPING

Simulate the world when H_0 is true

```
1 number_group1 <- nrow(subset(data, data$col == "Group1"))
2 number_group2 <- nrow(subset(data, data$col == "Group2"))

3 # for reproducibility
4 set.seed(123)

5 # for recording
6 bootstrap_median_record <- vector()

7 # repeats
8 for (a in 1:100) {
9   # no filtering, but with same number
10  bootstrap_group1 <- median(sample(data$col, number_group1, replace = TRUE))
11  bootstrap_group2 <- median(sample(data$col, number_group2, replace = TRUE))
12  bootstrap_median <- bootstrap_group1 - bootstrap_group2
13  bootstrap_median_record <- c(bootstrap_median_record, bootstrap_median)
14 }

15 # visualization
16 hist(bootstrap_median_record)
17 abline(v = median_diff, col = 'red') # for comparison
```

To get the p-value for reaching a definitive conclusion

```
1 length(subset(bootstrap_median_record, bootstrap_median_record >= median_diff))/100
```

- I then obtain a p-value for this test by determining how many of my bootstrap values have a more extreme value than my observed value. In this case, this length is [] which corresponds to $p=n$.
- $n > 0.05$, which prevents me from rejecting the H_0 at the standard p-value = 0.05 threshold.
- I can therefore conclude that there is no difference in []

Bootstrap for confidence interval

The question: what is the possible range for a statistic for several groups?

Formulate hypotheses

- H_0 : there is no difference in values between categories
- H_A : there is a difference in values between categories

Visualization with dot plot

```
1 plot(
2   data$value,
3   # ylab, xlab
4   xaxt = "n"
5 )

6 axis(
7   side = 1,
8   at = seq(1, nrow(data), 1),
9   labels = data$category,
10  las = 2
11 )
```

Repeat the bootstrapping for the entire dataset (also a duo-comparison below).

```
1 total_values <- sum(data$values)
2 this_value_record <- vector()
3 lower_bound_record <- vector()
4 upper_bound_record <- vector()

5 for (row in 1:nrow(data)) {
```

BOOTSTRAPPING

```
6 # define
7 this_category <- data[row, 1]
8 this_value <- data[row, 2]
9 not_this_value <- total_values - this_value
10
11 # create a new dataframe with:
12 # category repeat for the value times 8
13 # categorize into only two classes
13 re-classified <- c(rep(this_category, this_value),
14                         rep("not_this_category", not_this_value))
15
16 # initialize bootstrapping record
17 bootstrap_record <- vector()
18
19 for (replication_number in 1:100) {
20
21     # the actual bootstrapping
22     bootstrap_sample <- sample(re-classified, length(ore-classified), replace = T)
23
24     # write record
25     bootstrap_record <- c(bootstrap_record,
26                             length(
27                                 subset(bootstrap_sample, bootstrap_sample == this_category)
28                             )
29                         )
30
31     lower_bound <- quantile(bootstrap_record, 0.025)
32     upper_bound <- quantile(bootstrap_record, 0.975)
33
34     this_value_record <- c(this_value_record, this_value)
35     lower_cis <- c(lower_bound_record, lower_bound)
36     upper_cis <- c(upper_bound_record, upper_bound)
37 }
38 }
```

Visualize the results in a useful way

```
1 ymax <- ceiling(max(upper_bound_record) * 100) / 100
2
3 plot(
4     this_value_record,
5     xaxt = "n",
6     ylim = c(0, ymax),
7     # xlab, ylab
8     pch = "."
9 )
10
11 axis(
12     side = 1,
13     at = seq(1, nrow(data), 1),
14     labels = data$category,
15     las = 2
16 )
17
18 for (a in 1:length(lower_bound_record)) {
19     lines(x = c(a, a), y = c(lower_bound_record[a], upper_bound_record[a]))
20     lines(x = c(a - 0.1, a + 0.1),
21           y = c(lower_bound_record[a], lower_bound_record[a]))
22     lines(x = c(a - 0.1, a + 0.1),
23           y = c(upper_bound_record[a], upper_bound_record[a]))
24 }
25
26 points(
```

BOOTSTRAPPING

```
23     x = seq(1, nrow(data), 1),
24     y = this_value_record,
25     pch = 20
26 )
```

- As the confidence intervals defined here represent a [95%] certainty that the true proportion of new cases is within this interval
- I can therefore conclude that pairs of categories with non-overlapping confidence intervals are significantly different at the standard p-value threshold of [0.05].

```
1 # Initialize the values & vectors
2 first_satisfied <- 864
3 first_unsatisfied <- 714
4 second_satisfied <- 980
5 second_unsatisfied <- 473
6 first_bootstraps <- vector()
7 second_bootstraps <- vector()

8 # Create tag vectors
9 first_results <- c(rep(1, first_satisfied), rep(0, first_unsatisfied))
10 second_results <- c(rep(1, second_satisfied), rep(0, second_unsatisfied))

11 # Bootstrap for 100 times
12 for (a in 1:100) {
13   first_sample <- mean(sample(first_results, length(first_results),
14                         replace = TRUE))
15   second_sample <- mean(sample(second_results, length(second_results),
16                         replace = TRUE))
17   first_bootstraps <- c(first_bootstraps, first_sample)
18   second_bootstraps <- c(second_bootstraps, second_sample)
19 }

20 # Generate CI
21 first_upper <- quantile(first_bootstraps, probs = c(0.975))
22 second_lower <- quantile(second_bootstraps, probs = c(0.025))

23 # Gives conclusion
24 first_upper < second_lower
```

Visualization

```
1 boxplot(first_bootstraps, second_bootstraps, notch = TRUE,
2           names = c("early", "late"),
3           ylab = "Prop. of satisfied button presses")
```

Correlation and Regression

Regression (assumptions are met)

Formulating a linear model and exploring it

```
1 model <- lm(y ~ x, dataframe)
```

Check assumptions

- Independence of observations
- Linear relationship between variables
- Homoscedasticity of residuals
- Normal distribution of residuals

1. Independence of observations

If not dependent measurements such as time series: “Independence of observations can be assumed from the task itself.”

2. Linear relationship between variables

The relationship between the variables seems to be quite linear according to the scatter plot.

```
1 plot(x, y, main = " ", xlab = "Age, yrs", ylab = "Weight, kg")
```

3. Homoscedasticity of residuals

Choose one of the following:

- Visual examination: `plot(model, 1)`. Discussion: “The model residuals are not strictly homoscedastic, but the reason to reject homoscedasticity is not sufficient either.”
- Statistical test: Breusch-Pagan test with `lmtest::bptest(model)`. If p-value > 0.05, the null hypothesis of homoscedasticity is not rejected (ideal).

4. Normal distribution of residuals

Choose one of the following:

- Visual examination: `plot(model, 2)` or `hist(residuals(model))`. Discussion: The model residuals have some values that do not strictly correspond to the expected values, but the reason to reject normality is not sufficient either.
- Statistical test: Shapiro-Wilk test with `shapiro.test(residuals(model))`. If p-value > 0.05, the null hypothesis of normality is not rejected (ideal).

“Overall, the model is usable. However, some values seem to be out of place. For instance, there are really few measurements from the cars moving at a lower speed. It skews the model.”

Check for outliers

```
1 plot(cooks.distance(model), type = "h",
2      main = "Cook's Distance", ylab = "Distance")
```

Since there’s no clear reason for why the data at [data point] is abnormal, it is better to be safe than sorry.

Discuss the model

Use `str(model)` or `model` or the following (most often)

```
1 summary(model)
```

Interpretation:

- F-statistic: if our model is significant.
 - If the p-value is less than 0.05, we can reject the null hypothesis that the model is not significant.
 - So we can conclude that x is not a significant predictor of y.
- R-squared: how much of the variance in y is explained by x.
 - If R-squared is 0.5, it means that 50% of the variance in y is explained by x.

They are both indicators of the model’s goodness of fit.

The model [is/is not] statistically significant (p-value < 0.05 → significant). The variables correlate at [$r \approx 0.8$] ($R^2 = 0.65$). The model explains about [65%] of the total variance ($R^2 = 0.65$). Thus, the model is quite good. But there [is/is not] still quite a lot of unexplained variance. It would be better to account for more factors.

Visualize the model (optional)

```
1 ggplot(df, aes(x = Age, y = Weight)) +
2   geom_point(aes(shape = "Primary points"), color = "black") +
```

CORRELATION AND REGRESSION

```
3 geom_smooth(aes(color = "Model"), method = "lm", se = FALSE) +
4 scale_color_manual(values = c("Model" = "red")) +
5 scale_shape_manual(values = c("Primary points" = 1)) +
6 labs(x = "Speed, km/h", y = "Distance to stop, m") +
7 theme_minimal(base_size = 14) +
8 theme(
9   legend.title = element_blank(),
10  legend.position = "right" # bottom/top
11 )
```

Use it to predict

```
1 predict.lm(
2   object = model, # The regression model must be here
3   newdata = data.frame(Age = c(10, 12, 14, 16)) # Predict multiple at once
4 ) # Input values in a data frame
```

Perform and discuss correlation

```
1 cor.test(x, y, method = "pearson")
```

Interpretation:

- $0 \leq |r| < 0.3$ — Weak correlation
- $0.3 \leq |r| < 0.5$ — Moderate correlation
- $0.5 \leq |r| < 0.7$ — Strong correlation
- $0.7 \leq |r| \leq 1$ — Very strong correlation

Regression (assumptions NOT met)

Options:

- Transform the data: log, \log_{10} , square root
- Remove outliers: identify with `plot(model, 4)` explain and (replace)
- Non-parametric test: Spearman's correlation, non-linear models

Perform and discuss correlation

```
1 cor.test(x, y, method = "spearman")
```

Optimize a multi-variable model

Formulate start and end model

```
1 model_1 <- lm(weight ~ age, data = df)
2 model_full <- lm(weight ~ age + height + sex + diet, data = df)
```

Perform step-wise regression

```
1 step(
2   object = model_1,
3   direction = "forward",
4   scope = formula(model_full)
5 )
```

Time Series

Core

- Definition: a random variable is measured during a **prolonged period** of time in **regular intervals**
- Examples: the birth rate in a country over some time; deaths from a certain reason over some time; concentration of a certain compound (in the air/solution) over a certain time period; heart electrical activity; hormone level in the blood
- Components
 - trend m_t
 - * linear regression: $x_t = \beta_0 + \beta_1 * t$ - get β_1
 - * moving average: $m_t = \frac{1}{k} \sum x_t$ - get m_t
 - cyclicity/seasonality s_t
 - * additive model: $s_t = x_t - m_t$ - get s_t
 - error z_t
 - * calculated from equation

Construct time series object

```

1 library(stats)
2 library(forecast)

3 ts_data <- ts(orig_data,
4   start = c(year, month),
5   frequency = n
6 )

```

- $n = 12$: monthly data
- $n = 4$: quarterly data
- $n = 1$: yearly data

Visualize time series

```

1 plot(ts_data,
2   main = "Time Series Plot",
3   ylab = "Dependent variable",
4   xlab = "Time, yrs",
5   col = "blue",
6   lwd = 2
7 )

```

Data transformation (optional)

Methods:

- \log/\log_{10}
- standard practice for exponential growth data: `after <- log10(before)`
- Differencing (must integrated after modeling)
- good for random walk data (random walk to white noise)
- may not good for non-random walk data

Check assumptions

- Absence of missing values
- Normality of residuals
- Stationarity of the time series
- Independence of residuals
- Observations made in regular intervals

Navigate the time series component

Is there a trend?

- Fit
- Plot
- Quantify uncertainty

1. Fit the model

TIME SERIES

```
1 ts_data_lm <- lm(ts_data ~ time(ts_data))
2 summary(ts_data_lm)
```

2. Plot the trend (if using linear regression)

```
1 plot(
2   time(ts_data),
3   ts_data,
4   main = "Trend examination",
5   ylab = "Dependent variable",
6   xlab = "Time, yrs"
7 )
8 abline(ts_data_lm, col = "red", lwd = 2)
```

Interpretation: If p-value <0.05: Time is a statistically significant regressor.

3. Confidence interval for the trend to quantify uncertainty

```
1 # Extract coefficient and sd
2 coef_estimate <- coef(summary(ts_data_lm))["time", "Estimate"]
3 se_estimate <- coef(summary(ts_data_lm))["time", "Std. Error"]

4 # Calculate 95% CI
5 coef_estimate - qt(0.975, df = ts_data_lm$df.residual) * se_estimate # lower
6 coef_estimate + qt(0.975, df = ts_data_lm$df.residual) * se_estimate # upper
```

Is there a seasonality?

- ACF
- Decompose
- Plot

1. Let's use Autocorrelation Function (ACF) plot to check for repeating patterns.

```
1 acf(ts_data, main = "ACF Plot of ts_data")
```

2. Decompose the series (only works for regular frequency series)

```
1 ts_data_decomp <- decompose(ts_data)
```

3. Plot the decomposed series

```
1 plot(ts_data_decomp, main = "Decomposed Time Series")
```

Predictions

- Create time pin
- Predict
- Convert (optional)
- Plot

1. Create new time points

```
1 future_time <- time(ts_data)[length(ts_data)] + seq(1, 10)
```

2. Predict on log-scale in this case

```
1 future_preds_log <- predict(ts_data_lm,
2   newdata = data.frame(ts_data_log = future_time)
3 )
```

TIME SERIES

3. Convert back to original scale (optional)

```
1 future_preds <- 10^(future_preds_log)
```

4. Plot predictions

```
1 plot(uspop,
2       main = "Forecast",
3       ylab = "Dependent variable",
4       xlab = "Time, yrs"
5     )
6
6 lines(future_time, future_preds, col = "red", lwd = 2)
7 points(future_time, future_preds, col = "red", pch = 19)
```

Improve model and prediction

```
1 # Residuals from trend model
2 residuals_trend <- resid(ts_data_lm)
3
3 # Check autocorrelation
4 acf(residuals_trend)
5
5 # Fit ARIMA model to residuals
6 residuals_arima <- auto.arima(residuals_trend)
7 summary(residuals_arima)
8
8 # Forecast future trend
9 future_trend <- predict(ts_data_lm, newdata = data.frame(time.ts_data = future_time))
10
10 # Forecast residuals
11 future_residuals <- forecast(residuals_arima, h = 10)$mean
12
12 # Combine trend + residuals
13 future_log_combined <- future_trend + future_residuals
14 future_combined <- 10^(future_log_combined)
15
15 # Plot combined forecast
16 plot(ts_data,
17       main = "Improved Forecast",
18       ylab = "Dependent variable",
19       xlab = "Time, yrs"
20     )
21 lines(future_time, future_combined, col = "green", lwd = 2)
22 points(future_time, future_combined, col = "green", pch = 19)
```

Ordinary differential equations modeling

Concept overview

- Purpose: to describe **how a quantity changes over time**.
- General form: $\frac{dy}{dt} = f(y, t)$, which means: the rate of change of variable y over time t can be expressed with a function.

Tumor growth

- N : tumor size
- $r_{Logistic}$: intrinsic growth rate
- $r_{Gompertz}$: growth rate constant
- $K_{Logistic}$: carrying capacity
- $K_{Gompertz}$: asymptotic limit (interpreted as maximum tumor size)
- α : treatment effect

* r and K interpreted differently, but same use & value in equations.

Model	Equation	Characteristics
Exponential Model	$\frac{dN}{dt} = rN$	- Rapid, unchecked growth
Logistic Model	$\frac{dN}{dt} = rN(1 - \frac{N}{K})$	<ul style="list-style-type: none"> - Not realistic in the long term - Growth slows as it nears carrying capacity K - More realistic
Logistic with treatment <i>only when treatment is mentioned</i>	$\frac{dN}{dt} = rN(1 - \frac{N}{K}) - \alpha N$	<ul style="list-style-type: none"> - Adds treatment effect α - Simulates therapy scenarios - Slower growth over time - Asymmetric S-curve - Default tumor fit
Gompertz Model	$\frac{dN}{dt} = rN \ln(\frac{K}{N})$	

```

1 # Construct Logistic model
2 logistic_model <- function(t, state, parameters) {
3   with(as.list(c(state, parameters)), {
4     dN <- r * N * (1 - N / K)
5     return(list(c(dN)))
6   })
7 }

8 # Construct Gompertz model
9 gompertz_model <- function(t, state, parameters) {
10  with(as.list(c(state, parameters)), {
11    dN <- r * N * log(K / N)
12    return(list(c(dN)))
13  })
14}

15 # Common parameters
16 parameters <- c(r = 0.05, K = 1000)
17 state <- c(N = 10)
18 times <- seq(0, 200, by = 1)

19 # Solve both models
20 out_logistic <- ode(y = state, times = times, func = logistic_model,
21                      parms = parameters) %>% as.data.frame()
22 out_gompertz <- ode(y = state, times = times, func = gompertz_model,
23                      parms = parameters) %>% as.data.frame()

24 # Differentiate them
25 out_logistic$model <- "Logistic"
26 out_gompertz$model <- "Gompertz"

27 # Combine
28 out_combined <- rbind(out_logistic, out_gompertz)

29 # Plot

```

ORDINARY DIFFERENTIAL EQUATIONS MODELING

```

30 ggplot(out_combined, aes(x = time, y = N, color = model)) +
31   geom_line(linewidth = 1.2) +
32   labs(title = "Comparison of Tumour Growth Models", x = "Time", y = "Tumour Size") +
33   theme_minimal()

```

Interpretation:

- The Gompertz model reaches its inflection point earlier, and it predicts a slower, more gradual approach to the carrying capacity than the Logistic model.
- The Gompertz model is better because it accounts for the decelerating growth seen in tumors as they experience environmental constraints like lack of nutrients, hypoxia, or immune response. Logistic growth is symmetric and may overestimate growth in later stages.

Bacteria growth

- B : #bacteria
- r : replication rate
- K : carrying capacity

Model	Equation	Characteristics
Logistic model	$\frac{dB}{dt} = rB \left(1 - \frac{B}{K}\right)$	- Used in lab environments - Growth limited by nutrients, space, etc.

```

1 # Construct growth model for later use
2 # t@state: variable
3 # parameters: parameter
4 bacteria_model <- function(t, state, parameters) {
5   with(as.list(c(state, parameters)), { # necessary
6     dB <- r * B * (1 - B / K)
7     return(list(c(dB))) # necessary
8   })
9 }

10 # Set variables and parameters' values
11 parameters <- c(r = 0.4, K = 500)
12 state <- c(B = 5)
13 times <- seq(0, 50, by = 0.5)

14 # Call ODE to embed the model
15 out <- ode(y = state, times = times, func = bacteria_model, parms = parameters)
16 out <- as.data.frame(out)

17 # Plot
18 ggplot(out, aes(x = time, y = B)) +
19   geom_line(color = "cyan", linewidth = 1) +
20   labs(title = "Bacterial Growth", x = "Time", y = "Bacterial Population") +
21   theme_minimal()

```

Full SIR model - non Markovian, continuous

- β : transmission rate
- γ : recovery rate
- $R_0 = \frac{\beta}{\gamma}$: the basic reproduction number ($R_0 > 1$: outbreak spreads, $R_0 < 1$: outbreak dies out)

Component	Equation	Interpretation
#%Susceptible (S)	$\frac{dS}{dt} = -\beta SI$	- Individuals who CAN catch the disease (but still healthy)
#%Infected (I)	$\frac{dI}{dt} = \beta SI - \gamma I$	- Individuals who HAVE the disease and thus have the ability to transmit it - Calculated as new infections minus recoveries
#%Recovered (R)	$\frac{dR}{dt} = \gamma I$	- Individuals who have recovered (OR DIED) and are NO LONGER infectious

```

1 # Construct SIR model for later use
2 sir_model <- function(t, state, parameters) {
3   with(as.list(c(state, parameters)), {
4     dS <- -beta * S * I
5     dI <- beta * S * I - gamma * I
6     dR <- gamma * I
7     return(list(c(dS, dI, dR)))
8   })
9 }

10 # Set variables and parameters' values
11 parameters <- c(beta = 0.3, gamma = 0.1)
12 state <- c(S = 0.99, I = 0.01, R = 0.0)
13 times <- seq(0, 160, by = 1)

14 # Call ODE to embed the model
15 out <- ode(y = state, times = times, func = sir_model, parms = parameters)
16 out <- as.data.frame(out)
17 out_long <- pivot_longer(out, cols = -time, names_to = "Compartment", values_to = "Proportion")

18 # Plot the results
19 ggplot(out_long, aes(x = time, y = Proportion, color = Compartment)) +
20   geom_line(linewidth = 1) +
21   labs(title = "SIR Model Simulation", x = "Time", y = "Proportion of Population") +
22   theme_minimal()

```

Interpretations:

- The number of infected individuals peak when the rate of new infections equals the rate of recoveries, that is, when $\frac{dI}{dt} = 0$.
- Increasing γ flattens the curve, reduce peak, make the end faster.

Matrix-based modeling

Concept overview

- **Matrix-based modeling:** use matrix/-cies to represent relationships between variables
 - **Transition matrix (P):** express the relationship between variables, e.g., $P = \begin{pmatrix} 0.8 & 0.2 \\ 0.1 & 0.9 \end{pmatrix}$ (columns represent transitions *from* states A and B, rows represent transitions *to* states A and B)
 - **Initial state (S_0):** initial value of all variables in the system, e.g., $S_0 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$
 - **After/on the point of n transition (S_n):** value of all variables in the system after n transitions, e.g., $S_1 = \begin{pmatrix} 0.75 \\ 0.25 \end{pmatrix}$
 - **Transition matrix equation:** equation expressed by matrix-vector multiplication as an alternative for “tree diagram thinking”. $S_1 = P \cdot S_0 \rightarrow S_2 = P \cdot S_1 = P^2 \cdot S_0 \rightarrow \text{Generalization: } S_n = P^n \cdot S_0$ for $n \geq 0$
- **Markov chains:** process moves between states with probabilities and the next sole depends on current
- **Difference equations**
 - Note: *difference equations* and *differential equations* are different ways to model a (potentially the same) question, with different ways to approach *time*

Difference Equation		Differential Equation
Purpose	Model the difference between consecutive steps	Quantify instantaneous change rate
Expression	$x_{t+1} = Ax_t$ or $x_{t+1} - x_t = f(x_t)$	$\frac{dx}{dt} = f(x(t))$
Time	Discrete time steps ($t = 0, 1, 2, 3, \dots$)	Continuous time flow
Application	Population growth, disease spread, blood flow...	See above

- Division of difference equation(DE)

Description		Equation & example
Simple DE	Single variable	<ul style="list-style-type: none"> • $P_{t+1} = A \cdot P_t$ where $A = (1 + \gamma)$ • P_t is the population size • γ is the growth rate • $\begin{bmatrix} x_{t+1} \\ y_{t+1} \end{bmatrix} = A \cdot \begin{bmatrix} x_t \\ y_t \end{bmatrix}$ • A is a matrix. It captures the interactions • Two compartment DE system: <ul style="list-style-type: none"> – $x_{t+1} = 0.8x_t + 0.2y_t$ – $y_{t+1} = 0.1x_t + 0.9y_t$ – initial condition: $x_0 = 100, y_0 = 50$
System of DE	Multiple interacting variables forms a system	

SIR early-stage model

Early epidemic phase is often modeled with linearized SIR model, since $I(t)$ is small; $S(t) \approx S_0$ is constant; and $R(t) \approx 0$. Thus, $\frac{dI}{dt} = (\beta S_0 - \gamma)I$ since $\frac{dI}{dt} = \beta S(t)I - \gamma I$ and $(\beta S_0 - \gamma)$ is a constant, same role as r .

Single population SIR:

- $\frac{d}{dt}[I] = (\beta S_0 - \gamma)[I]$ as it is

```

1 # Parameters
2 beta <- 0.9
3 gamma <- 0.2
4 S0 <- 0.90
5 r <- beta * S0 - gamma

6 # Initial infected proportion
7 I0 <- 0.10

8 # Time settings
9 time_steps <- 10 # total time course
10 infected <- numeric(time_steps)
11 infected[1] <- I0 # initialize a vector

12 # linearised infection growth
13 for (t in 2:time_steps) {
14   infected[t] <- (1 + r) * infected[t-1] # use the formula

```

```

15  }

16  # Susceptible stays constant, no recovery yet
17  # So same throughout the time course
18 susceptible <- rep(S0, time_steps)
19 recovered <- rep(0, time_steps)

20 # Plot
21 plot(1:time_steps, infected, type = "b", col = "red",
22       xlab = "Time step", ylab = "Proportion",
23       ylim = c(0, 1),
24       main = "linearised SIR: Infected over Time")
25 lines(1:time_steps, susceptible, type = "l", col = "blue")
26 lines(1:time_steps, recovered, type = "l", col = "green")
27 legend("topright", legend = c("Infected (I)", "Susceptible (S)", "Recovered (R)"),
28        col = c("red", "blue", "green"), lty = 1, bty = "n")

```

Multiple population SIR:

- Variable I is not singular; parameter β and γ is not singular

$$\vec{I} = \begin{bmatrix} I_A \\ I_B \end{bmatrix}$$

- Transmission rate β table:
- Recovery rate γ : γ_A, γ_B

From / To	Infecting A	Infecting B
A	β_{AA}	β_{AB}
B	β_{BA}	β_{BB}

- Initial state is also not singular, so the constant $\beta S_0 - \gamma$ is not singular either, thus it should be expressed as a matrix:

$$A = \begin{bmatrix} \beta_{AA} S_A - \gamma_A & \beta_{BA} S_B \\ \beta_{AB} S_A & \beta_{BB} S_B - \gamma_B \end{bmatrix}$$

- Complete equation: $\frac{d}{dt} \begin{bmatrix} I_A \\ I_B \end{bmatrix} = A \cdot \begin{bmatrix} I_A \\ I_B \end{bmatrix}$

- $A\vec{v} = \lambda\vec{v}$

- Eigenvalue λ : the sign and size of the *dominant (largest in magnitude) eigenvalue* informs to what extent the epidemic is growing or fading
- If $|\lambda_{\max}| < 1$: infection dies out
- If $|\lambda_{\max}| > 1$: infection spreads
- If $|\lambda_{\max}| = 1$: infection is in equilibrium
- Eigenvector \vec{v} : shows the mode of disease spread across groups
 - * If $A < B$: B is the main infected group, dominant the spread direction
 - * If $A > B$: A is the main infected group, dominant the spread direction

Summary:

	Single-variable model	Multi-group model
Expression	$\frac{dI}{dt} = rI$	$\frac{d\vec{I}}{dt} = A\vec{I}$
Solution	$I(t) = I_0 e^{rt}$	$\vec{I}(t) = \sum c_i e^{\lambda_i t} \vec{v}_i$
Growth?	$r > 0$	$\lambda_{\max} > 0$
Spreading direction	No directionality	Depends on eigenvector \vec{v}

Full SIR - Markovian, discrete

- Expressed in Markovian form: $x_{t+1} = P^T x_t \mid x_t = (P^T)^t x_0$ or $x_{t+1} = P x_t \mid x_t = P^t x_0$
 - When $P * x = x$ or $P^T * x = x$
 - * Most people recover, no new infections
 - * The system is in equilibrium, disease is ended

Interpretation of the matrix:

$$\begin{bmatrix} 0.80 & 0.20 & 0.00 \\ 0.00 & 0.70 & 0.30 \\ 0.00 & 0.00 & 1.00 \end{bmatrix}$$

- Each row represent probabilities of transitioning between states, where row 1 represents transition starts with the susceptible (1_{st} column), to the infected (2_{nd} column) and directly jump to the recovered (3_{rd} column).

MATRIX-BASED MODELING

- The matrix is not meaningful to read as column.
- The sum of each row is 1 because we have to ensure conservation of probability.

Construct matrix

```
1 P <- matrix(c(
2   0.80, 0.20, 0.00,
3   0.00, 0.70, 0.30,
4   0.00, 0.00, 1.00
5 ), nrow = 3, byrow = TRUE)
6 P <- t(P) # very important!!!
```

State distribution after 1 time step

```
1 # Initial state vector: conditions at time step 0
2 state_vec <- c(0.99, 0.01, 0.00) # The sum is 100% of the population
3 state_after_one <- P %*% state_vec
4 print(state_after_one)
```

The proportion of infected individuals [increases slightly] after one time step.

Steady-state distribution

```
1 # Initial state vector
2 state_vec <- c(0.99, 0.01, 0.00)
3 steps <- 150
4 states_vec_over_time <- matrix(0, nrow = steps, ncol = 3)
5 states_vec_over_time[1, ] <- state_vec

6 for (t in 2:steps) {
7   states_vec_over_time[t, ] <- P %*% states_vec_over_time[t-1, ]
8 }

9 # Plot
10 matplot(1:steps, states_vec_over_time, type = "l", lty = 1, col = c("blue", "red", "green"),
11           xlab = "Time Step", ylab = "Proportion",
12           main = "Markov Chain Infection Modeling")
13 legend("right", legend = c("Susceptible", "Infected", "Recovered"),
14        col = c("blue", "red", "green"), lty = 1, bty = "n")

15 # Calculate state at time step 20 using X0 %*% P^20
16 state_at_10 <- state_vec %*% (P %^% 10)
```

Steady-state largest eigenvalues and eigenvectors

$$\vec{X} = P^t \cdot \vec{X}$$

```
1 res <- eigen(P)
2 res$values
3 res$vectors[, 1] # first eigenvector
```

So the steady state distribution is given by the first eigenvector, which corresponds to the largest eigenvalue (1 in this case). Actually in this scenario, as long as the R in SIR equals 1, the steady state distribution will always be $(0, 0, 1)$, regardless of the transition probabilities in the matrix.

Interpretation:

- Changing the transition probabilities in the matrix affects
 - How quickly individuals move from Susceptible → Infected → Recovered
 - The steady-state distribution
- But the proportion of infected individuals eventually stabilizes, indicating the steady-state condition.

Vessel blood flow model

- **Braching point:** several flow in & out (driven by *pressure difference*) equation. Poiseuille's law (analogy to Ohm's law)
 - $Q = \frac{\Delta p}{R}$
 - * Q : flow
 - * Δp : pressure difference = $p_{in} - p_{out}$
 - * R : resistance
 - e.g.,
 - * $p_1 - p_2 = R_1 Q_1$
 - * $p_2 - p_3 = R_2 Q_2$
 - * $p_2 - p_4 = R_3 Q_3$
- **Vessel:** one mass conservation equation
 - $Q_{Ins} = Q_{Outs}$
 - e.g., $Q_1 = Q_2 + Q_3$
- **Taken together:**
 - $p_2 + R_1 Q_1 = p_1$
 - $p_2 - R_2 Q_2 = p_3$
 - $p_2 - R_3 Q_3 = p_4$
 - $Q_1 - Q_2 - Q_3 = 0$
- **Matrix form:** $A\vec{x} = \vec{b}$
 - A : matrix of coefficients; \vec{x} : vector of unknowns; \vec{b} : vector of constants
 - * $A = \begin{pmatrix} 1 & 0 & 0 & 0 \\ -R_1 & 1 & -1 & 0 \\ 0 & -R_2 & 1 & 0 \\ 0 & -R_3 & 0 & 1 \end{pmatrix}$
 - * $\vec{x} = \begin{pmatrix} p_1 \\ p_2 \\ p_3 \\ p_4 \end{pmatrix}$
 - * $\vec{b} = \begin{pmatrix} Q_1 \\ Q_2 \\ Q_3 \\ 0 \end{pmatrix}$
 - **Solve:** $\vec{x} = A^{-1}\vec{b}$
- **Flow:** mass conservation equation
 - $Q = \frac{dV}{dt}$
 - * Q : flow
 - * V : volume
 - * t : time

Code implementation:

$$\vec{x} = A^{-1}\vec{b} \text{ now is } A \cdot \begin{bmatrix} p_2 \\ Q_1 \\ Q_2 \\ Q_3 \end{bmatrix} = \begin{bmatrix} p_1 \\ p_3 \\ p_4 \\ 0 \end{bmatrix}$$

```

1  ==== Step 1: Define pressure and resistance (known parameters) ===#
2  # Known pressures (in mmHg)
3  p1 <- 10 # Node 1 (inlet)
4  p3 <- 4 # Node 3 (outlet)
5  p4 <- 4.5 # Node 4 (outlet)
6  # p2 unknown!!!

7  # Resistances (arbitrary units)
8  R1 <- 2 # 1 to 2
9  R2 <- 4 # 2 to 3
10 R3 <- 3 # 2 to 4

11 ==== Step 2: Set up equation system ===#
12 # Coefficient matrix A
13 A <- matrix(
14   c(
15     1, R1, 0, 0,    # (p1 - p2) = R1 * Q1
16     1, 0, -R2, 0,   # (p2 - p3) = R2 * Q2
17     1, 0, 0, -R3,  # (p2 - p4) = R3 * Q3
18     0, 1, -1, -1   # Q1 = Q2 + Q3

```

```

19      ), nrow=4, byrow=TRUE)

20  # Known parameters grouped into vector b
21  b <- matrix(c(
22    p1,
23    p3,
24    p4,
25    0
26 ), nrow=4, byrow=TRUE)

27  # Solve the linear system for [p2, Q1, Q2, Q3]
28  solution <- solve(A, b)
29  solution

30  # Extract values
31  p2 <- solution[1]
32  Q1 <- solution[2]
33  Q2 <- solution[3]
34  Q3 <- solution[4]

35  list(
36    p2 = p2,
37    Q1 = Q1,
38    Q2 = Q2,
39    Q3 = Q3
40  )

41  ##### Step 3: Visualize the flows ===#
42  barplot(c(Q1, Q2, Q3),
43           names.arg = c("Q1 (1->2)", "Q2 (2->3)", "Q3 (2->4)"),
44           main = "Flow Rates in 4-Node Bifurcation",
45           ylab = "Flow Rate",
46           col = "blue")

```

What about adding one more branch?

```

1  # New knowns
2  p5 <- 5
3  R4 <- 5

4  # New Coefficient Matrix A_ext (5 equations, 5 unknowns: p2, Q1, Q2, Q3, Q4)
5  A_ext <- matrix(c(
6    1, R1, 0, 0, 0,      # (p1 - p2) = R1 * Q1
7    1, 0, -R2, 0, 0,      # (p2 - p3) = R2 * Q2
8    1, 0, 0, -R3, 0,      # (p2 - p4) = R3 * Q3
9    1, 0, 0, 0, -R4,      # (p2 - p5) = R4 * Q4
10   0, 1, -1, -1, -1      # Q1 = Q2 + Q3 + Q4
11 ), nrow=5, byrow=TRUE)

12 # Right-hand side vector b_ext
13 b_ext <- matrix(c(
14   p1,
15   p3,
16   p4,
17   p5,
18   0
19 ), nrow=5, byrow=TRUE)

20 # Solve for [p2, Q1, Q2, Q3, Q4]
21 solution_ext <- solve(A_ext, b_ext)
22 solution_ext

23 p2_ext <- solution_ext[1]
24 Q1_ext <- solution_ext[2]

```

```
25 Q2_ext <- solution_ext[3]
26 Q3_ext <- solution_ext[4]
27 Q4_ext <- solution_ext[5]

28 list(
29   p2 = p2_ext,
30   Q1 = Q1_ext,
31   Q2 = Q2_ext,
32   Q3 = Q3_ext,
33   Q4 = Q4_ext
34 )

35 barplot(c(Q1_ext, Q2_ext, Q3_ext, Q4_ext),
36           names.arg = c("Q1 (1->2)", "Q2 (2->3)", "Q3 (2->4)", "Q4 (2->5)"),
37           main = "Flow Rates in 5-Node Extended Network",
38           ylab = "Flow Rate",
39           col = "green")
```

Conditional probability

Conditional probability is for two event-analysis.

Basics

Probability and Odds

- Probability: the chance that A happens
- Odds: the ratio of $\frac{\text{The chance that } A \text{ happens}}{\text{The chance that } A \text{ not happens}}$

Code implementation

Principles:

Function type	Meaning	Purpose
<code>r...</code>	Random	Generate a series of random numbers
<code>d...</code>	Density (value)	Probability or density at a <i>specific</i> value
<code>p...</code>	Cumulative	Cumulative probability up to a value
<code>q...</code>	Quantile	Value required to a defined cumulative probability

Binomial distribution and Normal distribution:

Function	Explanation
<code>rbinom(n, size, prob)</code>	Generate n random values from a binomial distribution with <code>size</code> trials and success <code>prob</code> .
<code>dbinom(x, size, prob)</code>	Calculate the probability of exactly x successes in <code>size</code> trials with success <code>prob</code> .
<code>pbinom(q, size, prob)</code>	Compute the cumulative probability of q successes in <code>size</code> trials with success probability <code>prob</code> .
<code>qbinom(p, size, prob)</code>	Find the number of successes x such that the cumulative probability is at least p (i.e., $P(X \geq x) \geq p$).
<code>rnorm(n, mean, sd)</code>	Generate n random numbers from a normal distribution with given <code>mean</code> and standard deviation <code>sd</code> .
<code>dnorm(x, mean, sd)</code>	Compute the height of the normal distribution (density) at value x with given <code>mean</code> and <code>sd</code> .
<code>pnorm(q, mean, sd)</code>	Calculate the cumulative probability that a normally distributed value is q , with given <code>mean</code> and <code>sd</code> .
<code>qnorm(p, mean, sd)</code>	Find the value x such that the cumulative probability $P(X \leq x)$ is equal to p , given the <code>mean</code> and <code>sd</code> .

Logic notations

Symbol	Description	Truth Condition
$\neg A$	“not A ”	True if A is false
$A \wedge B$ or $A \& B$	“ A and B ”	True if both A and B are true
$A \vee B$	“ A or B ”	True if at least one of A or B is true
$A \rightarrow B$	“If A then B ”	True unless A is true and B is false
$A \leftrightarrow B$	“ A if and only if B ”	True if A and B have the same truth value

Probability notations

- Conditional probability: $P(A | B)$
- Joint probability: $P(A \cap B)$ or $P(A \& B)$

Probability calculation

	A-	A+	P(B)
B-	A-B-	A+B-	Total B-
B+	A-B+	A+B+	Total B+
P(A)	Total A-	Total A+	100%

- $P(A \cap B) = P(B \cap A)$
 - $P(A \cap B) = P(A | B) \cdot P(B)$
 - $P(B \cap A) = P(B | A) \cdot P(A)$
- $P(A | B) \cdot P(B) = P(B | A) \cdot P(A)$
- $P(B) = \sum P(B_i | A_i) \cdot P(A_i)$
- Joint probability
 - Calculated using *conditional probability*: $P(A \cap B) = P(A) \cdot P(B | A)$
 - Calculated using *individual probability*: $P(A \cap B) = P(A) \cdot P(B)$

Note

1. If A and B are independent, then $P(B | A) = P(B | \neg A) = P(B)$
2. $P(A | B) \neq P(B | A)$

CONDITIONAL PROBABILITY

Tools

- **Venn plot = contingency table:** for depicting event probability landscape
- **Probability trees:** for branching problems

Conditional probability

$$\bullet \ P(A \cap B | B) = \frac{P(A \cap B)}{P(B)}$$

$$\bullet \ P(A | B) = \frac{P(A \cap B)}{P(B)}$$

$$\bullet \ P(\text{What we want to know}) = \frac{P(\text{What we also know})}{P(\text{What we know})}$$

Interpretation: the probability of an event to happen, scaled by the knowledge we already have.

The A and B can be calculated using: 1. Raw counts 2. Original probability

Bayesianism

Bayes' Theorem

Expression

Since $A \cap B = B \cap A$, so if conditional probability is used twice, it can be deducted that

- $P(A \cap B | B) = \frac{P(A \cap B | A) * P(A)}{P(B)}$
- $P(A \cap B | A) = \frac{P(A \cap B | B) * P(B)}{P(A)}$

Interpretation

Bayes' Theorem is useful because

1. It tell us that the conditional probability given that we know *one thing* about an event can be derived from knowing *the other thing* about the event. In the end there is only one event. And sometimes $A \cap B | A$ is easier to calculate than $A \cap B | B$ and *vice versa*.
2. Always, the other component of the formula does not involve conditional probability calculation or very easily calculated or is given information.
3. It can be applied to the temporal dimension of an event, that is it unveils the relationship between prior belief and posterior belief.
4. The conditional probability given that we know one thing about an event can be derived from knowing the other thing about the event.

Bayesianism vs. Frequentism

	Frequentism	Bayesianism
Probability is	Long-run calculation	A degree of belief (how you see something based on the fact in hand)
Expression	$P(X) = \frac{X}{X+X^c}$	$P(X) \propto \text{Likelihood} \cdot \text{Prior knowledge}$
Want to add new data	Start a whole new long-run	By Bayesian updating

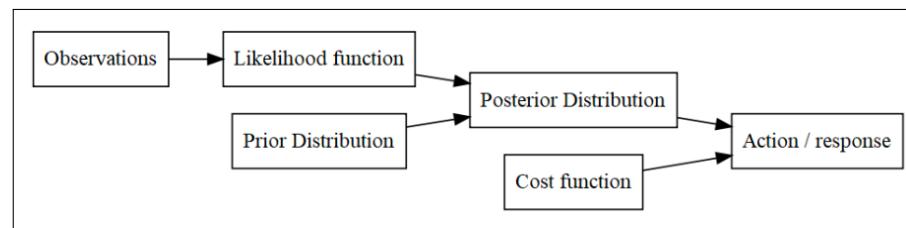
Bayesian logic components

Bayesianism inference is not a probability type that parallel with, for example, Bernoulli. It is a *way of thinking, an approach of seeing the world*.

- **Prior $P(H)$:** the *hypothesis*
- **Likelihood $P(\text{data} | H)$:** the *truth/fact*
- **Posterior $P(H | \text{data})$:** the *compromise between the hypothesis and the truth/fact (weighted multiplied average of prior and likelihood)*, the more we do the experiment, the more we know about the truth, the more impact likelihood has on the posterior
- **Marginal likelihood $P(\text{data})$:** the *total probability of the data*. This is a θ -independent constant so it does not influence the posterior shape (peak position & width).

Some use θ for prior, as a representation of prior's parameter.

If prior is not set for a new question, it is acceptable to formulate your own prior even if it is not based on the common sense; however, you should note that how you formulate prior will definitely affect the posterior accuracy.



Bayesian factor

- Purpose: to quantify and compare how strongly the data support two (and only two) competing hypotheses. So if there are too many hypotheses to compare, it is suggested to preclude some hypotheses that are not likely to be true.
- Definition: $\frac{\text{Posterior Odds}}{\text{Prior Odds}}$
- Note: not necessarily have to use. Not a standard procedure. Only used if it comes to comparing the effectiveness of models.
- Expression:

$$\text{Ratio of priors} = \text{Prior odds} = \frac{P(H_1)}{P(H_2)}, \text{ where } P(H_1) + P(H_2) + \dots + P(H_i) = 1$$

Ratio of posteriors = Posteriors odds = $\frac{P(H_1 | data)}{P(H_2 | data)}$, where $P(H_1 | data) + \dots + P(H_i | data) = 1$

$$\frac{\text{Likelihood under } H_1}{\text{Likelihood under } H_2} = BF(H_1 : H_2) = \frac{P(data | H_2)}{P(data | H_1)} = \frac{\frac{P(H_1 | data)}{P(H_2 | data)}}{\frac{P(H_1)}{P(H_2)}}$$

- Interpretation: $BF = 10$ suggest 10 times more likely to observe the data under H_1 than under H_2

A dice game

Your friend got 7 sixes in 20 rolls of a dice. Is this dice fair? What number of six do you think is most likely to be?

For this question, Bayesian is the way of thinking while Binomial/Bernoulli is the way of calculating.

Hypotheses:

- H_0 : the die is fair (precluded)
- H_1 : the die has 1 six
- H_2 : the die has 2 six
- ...
- H_6 : the die has 6 six (precluded)

```

1 # 1. Prior can be set as a uniform distribution or whatever you want
2 P_hypotheses <- rep(0.2, 5)

3 # 2. Likelihood for each hypothesis: binomial formula
4 Probabilities_sixes <- c(1 / 6, 2 / 6, 3 / 6, 4 / 6, 5 / 6)
5 P_givenData_Expectation <-
6   dbinom(x = 7, size = 20, prob = Probabilities_sixes)
7 dice <- cbind(Probabilities_sixes, P_givenData_Expectation)

8 # 3. Marginal probability
9 P_givenData <- sum(P_hypotheses * dice[, 2])

10 # 4. Posterior (for hypothesis 1)
11 (P_givenData_Expectation[1] * P_hypotheses[1]) / P_givenData

```

R as a calculator

Basic arithmetic operations

```

1 # Addition, subtraction, multiplication, division
2 6 * 7      # 42
3 20 / 5     # 4

4 # Exponentiation and modulo
5 2 ^ 3      # 8
6 10 %% 3    # 1    (remainder)
7 10 %/% 3   # 3    (integer division)

8 # Negatives and absolute value
9 -5 * 2     # -10
10 abs(-7.3) # 7.3

```

Mathematical functions

```

1 # Square root, exponentials, logarithms
2 sqrt(16)      # 4
3 exp(1)        # e^1 = 2.718282
4 log(10)       # natural log ln(10) = 2.302585
5 log10(1000)   # log base 10 of 1000 = 3
6 log(x, base = 2) # log2(x)

7 # Rounding functions
8 round(3.14159, 2) # 3.14
9 floor(3.9)        # 3
10 ceiling(3.1)     # 4
11 trunc(3.9)       # 3 (drops fractional part)

```

Vectorized operations and vector functions

```

1 # Define two vectors
2 x <- c(1, 2, 3, 4, 5)
3 y <- c(10, 20, 30, 40, 50)

4 # Common vector functions
5 sum(x)      # 15
6 prod(x)     # 120
7 mean(x)     # 3
8 median(x)   # 3
9 var(x)      # 2.5
10 sd(x)       # sqrt(2.5) = 1.5811
11 range(x)    # 1, 5
12 min(x)      # 1
13 max(x)      # 5

```

Matrix operations

```

1 # Create a 3x3 matrix (filled by column)
2 m <- matrix(1:9, nrow = 3, ncol = 3)
3 #   [,1] [,2] [,3]
4 # [1,]    1    4    7
5 # [2,]    2    5    8
6 # [3,]    3    6    9

7 # Indexing
8 m[1, 2] # 4

```

R AS A CALCULATOR

```
9  m[, 3]      # 7, 8, 9 (entire 3rd column)
10 m[2, ]       # 2, 5, 8 (entire 2nd row)

11 # Common matrix functions
12 t(m)         # transpose
13 diag(m)      # diagonal elements: 1, 5, 9
14 det(m)       # determinant
15 solve(m)     # inverse (if invertible)
16 eigen(m)     # eigenvalues and eigenvectors

17 # Matrix multiplication vs. element-wise
18 A <- matrix(c(1, 2, 3, 4), nrow = 2)
19 B <- matrix(c(5, 6, 7, 8), nrow = 2)
20 A %*% B      # matrix multiplication
21 A * B        # element-wise multiplication

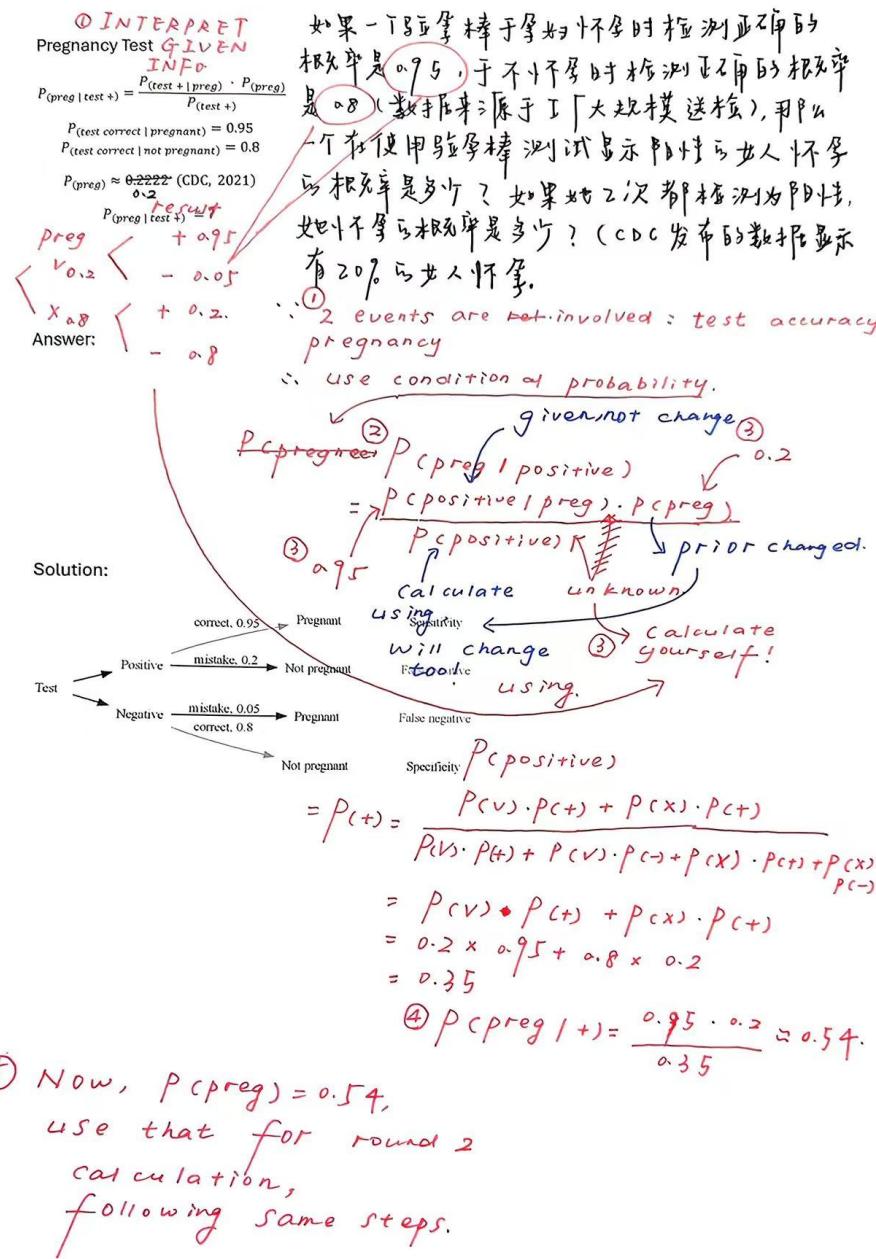
22 # Row/column summaries
23 rowSums(m)   # 12, 15, 18
24 colSums(m)   # 6, 15, 24
25 rowMeans(m)  # 4, 5, 6
26 colMeans(m)  # 2, 5, 8
```

Solving equations

```
1 # Solve a single equation, e.g., x^2 - 3x + 2 = 0
2 f <- function(x) x^2 - 3*x + 2
3 uniroot(f, lower = 0, upper = 3) # finds a root at x = 1 or x = 2
```

Probability practice problems

1 IVF test



Q: What can you say about this test kit? Is it good/specific/sensitive? Support your claims with odds and odds ratios.

- **Specific** = true negative = when test shows negative, the woman is actually not pregnant
- **Sensitive** = true positive = when test shows positive, the woman is actually pregnant
- **Good**: an overall assessment of specificity and sensitivity

Actually specificity and sensitivity cannot be assessed in this way, but in this case, we can just take it.

- **Is positive test result reliable?** Odds_(pregnant vs non-pregnant AND test⁺) = $\frac{0.211}{0.1556} \approx 1.356$
- **Is negative test result reliable?** Odds_(non-pregnant vs pregnant AND test⁻) = $\frac{0.6222}{0.0111} \approx 56.054$
- **What test result is more reliable?** Odds Ratio = $\frac{1.356}{56.054} \approx 0.024$

In conclusion:

- The chances of being correct are too low for a positive test result, that is there is too much false positive.
- The chances of being correct are too high for a negative test result (more being specific than being sensitive).
- So the overall quality is bad.

PROBABILITY PRACTICE PROBLEMS

2 Varicocele surgery

Varicocele causes infertility in men. The conception rate of patients within 1 year after surgery is around 24.7%. A new procedure emerged. 18 out of 53 patients reported conception within 1 year after that new surgery.

How likely is it that this novel procedure is better than the traditional surgery?

What would happen if the sample size is increased to 80 patients with the same proportion of the positive outcome?

```
1 # -----
2 # STEP 1: Frequentism probability
3 # -----
4 
5 # Probability of 18 conceptions out of 53 if true rate is 0.247
6 pbinom(q = 17, size = 53, prob = 0.247, lower.tail = FALSE)
7 # [1] 0.08326483
8 
9 # Larger sample case (80 patients with same proportion)
10 bigger_sample_success <- round((18 / 53) * 80)
11 pbinom(
12   q = bigger_sample_success - 1,
13   size = 80,
14   prob = 0.247,
15   lower.tail = FALSE
16 )
17 # [1] 0.04360466
18 
19 # -----
20 # STEP 2: Bayesian Inference Setup
21 # -----
22 
23 # Hypothesis space
24 model <- seq(0, 0.5, by = 0.025) # probabilities of success to test
25 
26 # Prior belief: using Beta(7, 21) as described in the slide (prior mean ~0.25)
27 prior <- dbeta(model, shape1 = 7, shape2 = 21)
28 
29 # Likelihood of observing 18 successes out of 53
30 likelihood <- dbinom(x = 18, size = 53, prob = model)
31 
32 # Posterior (unnormalized)
33 posterior_unnorm <- likelihood * prior
34 
35 # Normalized posterior
36 posterior <- posterior_unnorm / sum(posterior_unnorm)
37 
38 # -----
39 # STEP 3: Plotting the results
40 # -----
41 
42 # Plot prior
43 plot(model, prior, type = "h", col = "gray30", lwd = 2, ylim = c(0, max(c(prior, likelihood, posterior))),
44   main = "Bayesian Inference: Varicocele Case",
45   ylab = "Density / Likelihood",
46   xlab = "Probability of Success")
47 # Add likelihood
48 lines(model, likelihood, type = "h", col = "red", lwd = 2)
49 # Add posterior
50 lines(model, posterior, type = "h", col = "blue", lwd = 2)
51 legend("topright",
52   legend = c("Prior", "Likelihood", "Posterior"),
53   col = c("gray30", "red", "blue"),
54   lwd = 2)
55 
56 # -----
```

PROBABILITY PRACTICE PROBLEMS

```
46 # STEP 4: Comparing Hypotheses with Bayes Factor
47 # -----
48 # Hypotheses from the slide: H0 = 0.25, H1 = 0.35
49 H0 <- 0.25
50 H1 <- 0.35
51 P_data_H0 <- dbinom(18, size = 53, prob = H0)
52 P_data_H1 <- dbinom(18, size = 53, prob = H1)
53 # Bayes Factor BF[H1:H0] (Evidence in favor of H1 over H0)
54 BF_H1_H0 <- P_data_H1 / P_data_H0
55 BF_H1_H0 # ~0.35
56 # Reciprocal if needed (BF[H0:H1])
57 BF_H0_H1 <- 1 / BF_H1_H0
58 BF_H0_H1 # ~2.86
```

Critical thinking & future directions

A general rule for the question of “What would you suggest doing next?” is to think like a biologist. Think about the biological principles, not only statistics.

- Use **stratified sampling** to ensure that the sample includes **diverse** groups to account for potential **confounding variables** (e.g., age, gender, genetic background).
- Include different population subsets: for human-based studies, ensure that you have subgroups that represent the diversity of the population, such as by age, gender, or health condition, to improve the generalizability of the findings.
- Refine the control group to ensure it is as **similar** as possible to the experimental group, except for the tested key variable.
- Consider adding **more control groups** (e.g., positive control, negative control) to help validate the experiment.
- Use **randomization**. Randomly assign subjects to the experimental or control group to reduce selection bias.
- Use **double blinding**. Reduce bias in data collection and analysis, where the experimenters are unaware of the group assignments.
- **Longitudinal study design:** the study is cross-sectional, consider a longitudinal design to track changes over time in the same subjects. This can provide insights into the dynamics of the effect and its persistence.
- **Dose-response relationships:** Include multiple doses/concentrations of a treatment to clarify how the magnitude of the effect varies with doses/concentrations. This can provide insights into the pharmaco-dynamics or effects of a compound.
- The experiment is unable to address the long-term effects. So **follow-up experiments** needs to be designed to address these gaps.
- The experiment is unable to offer specific **mechanistic insights**, such as what is the underlying molecular signalling pathway responsible for it?
 - Consider using techniques such as Western Blotting, immunohistochemistry to explore underpinning pathways.
 - Other Options:
 - * Pathway-specific inhibitors
 - * RNA interference (RNAi)
 - * CRISPR/Cas9 gene knockouts: manipulate genes or proteins.
 - * Genomics, transcriptomics, or proteomics: identify changes at a molecular level (scRNA-seq find DEGs).
- Consider using a **different model**. Testing in other organisms or systems (e.g., cell culture vs. animal models) to broaden the findings.
- **Use other tests**
- The ethical issues when it comes to **choosing power**:
 - Too small a sample size with an underpowered study:
 - * Waste resources
 - * Can't reject H_0
 - * Misleading conclusions if results are nonsignificant
 - * Unethical if the conclusions lead to inferior treatment clinically
 - Too big a sample size with an overpowered study:
 - * Waste resources; especially for needless sacrifice of animals
 - * Pick up essentially trivial results which are meaningless
 - * Cost of collecting data > potential benefits
 - Stopping rules: see the slides
- The use of **bootstrapping**
 - The data we have is incomplete
 - Fundamentally not an exact method
 - (“*But it is all we have!*”)
- Power analysis of ANOVA
 - `library(pwr)`
 - `pwr.anova.test(k=GroupNumber, n=SampleNumberPerGroup, f = Cohen'sf, sig.level=alpha)` a power of 0.8 is sufficient.
 - `pwr.anova.test(k=GroupNumber, n=SampleNumberPerGroup, f = Cohen'sf, sig.level=alpha, power=0.8)` to see how many samples per group are needed to achieve a power of 0.8.
- If there are animal died/participants dropped out, it is reasonable to run toxicological tests to find out the most affected organs and check what causes this toxicity.

Concept Table

Distribution terms

Concept	Definition/Purpose	Equation
Uniform distribution	All outcomes within a specified range are equally likely, e.g., rolling a dice	x
Bernoulli distribution	A discrete distribution representing <i>a single trial</i> with two possible outcomes (success or failure), with a probability p of success	x
Binomial distribution	A discrete distribution representing <i>the number of successes</i> in n independent trials, each with a probability p of success (that is, a series of Bernoulli distribution)	x
Normal distribution	A continuous, symmetric distribution characterized by its bell-shaped curve, described by its mean μ and standard deviation σ	$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$ <ul style="list-style-type: none"> Mean: μ Variance: σ^2 Characteristics: the 68-95-99.7 rule x
Sampling distribution	Take samples of size n from a population which follows a particular distribution pattern and record the mean for multiple times. The distribution of the sample mean is sampling distribution <ul style="list-style-type: none"> the more sampling time increases the more sampling distribution resembles the sample the more sample size increases the more towards-the-center the sampling distribution will get 	x

Sampling terms

Concept	Definition/Purpose
Single random sampling	Every member of the population has an equal chance of being selected
Stratified sampling	The population is divided into strata based on certain characteristics, and random samples are taken from each strata according to strata size <ul style="list-style-type: none"> The population is divided into clusters groups Some clusters are <i>randomly</i> selected Individuals within those selected clusters are randomly sampled A cost-effective method for large populations spread over wide geographical areas
Two-stage cluster sampling	Select the most extreme or deviant cases (very high or very low) in a population <ul style="list-style-type: none"> Not representative Provide an upper/lower bound The difference between sample estimate and actual population Due to randomness, cannot be eliminated All samples are wrong, some are more useful than others so we can: <ul style="list-style-type: none"> Collect more samples Collect better samples
Extreme sampling	Select the most extreme or deviant cases (very high or very low) in a population <ul style="list-style-type: none"> Not representative Provide an upper/lower bound The difference between sample estimate and actual population Due to randomness, cannot be eliminated All samples are wrong, some are more useful than others so we can: <ul style="list-style-type: none"> Collect more samples Collect better samples
Sampling error	• Due to non-random selection of items, whether intentionally or not <ul style="list-style-type: none"> Due to sampling method, can be eliminated
Sampling bias	Even if a population is not normally distributed, the sampling distribution for large enough sample sizes n will tend to be normal.
Central limit theorem	Summary of n: the more n increases... <ul style="list-style-type: none"> the more towards-the-center will the sampling distribution get the more normal will the sampling distribution be the less SEM will be

General statistical terms

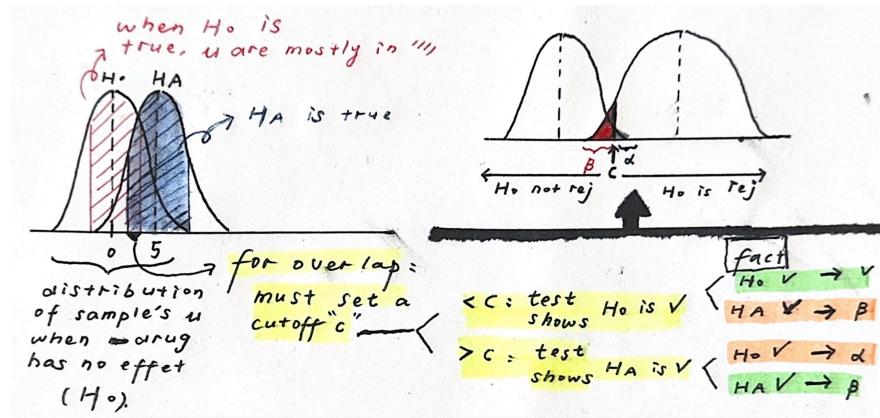
CONCEPT TABLE

Concept	Definition/Purpose	Equation
Standard score (Z-score)	Z-score provides a standard measure of how many standard deviations a data point is from the mean (distance from the mean), this concept is the basis of statistic such as F-statistic. <ul style="list-style-type: none"> • $Z = 0$ indicates that the data point is exactly at the mean • Positive Z indicates values above the mean • Reported as “the data point is 2 standard deviations above the mean” 	$Z = \frac{X-\mu}{\sigma}$ <ul style="list-style-type: none"> • The Z-score • X: The raw score (data point) • μ: The mean of the population • σ: The standard deviation of the population
Standard error of the mean (SEM)	A measure of how well the sample mean estimates the population mean <ul style="list-style-type: none"> • the larger SEM is, the better the sample is, the more appropriate the sampling method is 	$SEM = \frac{\sigma}{\sqrt{n}}$ <ul style="list-style-type: none"> • σ represents the standard deviation of the population • n is the sample size
Compare 2 Samples non-statistically	Depends on <ol style="list-style-type: none"> 1. Effect Size 2. Sample Size 	$t = \frac{\bar{x}-\mu}{\frac{s}{\sqrt{n}}}$
Effect size (Cohen's d)	This is the base of the rationale of t-statistic A quantitative measure of the magnitude of the difference between groups. It helps to uncover the practical meaning of findings, beyond just statistical significance, but magnitude of the effect	Cohen's d: Used to measure the standardized mean difference between two groups: $d = \frac{\bar{X}_1-\bar{X}_2}{s}$ <ul style="list-style-type: none"> • \bar{X}_1 and \bar{X}_2 are the means of the two groups. • s is the pooled standard deviation. <ul style="list-style-type: none"> – Small effect: $d = 0.2$ – Medium effect: $d = 0.5$ – Large effect: $d = 0.8$
Confidence Interval	- Confidence level: 95% - Explanation: I am 95% sure that the real value lies in the range $(\underline{ }, \overline{ })$	95% CI = estimate $\pm t_{0.975, df} \cdot SE_{\text{of Estimate}}$
Statistical inference	Infer properties of the population from sampled data using statistical tests Infer parameter (non-hat) from the manipulation of estimation (hat)	estimation (means sampling and get statistics here) + hypothesis testing
Variability	<ul style="list-style-type: none"> • Without variability , there is no need for statistics • The lower the variability, the higher probability that we will reject the null hypothesis 	<p>Same mean values, same sample sizes, different variabilities</p> <p>medium variability The mean difference is the same for all three cases high variability</p> <p>low variability</p>

Hypothesis testing terms

	H_0 is true	H_a is true
Test shows significance	<u>V (not wanted)</u> $= \alpha(\text{SignificanceLevel})$ TYPE 1 ERROR FALSE POSITIVE	<u>S (wanted)</u> $= Sensitivity = Power = 1 - \beta$ CORRECT DECISION TRUE POSITIVE
Test shows non-significant	<u>U (wanted)</u> $= Specificity$ CORRECT DECISION TRUE NEGATIVE	<u>T (not wanted)</u> $= \beta$ TYPE 2 ERROR FALSE NEGATIVE

Terms	Name	Formula
FDR	False Discovery Rate	$FDR = \frac{V}{V+S}$
TPR (Sensitivity)	True Positive Rate	$TPR = \frac{S}{S+T}$
FPR (Type 1 error, α)	False Positive Rate	$\alpha = \text{Type 1 error} = FPR = \frac{V}{V+U}$
Type 2 error, β	False Negative Rate	$\beta = \text{Type 2 error} = \frac{T}{S+T}$



As long as the distribution of H_0 and H_A overlapped, the cutoff is a trade-off between α and β . You have to sacrifice one of those to compensate the other.

For how to balance these two risks, see ABOVE.

Correlation terms

Concept	Definition/Purpose	Equation
Variance	<ul style="list-style-type: none"> To measure how much each value in the data set deviates from the mean Does NOT equal to variability 	$s_x^2 = \frac{\sum(X-\bar{X})^2}{n-1}$
Covariance	<ul style="list-style-type: none"> To find out whether two variables move in the same direction (positive covariance) or in opposite directions (negative covariance) 	$\text{Cov}(x, y) = \frac{\sum(X-\bar{X})(Y-\bar{Y})}{n-1}$
Correlation coefficient	<ul style="list-style-type: none"> Correlation is the standardized covariance Range: $[-1, 1]$ 	$r = \frac{\text{Cov}(x, y)}{s_x s_y}$
Linear regression	<ul style="list-style-type: none"> Y: dependent variable/response variable, X: independent variable/regressor/predictor variable, β_0: Y-intercept, β_1: curve slope, ϵ: error term (unexplained variance) 	$Y = \beta_0 + \beta_1 X + \epsilon$
Residuals	<ul style="list-style-type: none"> Random error that failed to be explained by the model, unwanted 	$\text{Residuals} = \text{ActualValue} - \text{PredictedValue}$