

A large blue circular graphic on the left side of the slide, partially cut off by the edge.

# Statistics for Data Science

## Optional Content - Hypothesis Testing

This file is meant for personal use by [hhung.inbox@gmail.com](mailto:hhung.inbox@gmail.com) only.  
Sharing or publishing the contents in part or full is liable for legal action.

# Agenda - Some Important Tests

1. Test for one mean
2. Test for equality of means
3. Test for equality of means - equal std dev
4. Test for equality of means - unequal std dev
5. Paired test for equality of means
6. Test for one proportion
7. Test for two proportions
8. Test for one variance
9. Test for equality of variances
10. Test of independence
11. ANOVA test



# Some important Statistical Tests

This file is meant for personal use by [hhung.inbox@gmail.com](mailto:hhung.inbox@gmail.com) only.  
Sharing or publishing the contents in part or full is liable for legal action.

# Hypothesis Testing Frameworks

Choice of test depends on test statistic and data availability

## Means

Compare the sample mean to the population mean when std dev is known

1-sample z-test

Compare the sample mean to the population mean when std dev is unknown

1-sample t-test

Compare the sample means from 2 independent populations when std devs are known

2-sample ind. z-test

Compare the sample means from 2 independent populations when std devs are unknown

2-sample ind. t-test

Compare the sample means from 2 related populations when std devs are unknown

Paired t-test

Compare the sample means from 2 or more independent populations

ANOVA Test

## Proportions

Compare the sample proportion to the population proportion

1-sample z-test

Compare the sample proportions from two populations

2-sample z-test

## Variances

Compare the sample variance to the population variance

Chi-Square test

Compare the sample variances from two populations

F-test

## Frequencies

Check whether the categorical variables from a population are independent

Chi-Square Test of Independence

# Test for one mean

## Example

A certain food aggregator ZYX is facing stiff competition from its main rival SWG during Corona period. To retain business, ZYX is advertising that, within a radius of 5 km from the restaurant where the order is placed, it can still deliver in 40 minutes or less on the average (and changed condition has not made any impact on them).

The delivery times in minutes of 25 randomly selected deliveries are given in a CSV file.

Assuming the delivery distribution is approximately normal, is there enough evidence that ZYX's claim is false?

**This is clearly a one-tailed hypothesis problem, concerning population mean  $\mu$ , the average delivery time.**

This file is meant for personal use by hhung.inbox@gmail.com only.  
Sharing or publishing the contents in part or full is liable for legal action.

# Test for One Mean - Unknown Std Dev

Significance of the test	Assumptions	Test Statistic Distribution
Test for population mean $H_0 : \mu = \mu_0$	<ul style="list-style-type: none"><li>• Continuous data</li><li>• Normally distributed population and sample size <math>&lt; 30</math></li><li>• Unknown population standard deviation</li><li>• Random sampling from the population</li></ul>	t distribution  (The test is also known as <b>One-sample t-test</b> )

# Test for equality of means (Known std dev)



## Example

To compare customer satisfaction levels of two competing media channels, 150 customers of Channel 1 and 300 customers of Channel 2 were randomly selected and were asked to rate their channels on a scale of 1-5, with 1 being least satisfied and 5 most satisfied. (The survey results are summarized in a CSV file)

Test at 0.05 level of significance whether the data provide sufficient evidence to conclude that channel 1 has a higher mean satisfaction rating than channel 2.

**This is a two-sample problem where the channel 1 and channel 2 populations are independent. Further, this is a one-tailed hypothesis problem, concerning population means  $\mu_1$  and  $\mu_2$ , the mean customer satisfaction for channel 1 and channel 2 respectively.**

This file is meant for personal use by hiring.inbox@gmail.com only.  
Sharing or publishing the contents in part or full is liable for legal action.

# Test for Equality of Means - Known Std Devs

Significance of the test	Assumptions	Test Statistic Distribution
Test for equality of two population means $H_0 : \mu_1 = \mu_2$	<ul style="list-style-type: none"><li>• Continuous data</li><li>• Normally distributed population or sample size <math>&gt; 30</math></li><li>• Independent populations</li><li>• Known population standard deviations <math>\sigma_1</math> and <math>\sigma_2</math></li><li>• Random sampling from the population</li></ul>	Standard Normal distribution  (The test is also known as <b>Two independent sample z-test</b> )

# Test for equality of means (Equal and unknown std dev)

# Example

In the lockdown period, because of working from home and increased screen time, many opted for listening to FM Radio for entertainment rather than watching Cable TV. An advertisement agency randomly collected daily usage time data (in minutes) from both type of users and stored it in a CSV file.

Assuming daily Radio and TV usage time are normally distributed, do we have enough evidence to conclude that there is any difference between daily TV and Radio usage time at 0.05 significance level?

This is a two-sample problem where FM Radio and Cable TV users are assumed independent. Further, this is a two-tailed hypothesis problem, concerning population means  $\mu_1$  and  $\mu_2$ , the daily mean usage time of Radio and TV respectively.

This file is meant for personal use by hhung.inbox@gmail.com only.  
Sharing or publishing the contents in part or full is liable for legal action.

# Test for Equality of Means : Equal Std Devs

Significance of the test	Assumptions	Test Statistic Distribution
Test for equality of two population means $H_0 : \mu_1 = \mu_2$	<ul style="list-style-type: none"><li>• Continuous data</li><li>• Normally distributed populations</li><li>• Independent populations</li><li>• Equal population standard deviations</li><li>• Random sampling from the population</li></ul>	t distribution  (The test is also known as <b>Two independent sample t-test</b> )

# Test for equality of means (Unequal and unknown std dev)

# Example

SAT verbal scores of two groups of students are given in a CSV file. The first group, **College**, contains scores of students whose parents have at least a bachelor's degree and the second group, **High School**, contains scores of students whose parents do not have any college degree.

The Education Department is interested to know whether the sample data support the theory that students show a higher population mean verbal score on SAT if their parents attain a higher level of education.

Assuming SAT verbal scores for two populations are normally distributed, do we have enough statistical evidence for this at 5% significance level?

This is a two-sample problem as the College and High School populations are different. Further, this is a one-tailed hypothesis problem, concerning population means  $\mu_1$  and  $\mu_2$ , the mean verbal score on SAT for College and High School groups.

This file is meant for personal use by hhung.inbox@gmail.com only.  
Sharing or publishing the contents in part or full is liable for legal action.

# Test for Equality of Means : Unequal Std Devs

Significance of the test	Assumptions	Test Statistic Distribution
Test for equality of two population means $H_0 : \mu_1 = \mu_2$	<ul style="list-style-type: none"><li>• Continuous data</li><li>• Normally distributed populations</li><li>• Independent populations</li><li>• Unequal population standard deviations</li><li>• Random sampling from the population</li></ul>	t distribution  (The test is also known as <b>Two independent sample t-test</b> )



# Paired Test for equality of means

## Example

Typical prices of single-family homes in Florida are given for a sample of 15 metropolitan areas (in 1000 USD) for 2002 and 2003 in a CSV file.

Assuming the house prices are normally distributed, do we have enough statistical evidence to say that there is an increase in the house price in one year at 0.05 significance level?

This is a paired sample problem as the two observations (for 2002 and 2003) are taken on one sampled unit (a metropolitan area). Further, this is a one-tailed hypothesis problem, concerning population means  $\mu_1$  and  $\mu_2$ , the mean house price in 2002 and 2003 respectively.

# Paired test for Equality of Means

Significance of the test	Assumptions	Test Statistic Distribution
Test for equality of two population means $H_0 : \mu_1 = \mu_2$	<ul style="list-style-type: none"><li>• Continuous data</li><li>• Normally distributed populations</li><li>• Independent observations</li><li>• Random sampling from the population</li></ul>	t distribution  (The test is also known as <b>Paired t-test</b> )

# Test for One Proportion

## Example

A researcher claims that Democratic party will win in the next United States Presidential election.

To test her belief the researcher randomly surveyed 90 people and 24 out of them said that they voted for Democratic party.

Is there enough evidence at  $\alpha = 0.05$  to support this claim?

This is clearly a one-tailed test, concerning population proportion  $p$ , the proportion of people voted from Democratic party.

# Test for One Proportion

Significance of the test	Assumptions	Test Statistic Distribution
Test for population proportion $H_0 : p = p_0$	<ul style="list-style-type: none"><li>• Binomially distributed population</li><li>• Random sampling from the population</li><li>• When both mean (<math>np</math>) and <math>n(1-p)</math> are greater than or equal to 10, the binomial distribution can be approximated by a normal distribution</li></ul>	Standard Normal distribution  (The test is also known as <b>One proportion z-test</b> )

# Test for Two Proportions

## Example

A car manufacturer aims to improve its products' quality by reducing the defects. So, the manufacturer randomly checks the efficiency of two assembly lines in the shop floor. In line 1, there are 20 defects out of 200 samples and In line 2, there are 25 defects out of 400 samples.

At 5% level of significance, do we have enough statistical evidence to conclude that the two assembly procedures are different?

This is clearly a two-tailed test, concerning two population proportion  $p_1$  and  $p_2$ , the proportion of defects in assembly line 1 and assembly line 2 respectively.



# Test for Two Proportions

Significance of the test	Assumptions	Test Statistic Distribution
Test for equality of two population proportions $H_0 : p_1 = p_2$	<ul style="list-style-type: none"><li>• Binomially distributed populations</li><li>• Independent populations</li><li>• Random sampling from the populations</li><li>• When both mean (<math>np</math>) and <math>n(1-p)</math> are greater than or equal to 10, the binomial distribution can be approximated by a normal distribution</li></ul>	Standard Normal distribution  (The test is also known as <b>Two proportions z-test</b> )

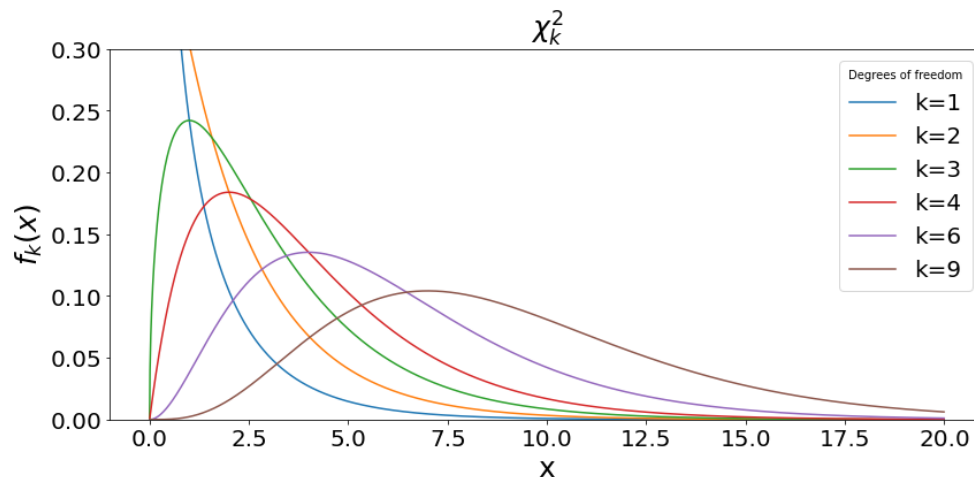
# Test for One Variance

# Test for Variance

Variance tests are used for a comparison of variability, often as a predecessor for other tests

Let us take many samples of the same size from a normal population and find the sample variances

They follow a **chi-square ( $\chi^2$ ) distribution**, which is dependent on the degrees of freedom



## Example

It is conjectured that the standard deviation for the annual return of mid cap mutual funds is 22.4%, when all such funds are considered and over a long period of time. The sample standard deviation of a certain mid cap mutual fund based on a random sample of size 32 is observed to be 26.4%.

Do we have enough evidence to claim that the standard deviation of the chosen mutual fund is greater than the conjectured standard deviation for mid cap mutual funds at 0.05 level of significance?

**This is clearly a one-tailed test, concerning population variance, the variance for mid cap mutual funds.**

# Test for One Variance

Significance of the test	Assumptions	Test Statistic Distribution
Test for population variance $H_0 : \sigma^2 = \sigma_0^2$	<ul style="list-style-type: none"><li>• Continuous data</li><li>• Normally distributed population</li><li>• Random sampling from the population</li></ul>	Chi Square distribution  (The test is also known as Chi-square test for variance)

# Test for Equality of Variances

## Example

The variance of a process is an important quality of the process. A large variance implies that the process needs better control and there is opportunity to improve.

The data (Bags.csv) includes weights for two different sets of bags manufactured from two different machines. It is assumed that the weights for two sets of bags follow normal distribution.

Do we have enough statistical evidence at 5% significance level to conclude that there is a significant difference between the variances of the bag weights for the two machines.

**This is clearly a two-tailed test, concerning two population variances, the variance for bag 1 weights and the variance for bag 2 weights.**

This file is meant for personal use by nhung.inbox@gmail.com only.  
Sharing or publishing the contents in part or full is liable for legal action.

# Test for Equality of Variances

Significance of the test	Assumptions	Test Statistic Distribution
Test for equality of two population variances $H_0: \sigma_1^2 = \sigma_2^2$	<ul style="list-style-type: none"><li>• Normally distributed populations</li><li>• Independent populations</li><li>• Larger variance should be placed in the numerator</li></ul>	F distribution  (The test is also known as <b>F-test for variances</b> )





# Test of Independence

This file is meant for personal use by [hhung.inbox@gmail.com](mailto:hhung.inbox@gmail.com) only.  
Sharing or publishing the contents in part or full is liable for legal action.

# Chi Square Test for Independence

2x2 contingency table that describes two variables (smoking and gender) at two levels each and stores the number of observations at each cell

	Male	Female	Total
Smoker	120	100	220
Non-smoker	60	140	200
Total	180	240	420

We are interested to know whether the **two variables are independent**

$H_0$ : Smoking and gender are independent.

$H_a$ : Smoking and gender are not independent.

## Example

The following table summarizes beverage preference across different age-groups.

	Beverage Preference		
Age	Tea/Coffee	Soft Drink	Others
21 - 34	25	90	20
35 - 55	40	35	25
> 55	24	15	30

Does beverage preference depend on age?

This is a problem of Chi-Square test of independence, concerning the two independent categorical variables, Age and Beverage Preference.

# Chi-Square Test for Independence

Significance of the test	Assumptions	Test Statistic Distribution
In a contingency table $H_0$ : The row and column variables are independent	<ul style="list-style-type: none"><li>• Categorical variables</li><li>• Expected value of the number of sample observations in each level of the variable is at least 5</li><li>• Random sampling from the population</li></ul>	Chi Square distribution  (The test is also known as <b>Chi-square test of independence</b> )

# Analysis of Variance (ANOVA)

# ANOVA Test : Introduction

**Analysis of Variance (ANOVA)** is used to determine whether the means of more than two independent populations are significantly different.

$$H_0: \mu_1 = \mu_2 = \mu_3 = \dots = \mu_k$$

$H_a$ : at least one of these means is not the same



Why do we call it ANOVA? - The mathematical tools to calculate the p-value rely heavily on using the variances of the populations.



ANOVA is used in various problems such as comparing the yields of the crop from several varieties of seeds, comparing the gasoline mileage of various types of automobiles, etc.

# ANOVA Test : Some important terms

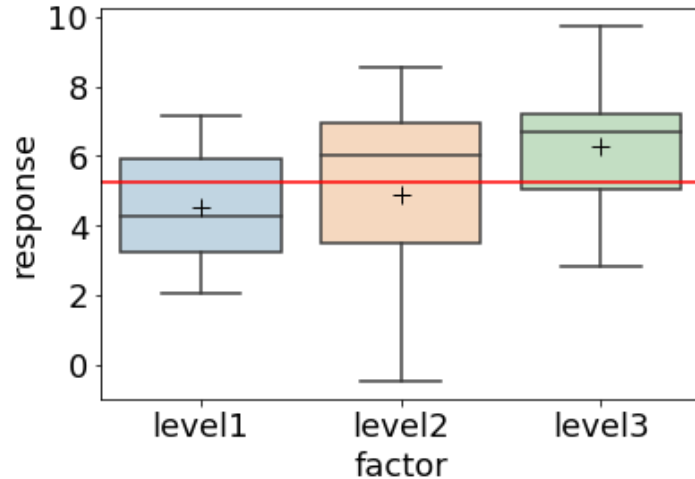
**Response:** Dependent variable which is continuous and assumed to follow a normal distribution

Consider, an example where interest lies in comparing the **weekly volume of sales** by different teams of sales executives.

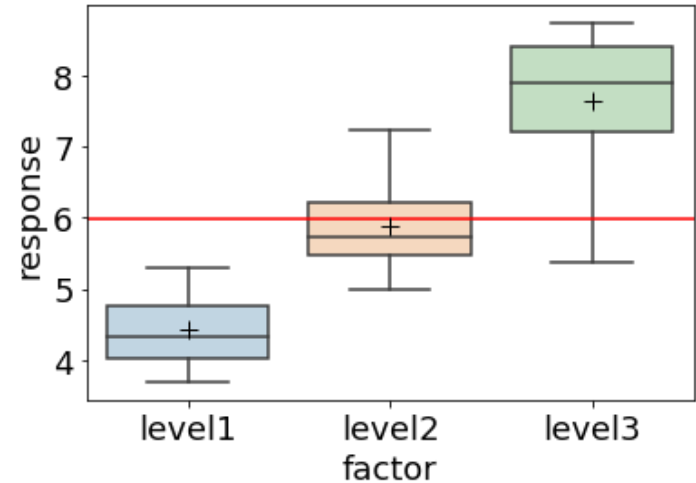
**Factor:** Independent explanatory variable with several levels

# ANOVA Test : One-way ANOVA

One-way ANOVA is used when the response variable depends on a single factor.



Between group variation is lower



Between group variation is higher



# ANOVA Test : How it works

F-Statistic is the ratio of the between group variations to within group variations.



$$F - statistic = \frac{\textit{Between group variations}}{\textit{Within group variations}}$$



A large value of F-Statistic indicates that there is more variation between groups than within groups.



Thus, it will provide evidence against the null hypothesis.

## Example

Traffic management inspector in a certain city wants to understand whether carbon emissions from different cars are different. The inspector has reasons to believe that Fuel type may be one important factor responsible for differences in carbon emission.

For this purpose, the inspector has taken random samples from all registered cars on the road in that city and would like to test if the amount of carbon emission release depends on fuel type at 5% significance level.

**Here, we will compare the means of emission for the three different fuel types.**

# ANOVA Test : One-way ANOVA

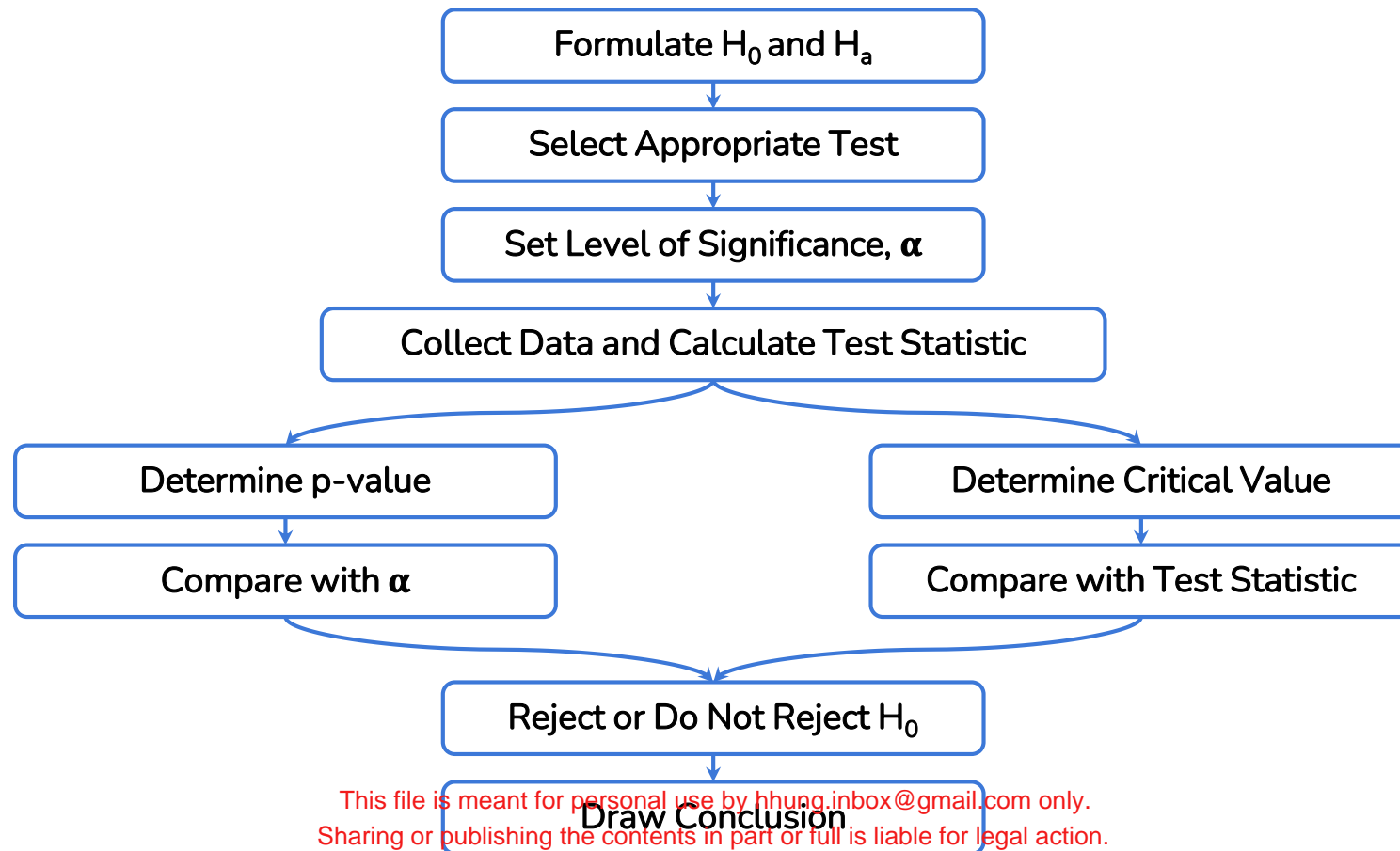
Significance of the test	Assumptions	Test Statistic Distribution
Test for means for more than two populations $H_0$ : All population means are equal	<ul style="list-style-type: none"><li>• The populations are normally distributed</li><li>• Samples are independent simple random samples</li><li>• Population variances are equal</li></ul>	F distribution  (The test is also known as <b>One-way ANOVA F-test</b> )



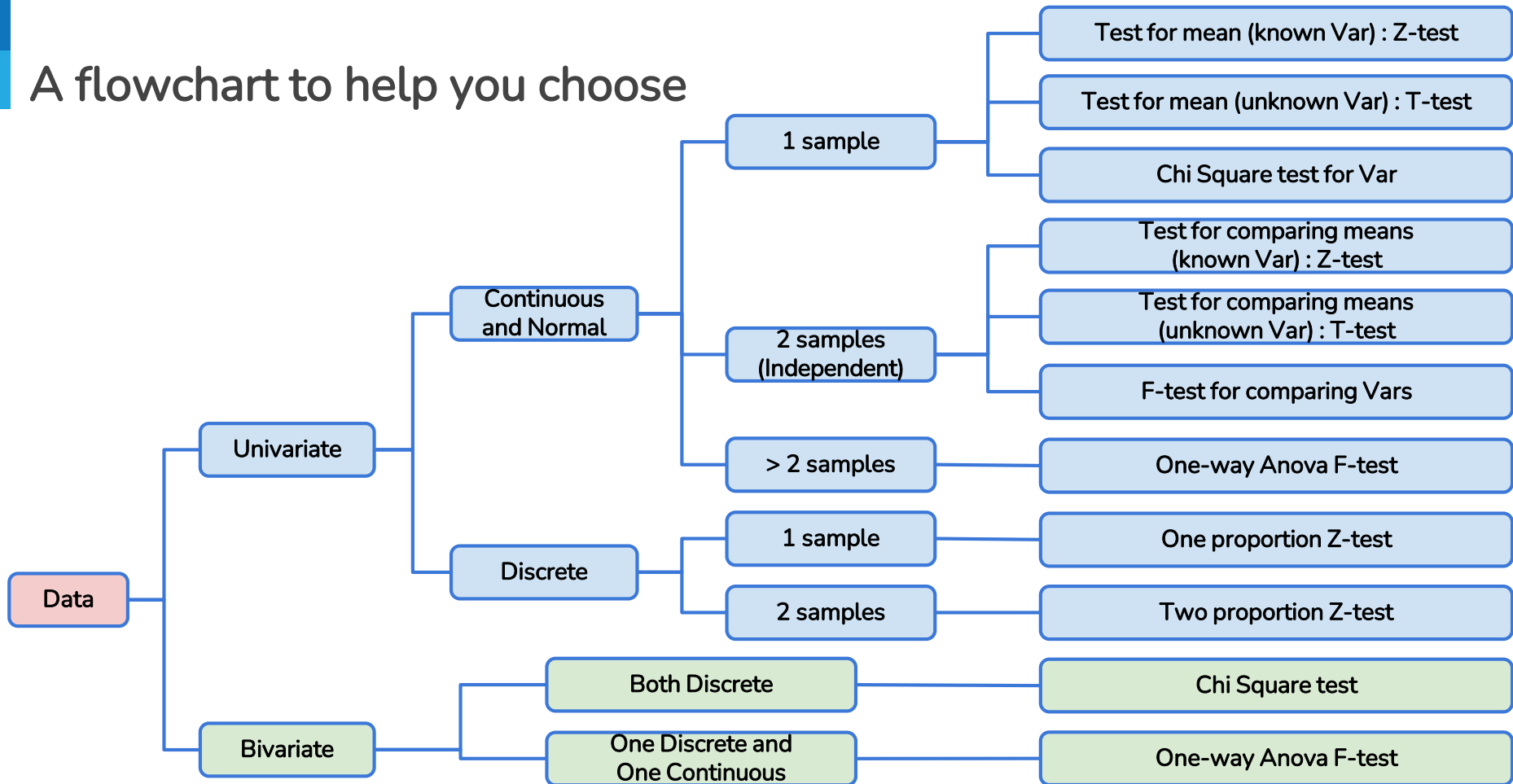
# Now let's summarize

This file is meant for personal use by [hhung.inbox@gmail.com](mailto:hhung.inbox@gmail.com) only.  
Sharing or publishing the contents in part or full is liable for legal action.

# Hypothesis Testing Steps



# A flowchart to help you choose



This file is meant for personal use by [hhung.inbox@gmail.com](mailto:hhung.inbox@gmail.com) only.  
Sharing or publishing the contents in part or full is liable for legal action.