

# Inferential Statistics

## L5 - Hypothesis Testing

Erlis Ruli (erlis.ruli@unipd.it)

Department of Statistics, University of Padova

# Contents

- 1 Motivation
- 2 Mathematical formulation
- 3 Methods for computing tests
  - The likelihood ratio test
  - The Wald test
  - Pearson's  $\chi^2$  test
- 4 The p-value
- 5 Methods for evaluating tests
- 6 Some notable tests

# Problem statement

Suppose that the average energy consumption of our population of WMs mounting a standard motor is  $\mu_0$ .

It's claimed that NG1 family motors would lead to more efficient WMs, i.e. would lead to average consumption  $\mu$ , with  $\mu < \mu_0$ .

There are two possibilities:

- the claim is false, so  $\mu \geq \mu_0$ ; this is called Null Hypothesis (“null” because it adds nothing to the current state of art)
- the claim is true, so  $\mu < \mu_0$ ; it's called Alternative Hypothesis.

## Problem statement (cont'd)

Concretely, suppose that  $\mu_0 = 20$ .

We equip 10 WM's with the NG1 motor and measure their E consumption.

Let these energy values be

19.1, 20.6, 17.3, 21.1, 19.5, 19.5, 21.4, 19.1, 20.5, 19.5.

Their average is 19.76. But  $19.76 < \mu_0 = 20$ , so the NG1 motor seems to lead on average to more efficient WMs!

Is that really so? Couldn't this be due to a pure luck?

# A Simple formulation

Suppose the sample above is a realisation of the iid random sample  $Y_1, \dots, Y_n$  with  $Y_i \sim N(\mu, 5)$ .

This is a reasonable assumption given that overall energy consumption of a WM is the sum of the consumptions due to the various components of a WM (motor, resistor, etc.).

For this specific example we have (from L4) that

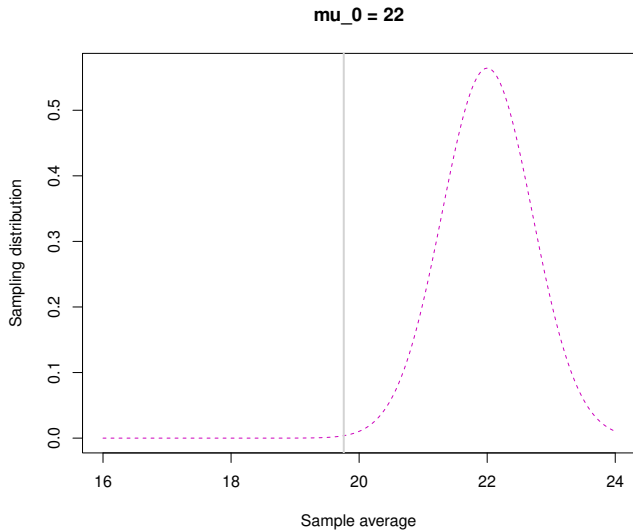
$$\frac{\sqrt{n}(\bar{Y} - \mu)}{\sqrt{5}} \sim N(0, 1),$$

and thus

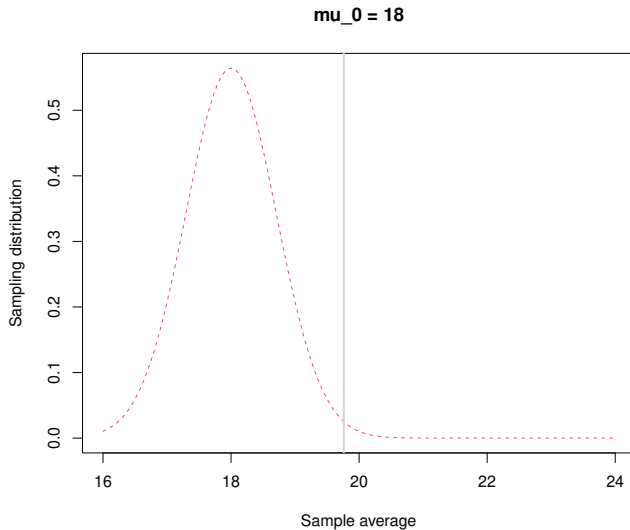
$$\bar{Y} \sim N(\mu, 0.5).$$

The following figure shows this distribution for several values of  $\mu$ .

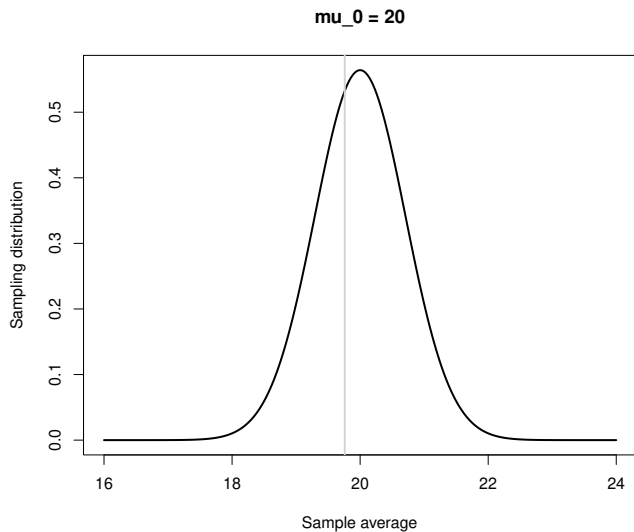
# Sampling distribution of $\bar{Y}$



# Sampling distribution of $\bar{Y}$

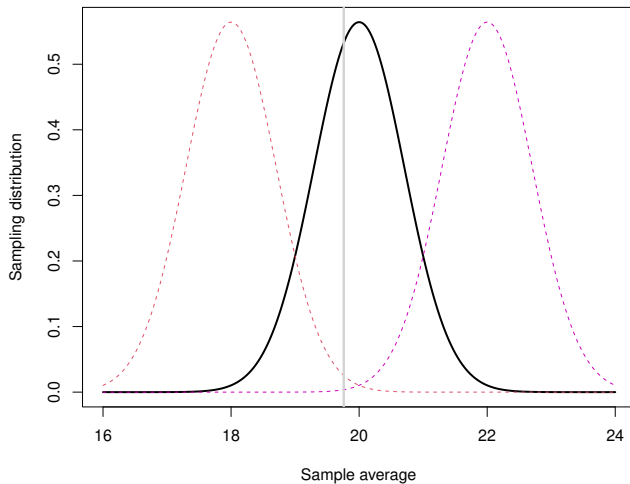


# Sampling distribution of $\bar{Y}$





# Sampling distributions of $\bar{Y}$



# Test hypothesis and decision

Suppose that we judge surprising all values that under the sampling distribution are very unlikely to happen, say all those values with prob less than 0.01.

Under the sampling distribution with  $\mu_0 = 20$ , this value is 19.42, since  $P_{\mu_0}(\bar{Y} \leq 19.42) = 0.01$ .

We got a decision rule: a sample average  $< 19.42$  is deemed surprising and should make us suspect about the worthiness of the null hypothesis, so we reject the null hypothesis. Otherwise we do not reject.

In the case above, the observed sample average was  $\bar{y} = 19.76 > 19.42$ , according to the rule, we should not be surprised and thus do not reject the null hypothesis.

# Terminology

In the problem above we tested the null hypothesis  $H_0 : \mu \geq \mu_0$  against the alternative  $H_1 : \mu < \mu_0$ .

In other situations we may be interested in testing

$$H_0 : \mu \leq \mu_0 \text{ against } H_1 : \mu > \mu_0 \quad (*)$$

or

$$H_0 : \mu = \mu_0 \text{ against } H_1 : \mu \neq \mu_0 \quad (**)$$

Hypotheses s.t. (\*) are called one-tailed, and (\*\*) are called two-tailed.

# That was too simple

In our motivating example we assumed population variance was known ( $\sigma^2 = 5$ ). In practice, this assumption is not realistic and must be relaxed.

Furthermore, we assumed an iid random sample with  $Y_i$  following a normal distribution. But, in many problems, the normal distribution is not suitable.

The point is that, relaxing  $\sigma^2 = 5$  or the distributional assumptions makes the above test useless; thus we have to look elsewhere. In other words, we need to know how to build a test for a given problem at hand.

In the rest of this lecture we will discuss two popular methods (both frequentist-parametric) and we will provide criteria for evaluating their performance.

# Setting the scene

Let  $Y_1, \dots, Y_n$  be an iid random sample with  $Y_i \sim F_\theta$ , where  $\theta \in \Theta$  is the unknown parameter with parameter space  $\Theta$ .

We assume that  $F_\theta$  has pdf  $f$ , indexed by the same parameter  $\theta$ .

The iid assumption can be relaxed, but let's keep it simple for the moment.

We denote by  $\mathcal{Y}$  the range of  $Y_i$  and by  $\mathcal{Y}^n = \mathcal{Y} \times \mathcal{Y} \times \dots \times \mathcal{Y}$ , the Cartesian product of  $\mathcal{Y}$   $n$  times.

# Setting the scene (cont'd)

Performing a statistical test essentially entails

building a decision rule from a sample to decide if reject

$$H_0 : \theta \in \Theta_0 \text{ in favour of } H_1 : \theta \in \Theta_0^c,$$

where  $\Theta_0 \subset \Theta$ .

Specifically, a test is a binary decision rule operating on a subset  $R \subset \mathcal{Y}^n$  as follows:

reject  $H_0$  if the observed sample  $\mathbf{y} = (y_1, \dots, y_n) \in R$ , and accept  $H_0$  otherwise.

# Methods for computing statistical tests

A statistical test is typically defined on the basis of a test statistic  $T(\mathbf{Y}) = T(Y_1, \dots, Y_n)$  which is a function of the sample and closely related to a statistic seen in L4.

The test statistic used for computing the statistical test determines the nature and the name of the statistical test itself.

# Likelihood Ratio Tests

The likelihood ratio test statistic for  $H_0 : \theta \in \Theta_0$  against  $H_1 : \theta \in \Theta_0^c$  is

$$\lambda(\mathbf{y}) = \frac{\sup_{\theta \in \Theta_0} L(\theta)}{\sup_{\theta \in \Theta} L(\theta)}.$$

A likelihood ratio test (LRT) is any test that has rejection region

$$R = \{\mathbf{y} : \lambda(\mathbf{y}) < c\},$$

where  $c$  is any number s.t.  $0 \leq c \leq 1$ .

Recalling that  $\hat{\theta}$  is the MLE of  $\theta$  and denoting by  $\hat{\theta}_0$  the constrained MLE of  $\theta$  when the parameter space is  $\Theta_0$ , then the LRT statistic is

$$\lambda(\mathbf{y}) = \frac{L(\hat{\theta}_0)}{L(\hat{\theta})}.$$



## Example 1

Let  $Y_i \sim N(\theta, 5)$ , be an iid sample of size  $n$  and consider testing  $H_0 : \theta = \theta_0$  versus  $H_1 : \theta \neq \theta_0$ . Here  $\theta_0$  is a number fixed prior to the experiment.

Under  $H_0$  we have only one possible value for  $\theta$ , thus the numerator of  $\lambda(\mathbf{y})$  is  $L(\theta_0)$ . On the other hand the (unrestricted) MLE for  $\theta$  is  $\hat{\theta} = \bar{y}$ , the LRT statistic is

$$\begin{aligned}\lambda(\mathbf{y}) &= \frac{(10\pi)^{-n/2} \exp[-\sum_i (y_i - \theta_0)^2 / 10]}{(10\pi)^{-n/2} \exp[-\sum_i (y_i - \bar{y})^2 / 10]} \\ &= \exp[-n(\bar{y} - \theta_0)^2 / 10].\end{aligned}$$

The LRT is thus a test that rejects  $H_0$  for small values of  $\lambda(\mathbf{y})$ , and the rejection region can be written as

$$\left\{ \mathbf{y} : |\bar{y} - \theta_0| \geq \sqrt{-10(\log c)/n}, \right\}$$

for some  $c \in (0, 1]$

## Example 2

Let all be as in the previous example but  $H_0 : \theta \geq \theta_0$  versus  $H_1 : \theta < \theta_0$ . These are the hypotheses we wanted to test in the motivating example in slide 3.

Discarding some additive constants, the log-likelihood function is

$$\begin{aligned}\ell(\theta) &= - \sum_{i=1}^n (y_i - \theta)^2 / 10 \\ &= - \sum_i (y_i - \bar{y})^2 / 10 - n(\bar{y} - \theta)^2 / 10,\end{aligned}$$

quadratic and concave. Thus, under  $H_0$ ,  $\sup \ell(\theta)$  is attained at

- (I)  $\theta_0$  if  $\theta_0 > \bar{y}$
- (II)  $\bar{y}$  if  $\theta_0 \leq \bar{y}$ .

The denominator of the LRT statistic is as before. So the likelihood ratio test statistic is

## Example 2 (cont'd)

$$\lambda(\mathbf{y}) = \begin{cases} \exp[-n(\bar{y} - \theta_0)^2/10] & \text{if } \theta_0 > \bar{y} \\ 1 & \text{otherwise} \end{cases}$$

and the rejection region for the LRT test is

$$R = \begin{cases} \{\mathbf{y} : \bar{y} < \theta_0 - \sqrt{-10(\log c)/n}\} & \text{if } \theta_0 > \bar{y} \\ \emptyset & \text{otherwise.} \end{cases}$$

Thus the test rejects if the sample average is lower than some threshold on the left of  $\theta_0$ , precisely as we saw in slides 5-9.

We'll discuss the choice of  $c$  in the next section, but if take  $c = 0.0977$ , and go back to the example in slides 5-7 where  $\hat{\theta}_0 = 20$ ,  $n = 10$ , we have that

$$\bar{y} = 19.41 \not< \theta_0 - \sqrt{-10(\log c)/n} = 18.47$$

We cannot reject  $H_0$ , e.g.

there is no evidence about the efficiency claim of NG1.

# Four possibilities

A decision may be wrong, indeed, there are four possibilities

Decision	Truth	
	$\theta \in \Theta_0$ ( $H_0$ is true)	$\theta \notin \Theta_0$ ( $H_1$ is true)
Reject $H_0$	type I error	ok
Don't reject $H_0$	ok	type II error

In a single test, we may either get a correct (ok) or a wrong decision. In the latter case, we could make a type I or type II error.

The size of type I error is defined by

$$\alpha' = \sup_{\theta \in \Theta_0} P(\text{reject } H_0 | H_0 \text{ is true}).$$

and the size of type II error is defined by

$$\beta(\theta) = 1 - P(\text{reject } H_0 | H_1 \text{ is true}) \quad \forall \theta \in \Theta_0^c.$$

# No free meal

Ideally, we'd like error-free decision, i.e.

$$\alpha' + \beta(\theta) = 0, \forall \theta,$$

but this is impossible.

Indeed, for  $\alpha' = 0$  we must never reject  $H_0$ . But then, if  $H_1$  is true, we make a type II error for sure, so  $\beta(\theta) = 1$ .

On the other hand, for  $\beta(\theta) = 0$  we have to always reject  $H_0$ . But so doing and when  $H_0$  is true, we make a type I error for sure, so  $\alpha' = 1$ .

# Choosing the threshold

The current practice is to fix  $\alpha$  at some small value (e.g. 0.01, or 0.05) and make sure that

$$\alpha' \leq \alpha,$$

without worrying about the value of  $\beta(\theta)$ .

The rationale for such a choice is that, often, making a type I error is more dangerous than making a type II error.

Fixing  $\alpha$  entails fixing the amount of type I error, which entails fixing the size of the rejection region  $R$ , or the value of  $c$  in the case of LRT.

### Example 3

Consider again Example 2 and let  $\alpha = .01$ . The rejection region is

$$R = \{\mathbf{y} : \bar{y} < \theta_0 - \sqrt{-10(\log c)/n}\}.$$

First note that

$$\bar{Y} \sim N(\theta_0, 10/n),$$

thus

$$\begin{aligned} P_{\theta_0}(\mathbf{Y} \in R) &= P_{\theta_0}(\bar{Y} < \theta_0 - \sqrt{-10(\log c)/n}) \\ &= P(Z \leq -\sqrt{-\log c}) \leq \alpha \end{aligned}$$

Solving the last inequality, gives  $z_\alpha = -2.326 = -\sqrt{-\log c}$ , so  $c = .0977$ ; here  $z_\alpha$  denotes the quantile of level  $\alpha$  of the standard normal distribution.

# The level and the size of a test

Sometimes we may want to distinguish between cases when the inequality on the size of type I error can be reached from cases in which it cannot be reached.

Indeed

- if  $\alpha' = \alpha$  the test is called a test of size  $\alpha$ .
- If  $\alpha' \leq \alpha$ , we call it a test of level  $\alpha$ .

Often, especially when  $F_\theta$  is discrete, a test of level  $\alpha$  is the best result we can get.

However, in most practical cases, we can only compute tests of size  $\alpha$  as  $n \rightarrow \infty$ ; these are called asymptotic tests.



## Example 4 (A Poisson test)

Let  $Y_1, \dots, Y_n$  be an iid random sample from the  $\text{Poi}(\theta)$ , with  $\theta$  unknown. Suppose we want to test  $H_0 : \theta \leq \theta_0$  against  $H_1 : \theta > \theta_0$ . The LRT statistic is

$$\begin{aligned}\lambda(\mathbf{y}) &= \frac{e^{-n\theta_0} \theta_0^{\sum_i y_i} \prod_i (y_i!)^{-1}}{e^{-n\bar{y}} \bar{y}^{\sum_i y_i} \prod_i (y_i!)^{-1}} \\ &= \exp \left( - \sum_i y_i - n\theta_0 \right) (\theta_0/\bar{y})^{\sum_i y_i}.\end{aligned}$$

Here, it's not possible to determine  $c$  exactly since  $P_{\theta_0}(\mathbf{Y} \in R)$  is not computable analytically. But we can build a different test as follows.

First, note that under  $H_0$ ,  $n\bar{Y} \sim \text{Poi}(n\theta)$ . Thus, if we fix a large threshold based on this distribution, we may judge “suspicious” values of  $n\bar{y}$  above this threshold.

## A Poisson test (cont'd)

We can set the threshold to the quantile of  $\text{Poi}(n\theta_0)$  of level  $\geq 1 - \alpha$  and reject  $H_0$  if  $n\bar{y}$  is greater than the threshold.

Concretely, let  $\theta_0 = 1$ ,  $\alpha = .05$  and let the observed sample be 0, 0, 3, 5, 7, hence  $n\bar{y} = 15$ .

Because  $\sum_i Y_i \sim \text{Poi}(5)$ , a threshold is 9 (We will see in the next section why choosing any value  $< 9$ , is not ok). Because  $15 > 9$ , we reject  $H_0$ .

Note that, in this particular example  $\alpha' = .031 < .05$  and thus the test is only of level .05.

This test may be improved to have size  $\alpha$  through a technique known as randomisation; but we won't see randomised tests in this course.

# Back to LRT

In many cases the LRT statistic doesn't have a known distribution, but it has a limiting distribution as  $n$  diverges.

## Theorem 5

Suppose  $\theta = (\theta_1, \dots, \theta_q, \theta_{q+1}, \dots, \theta_r)$  and let  $\Theta_0 \subset \Theta$  s.t.

$$\Theta_0 = \{\theta : \theta_{q+1} = \theta_{0,q+1}, \theta_{q+2} = \theta_{0,q+2}, \dots, \theta_r = \theta_{0,r}\}.$$

Under  $H_0 : \theta \in \Theta_0$  and suitable regularity conditions,

$$-2 \log \lambda(\mathbf{Y}) \xrightarrow{d} \chi_{r-q}^2 \quad \text{as } n \rightarrow \infty.$$

The degrees of freedom in the limiting distribution are  $r - q = \dim(\Theta) - \dim(\Theta_0)$ ;  $\dim(S)$  denotes the dimension of the space  $S$ .

## Few notes before going on

This is an asymptotic result, i.e. valid only in the limit as  $n \rightarrow \infty$ , so LRT test is guaranteed to have size  $\alpha$  in the limit.

In practice, we work with a finite  $n$ . Nevertheless, if  $n$  is "high enough",  $-2 \log \lambda$  will have a distribution close to  $\chi^2_{r-q}$ . Thus, in practice, we read " $\xrightarrow{d}$ " as " $\sim$ ".

The discrepancy between the distributions  $-2 \log \lambda$  and  $\chi^2_{r-q}$  for fixed  $n$ , depends on many factors (number of parameters to estimate, degree of dependence in the sample, etc.) Roughly speaking, the larger  $n/r$  the smaller the discrepancy.

## Example 6 (LRT in full action)

Back to Example 4 and using the Theorem, for large  $n$

$$\begin{aligned} P_{\theta_0}(\lambda(\mathbf{Y}) < c) &= P_{\theta_0}(-2 \log \lambda(\mathbf{Y}) > -2 \log(c)) \\ &\doteq P(\chi_1^2 > -2 \log(c)). \end{aligned}$$

Equating the last probability with  $\alpha$  we have

$$-2 \log c = \chi_{1,1-\alpha}^2 \implies c = \exp\left(-\frac{1}{2}\chi_{1,1-\alpha}^2\right),$$

where  $\chi_{1,1-\alpha}^2$  denotes the quantile of level  $1 - \alpha$  of the  $\chi_1^2$  distribution.  
For  $\alpha = 0.05$ ,  $\chi_{1,1-\alpha}^2 = 3.84$ , so  $c = 0.1465$ .

Because (check!)  $\lambda(\mathbf{y}) = 0.0015 < c$ , we reject  $H_0$  (again).

# The Wald test

Is useful for testing  $H_0 : \theta = \theta_0$  against  $H_1 : \theta \neq \theta_0$ , when there is an estimator for  $\hat{\theta}$  that has (at least a limiting) normal distribution.

We saw in L4, that, under regularity conditions, the MLE has a limiting normal distribution. Indeed, the Wald test is typically used in conjunction with the MLE.

Formally, for a scalar parameter

$$\hat{\theta} \sim N\left(\theta, I_n(\hat{\theta})^{-1}\right).$$

The Wald test of approximate level  $\alpha$  is then to reject  $H_0 : \theta = \theta_0$  if

$$|\text{Wald test statistic}| = |W| = \left| \frac{\hat{\theta} - \theta_0}{\widehat{\text{se}}} \right| \geq z_{1-\alpha/2},$$

where  $\widehat{\text{se}} = \sqrt{1/I_n(\hat{\theta})}$  is the estimated standard error of  $\hat{\theta}$ ; the asymptotically equivalent version with  $J$  in place of  $I$  may also be used.

# Example of Wald test, use the previous example

## Example 7

An IT-alert message test was sent to some of the residents in regione Veneto on 21st September 2023. It was claimed that the message was sent to roughly  $1/3$  of the population. Suppose we take a random sample of adults in the regione Veneto and ask them if they received this message or not.

So let  $X_1, \dots, X_n$  be an iid sample with  $X_i \sim \text{Ber}(\theta)$  ( $X_i = 1$  if message received). We wish to test  $H_0 : \theta = 1/3$  vs  $H_1 : \theta \neq 1/3$ .

The MLE of  $\theta$  is  $\hat{\theta} = \sum_i X_i / n$ . A Wald test of approximate level  $\alpha = 0.05$  has  $\hat{\text{se}} = \sqrt{\hat{\theta}(1 - \hat{\theta})/n}$ , thus

$$W = \frac{\sqrt{n}(\hat{\theta} - 1/3)}{\sqrt{\hat{\theta}(1 - \hat{\theta})}},$$

and we reject  $H_0$  provided  $|W| > 1.96$ . (Exercise. Conduct a small survey and use the data to test the above hypothesis.)

## Wald with multivariate $\Theta$

Let now  $\theta$  be a  $p \times 1$  parameter and the model  $F_\theta$  is regular and s.t. the MLE has a normal limiting distribution, i.e.

$$\hat{\theta} \sim N_p(\theta, \hat{I}^{-1}).$$

Furthermore, let  $\theta_i$  denotes the  $i$ th component of  $\theta$  and  $\hat{\theta}_i$  denotes the  $i$ th component of  $\hat{\theta}$ .

To test the hypotheses  $H_0 : \theta_i = \theta_{i0}$  against  $H_1 : \theta_i \neq \theta_{i0}$ , at the level  $\alpha$ , where  $\theta_{i0}$  is a scalar fixed before observing the data, by the Wald test is to

reject  $H_0$  if  $\left| \frac{\hat{\theta}_i - \theta_{i0}}{\widehat{\text{se}}_i} \right| > z_{1-\alpha/2}$ ;

here  $\widehat{\text{se}}_i = \sqrt{\widehat{I}^{ii}}$ , or its equivalent version based on  $\hat{J}$ .



## Example 8 (Two classification algorithms)

We wish to compare the performance of two ML algorithms. Algo1 was run on a test set of size  $m$  and Algo2 was run on a test set of size  $n$ .

Let  $X$  and  $Y$  be the number of misclassifications with Algo1 and Algo2, resp. Assuming the sample is independent, then  $X_1, \dots, X_m$  are iid  $\text{Ber}(\theta_1)$  and  $Y_1, \dots, Y_n$  are iid  $\text{Ber}(\theta_2)$ .

The hypotheses of interest are then  $H_0 : \theta_1 = \theta_2$  vs  $H_1 : \theta_1 \neq \theta_2$ . Let  $\delta = \theta_1 - \theta_2$ , then above hypotheses translate to  $H_0 : \delta = 0$  vs  $H_1 : \delta \neq 0$ .

Let's apply a Wald test to  $\delta$ . First, we estimate  $\theta_1, \theta_2$  by MLE and then estimate  $\delta$  (using the equivariance principle) by  $\hat{\delta} = \hat{\theta}_1 - \hat{\theta}_2$ .

By the properties of the MLE,  $\hat{\delta}$  is asympt. normal and with variance

$$\widehat{\text{se}}(\hat{\delta})^2 = \frac{\hat{\theta}_1(1-\hat{\theta}_1)}{m} + \frac{\hat{\theta}_2(1-\hat{\theta}_2)}{n}$$

## Example 8 (cont'd)

The Wald test of approx. size  $\alpha$  is then to reject  $H_0$  if  $|\hat{\delta}/\widehat{\text{se}}(\hat{\delta})| > z_{1-\alpha/2}$ .

In a practical experiment conducted on  $m = 20, n = 30$  with Algo1 and Algo2 we obtained 5 and 10 misclassifications. Are the two algorithms performing equally?

The observed Wald statistic is  $w = -0.64$ .

Since  $w < 1.96$  we do not reject  $H_0$  and conclude that there is no evidence to support that the two algorithms have different classification performances.

## Example 9 (Computing $\alpha'$ and power of a test)

Consider  $X_1, \dots, X_n$  an iid random sample with  $X_i \sim N(\mu, \sigma^2)$ , where  $\sigma^2$  is known and suppose we wish to test  $H_0 : \mu = \mu_0$  against  $H_1 : \mu \neq \mu_0$ . Let's compute  $\alpha'$  and the power of the LRT.

The LRT at significance level  $\alpha$  is to reject  $H_0$  if  $|\sqrt{n}(\bar{X} - \mu_0)/\sigma| > z_{1-\alpha/2}$ . Now

$$\begin{aligned}\alpha' &= \sup_{\mu=\mu_0} P_{\mu}(\mathbf{X} \in R) = \sup_{\mu=\mu_0} P_{\mu}(|\sqrt{n}(\bar{X} - \mu)/\sigma| > z_{1-\alpha/2}) \\ &= P(|Z| > z_{1-\alpha/2}) = \alpha.\end{aligned}$$

On the other hand, for any  $\mu \neq \mu_0$

$$\begin{aligned}\beta(\mu) &= P_{\mu}(\mathbf{X} \notin R) = P_{\mu}(|\sqrt{n}(\bar{X} - \mu_0)/\sigma| \leq z_{1-\alpha/2}) \\ &= \Phi(z_{1-\alpha/2} + \sqrt{n}(\mu_0 - \mu)/\sigma) - \Phi(-z_{1-\alpha/2} + \sqrt{n}(\mu_0 - \mu)/\sigma)\end{aligned}$$

## Example 9 (cont'd)

Thus, the power of this test,  $\gamma(\mu) = 1 - \beta(\mu)$  depends on

- (i) the sample size  $n$ ,
- (ii) the population variance  $\sigma^2$ ,
- (iii) the population average  $\mu$  (and, of course  $\alpha$ ).

So to have a high chance of correctly rejecting  $H_0$ , we need to

- (i) use a large sample size  $n$ ,
- (ii) have  $\mu$  far from  $\mu_0$
- (iii) use  $\alpha$  that's not too small

# Pearson's $\chi^2$ test

Let  $(X_1, \dots, X_k) \sim \text{Mn}(n, \theta_1, \dots, \theta_k)$ , then  $\hat{\theta} = (\hat{\theta}_1, \dots, \hat{\theta}_k)$ , where  $\hat{\theta}_i = X_i/n$ , where  $n = \sum_i X_i$ . Suppose we want to test

$$H_0 : \theta = \theta_0 \quad \text{against} \quad H_1 : \theta \neq \theta_0,$$

where  $\theta_0 = (\theta_{01}, \dots, \theta_{0k})$  is a fixed vector prior to observing the data.

Consider the statistic

$$T = \sum_{i=1}^k \frac{(X_i - n\theta_{0i})^2}{n\theta_{0i}}.$$

It can be shown that  $T \xrightarrow{d} \chi_{k-1}^2$  as  $n \rightarrow \infty$ . The Pearson  $\chi^2$  test of approximate size  $\alpha$  is then to reject  $H_0$  if the observed value of  $T$  is greater than  $\chi_{k-1, 1-\alpha}^2$ .

## Example 10

It is conjectured that in the human population, 48% have blood type O, 38% have type A, 10% have type B and 7% have type AB. To test this hypothesis at a level  $\alpha = 0.05$  a sample of  $n$  people is taken, where  $n_1$  have type O,  $n_2$  have type A,  $n_3$  have type B and  $n_4$  have type AB.

The observed test statistic is

$$T^{obs} = \frac{(n_1 - 0.48n)^2}{0.48n} + \frac{(n_2 - 0.38n)^2}{0.38n} + \frac{(n_3 - 0.10n)^2}{0.10n} + \frac{(n_4 - 0.07n)^2}{0.07n}.$$

We reject  $H_0$  if  $T^{obs} > 7.815$ .

# An alternative view

Summing up, a test statistics is performed in three steps:

- (i) identify a test statistic for the parameter of interest and build a suitable rejection region  $R$ ;
- (ii) compute the test statistic at the observed sample to get the observed test statistic, say  $T^{obs}$ ;
- (iii) if  $T^{obs}$  follows in  $R$ , reject  $H_0$  otherwise accept it.

Many scholars do not find this binary choice very informative; they prefer to compute a  $p$ -value and take an action based on this.

# The $p$ -value

The  $p$ -value is defined as the smallest  $\alpha$  that leads to reject  $H_0$ . More formally, for every  $\alpha \in (0, 1)$ , if  $R_\alpha$  is a rejection region of size  $\alpha$ , then

$$p\text{-value} = \inf\{\alpha : T(X_1, \dots, X_n) \in R_\alpha\}.$$

A 'small'  $p$ -value indicates that  $H_0$  is not compatible with the data and it must be rejected.

A  $p$ -value could be treated as “small” when, for instance, is lower than  $\alpha$ , the size of type I error.



Concretely, if  $T_n = T(X_1, \dots, X_n)$  is a test statistic with observed value  $T^{obs} = T(x_1, \dots, x_n)$  and a size  $\alpha$  test has rejection region of the form

$$\{X_1, \dots, X_n : T(X_1, \dots, X_n) \geq c\},$$

then the  $p$ -value is defined by

$$\sup_{\theta \in \Theta_0} P_{\theta} (T(X_1, \dots, X_n) \geq T^{obs}).$$

On the other hand, if the rejection region is of the form

$$\{X_1, \dots, X_n : T(X_1, \dots, X_n) \leq c\},$$

then the  $p$ -value is defined by

$$\sup_{\theta \in \Theta_0} P_{\theta} (T(X_1, \dots, X_n) \leq T^{obs}).$$

Finally, if the rejection region is of the form

$$\{X_1, \dots, X_n : T(X_1, \dots, X_n) \leq c_1\} \cup \{X_1, \dots, X_n : T(X_1, \dots, X_n) \geq c_2\},$$

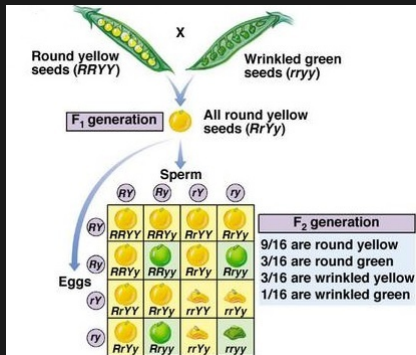
then

$$p\text{-value} = 2 \min \left( \sup_{\theta \in \Theta_0} P_{\theta}(T(X_1, \dots, X_n) \leq T^{obs}), \sup_{\theta \in \Theta_0} P_{\theta}(T(X_1, \dots, X_n) \geq T^{obs}) \right)$$

If  $\Theta_0 = \{\theta_0\}$  then replace  $\sup_{\theta \in \Theta_0} P_{\theta}$  by  $P_{\theta_0}$ .

## Example 11

Consider Mendel's experiment on peas, where round yellow seeds are breed with wrinkled green seeds. There are four types of progeny: round **yellow**, wrinkled **yellow**, round **green**, wrinkled **green**.



## Example 11 (cont'd)

Let  $(Y_1, Y_2, Y_3, Y_4)$  be vector with the numbers seeds of the four types. Then follows a multinomial distribution with parameter  $\theta = (\theta_1, \theta_2, \theta_3, \theta_4)$ .

Mendel's theory of inheritance predicts

$$\theta = \theta_0 = \left( \frac{9}{16}, \frac{3}{16}, \frac{3}{16}, \frac{1}{16} \right).$$

In  $n = 556$  his trials he obtained  $\mathbf{y} = (315, 101, 108, 32)$ .

Let's test if his theory is supported by the data using the LRT test. Let

$$H_0 : \theta = \theta_0 \text{ and } H_1 : \theta \neq \theta_0,$$

so at the observed data

$$\lambda(\mathbf{y}) = -2 \log \left( \frac{L(\hat{\theta}_0)}{L(\hat{\theta})} \right) = 0.48,$$

and  $p\text{-value} = P(\chi_3^2 \geq 0.48) = 0.92$ . (There are four parameters, but only three are free to vary and under  $H_0$  they are all fixed, thus  $df=3$ )

# Methods for evaluating tests

A testing problem can often be solved by means of several tests. The following two criteria are (the most widely) used for choosing the best one:

- prefer tests with smallest type I error, e.g. smallest  $\alpha'$  but still  $\alpha' \leq \alpha$ ;
- prefer tests with smallest type II error  $\beta(\theta)$ , or equivalently, prefer tests with highest power

$$\gamma(\theta) = 1 - \beta(\theta), \quad \forall \theta \in \Theta_0^c.$$

Due to the trade-off between the type I and type II errors, a test with  $\alpha'$  much smaller than  $\alpha$  will most likely have small power.

In some particular cases (simple null vs simple alternative) it is possible to show that a test based on the LRT statistic is the most powerful among all test; this result is known as the Newman-Pearson Lemma.

In general, this optimality result is hard or impossible to apply. However, it turns out that the LRT test and the Wald test tend to have decent power when  $n$  is sufficiently larger than the number of parameters.

To tell which test outperforms which in a given problem one often has to resort Monte Carlo methods as analytical calculations are typically impossible.

## Example 12 (The $t$ -test)

Let  $X_1, \dots, X_n$  be an iid random sample from  $N(\mu, \sigma^2)$ , with both parameters unknown; thus  $\theta = (\mu, \sigma^2)$ . We wish to test the hypothesis  $H_0 : \mu = \mu_0$  against  $H_1 : \mu \neq \mu_0$  via the LRT test.

Under  $H_0$  we have that

$$\sup_{\theta \in \Theta_0} L(\theta) = \frac{\exp\left[-\frac{1}{2\sigma_{\mu_0}^2} \sum_{i=1}^n (X_i - \mu_0)^2\right]}{(2\pi)^{n/2} \widehat{\sigma_{\mu_0}^2}^{n/2}},$$

where  $\widehat{\sigma_{\mu_0}^2} = \sum_{i=1}^n (X_i - \mu_0)^2 / n$ . Under  $H_1$  we have

$$\sup_{\theta \in \Theta} L(\theta) = \frac{e^{-n/2}}{(2\pi)^{n/2}} \left[ \frac{\sum_i (X_i - \bar{X})^2}{n} \right]^{-n/2}.$$

The LRT test is thus

## Example 12 (cont'd)

$$\begin{aligned}\frac{\sum_i (X_i - \mu_0)^2}{\sum_i (X_i - \bar{X})^2} \geq b &\iff \frac{n(\bar{X} - \mu_0)^2}{\sum_i (X_i - \bar{X})^2} \geq b - 1 \\ &\iff \frac{n(\bar{X} - \mu_0)^2(n-1)}{\sum_i (X_i - \bar{X})^2} \geq (b-1)(n-1) \\ &\iff \frac{n(\bar{X} - \mu_0)^2}{S^2} = T^2 \geq (b-1)(n-1) = d\end{aligned}$$

The rejection region for the LRT is of the type

$$\{\mathbf{X} : T_n^2 \geq d\} \equiv \{\mathbf{X} : |T_n| \geq \sqrt{d} = a\}$$

In order to define a size  $\alpha$  test we have to find  $a$  such that

$$\sup_{\theta \in \Theta_0} P_{\theta}(|T_n| \geq a) \leq \alpha.$$

But  $T_n \sim t_{n-1}$ , thus choosing  $a = t_{n-1, 1-\alpha/2}$  fills the bill.



## Example 12 (cont'd)

The test is thus:

Reject  $H_0$  if  $|T^{obs}|$  is greater than the quantile of level  $1 - \alpha/2$  of the  $t_{n-1}$  distribution,

where  $T^{obs} = \frac{\bar{X} - \mu_0}{S/\sqrt{n}}$  is the observed  $t$ -statistic. This is also known as the  $t$ -test.

The  $p$ -value for this test is  $2 \min(P(t_{n-1} > T^{obs}), P(t_{n-1} < T^{obs}))$ .

As a numerical example, suppose that an observed sample of  $n = 10$  WM's lead to  $\bar{x} = 201$  and  $s^2 = 5^2$  and suppose we wish to test  $H_0 : \mu = 200$  against  $H_1 : \mu \neq 200$  at the level  $\alpha = .05$ . Then

$$T^{obs} = \frac{\sqrt{10}(201-200)}{5} = 0.632.$$

Since  $t_{9,0.975} = 2.26 \not< 0.632$ , we do not reject  $H_0$  at level  $\alpha = 0.05$ . The  $p$ -value is  $2P(t_9 > 0.632) = 0.543$ , which suggests no evidence against  $H_0$ .

### Example 13 (LRT test for the variance)

Let  $X_1, \dots, X_n$  be an iid random sample from  $N(\mu, \sigma^2)$ , with both parameters unknown. We wish to test the hypothesis  $H_0 : \sigma^2 = \sigma_0^2$  against  $H_1 : \sigma^2 \neq \sigma_0^2$  via the LRT test. The LRT statistic is

$$\lambda(\mathbf{x}) = \frac{(2\pi\sigma_0^2)^{-n/2} \exp\left(-\frac{n\hat{\sigma}^2}{2\sigma_0^2}\right)}{(2\pi\hat{\sigma}^2)^{-n/2} \exp\left(-\frac{n}{2}\right)},$$

and the LRT test is to reject  $H_0$  if  $\lambda(\mathbf{x}) < c$ , or if

$$\left(\frac{\hat{\sigma}^2}{\sigma_0^2}\right)^{n/2} \exp\left(-\frac{n\hat{\sigma}^2}{2\sigma_0^2}\right) < b,$$

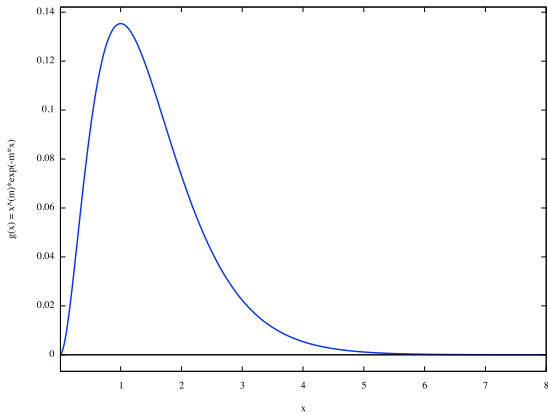
for some  $b$  depending on  $c$ .

## Example 13

Setting  $t = \hat{\sigma}^2 / \sigma_0^2$  this inequality is equivalent to

$$g(t) < b,$$

where  $g(t) = t^m e^{-mt}$ .



## Example 13

To compute the test we need a suitable value for  $b$  (and thus for  $c$ ).

If  $\alpha$  is the size of type I error, the usual way to fix  $b$  is by finding  $R$  s.t.

$$\sup_{\theta \in \Theta_0} P_{\theta}(\mathbf{X} \in R) \leq \alpha.$$

But

$$P_{\theta}(\mathbf{X} \in R) = P(g(\hat{\sigma}^2/\sigma_0^2) < b).$$

Note that  $g$  has unique maximum at  $t = 1$ . Furthermore,  $g$  is increasing for  $t < 1$  and decreasing for  $t > 1$ . Thus set  $R$  we are looking at is of the type

$$\{\mathbf{X} : \hat{\sigma}^2/\sigma_0^2 < b_1 \text{ or } \hat{\sigma}^2/\sigma_0^2 > b_2\},$$

which should have probability  $\alpha$  and  $b_1 < 1 < b_2$  s.t.  $g(b_1) = g(b_2) = b$ .

## Example 13

But then

$$\begin{aligned}\alpha &= P_{\theta_0}(X \in R) \\&= P_{\theta_0}(\mathbf{X} : \hat{\sigma}^2/\sigma_0^2 < b_1 \text{ or } \hat{\sigma}^2/\sigma_0^2 > b_2) \\&= P_{\theta_0}(\mathbf{X} : \hat{\sigma}^2/\sigma_0^2 < b_1) + P_{\theta_0}(\mathbf{X} : \hat{\sigma}^2/\sigma_0^2 > b_2) \\&= P_{\theta_0}(\mathbf{X} : n\hat{\sigma}^2/\sigma_0^2 < nb_1) + P_{\theta_0}(\mathbf{X} : n\hat{\sigma}^2/\sigma_0^2 > nb_2) \\&= P(\chi_{n-1}^2 < a_1) + P(\chi_{n-1}^2 > a_2),\end{aligned}$$

where the last equality is due to the fact that  $n\hat{\sigma}^2/\sigma^2 \sim \chi_{n-1}^2$  under a normal iid sample.

There are two ways to compute  $a_1, a_2$ :

- (i) using quantiles of the  $\chi_{n-1}^2$  (easy to compute)
- (ii) numerical inversion (harder, better power, homework).

## Example 13 (cont'd)

Setting  $a_1 = \chi_{n-1, \alpha/2}^2$  and  $a_2 = \chi_{n-1, 1-\alpha/2}^2$ , leads to the test

Reject  $H_0 : \sigma^2 = \hat{\sigma}_0^2$  if the observed value of  $n\hat{\sigma}^2/\sigma_0^2$  is  $< \chi_{n-1, \alpha/2}^2$  or  $> \chi_{n-1, 1-\alpha/2}^2$ .

Suppose the observed sample is

2.51, 1.31, 1.55, 1.87, 1.64, 4.00, 3.09, 2.81, 5.81, 3.72

and let's test

$H_0 : \sigma^2 = 3$  vs  $H_1 : \sigma^2 \neq 3$  at  $\alpha = 0.05$ .

Using the observed sample we have  $n\hat{\sigma}^2/\sigma_0^2 = 5.83$ . Also  $\chi_{9, 0.025}^2 = 2.7$  and  $\chi_{9, 0.975}^2 = 19.02$ , thus we don't reject  $H_0$ .

The p-value of the test is

$$2 \min(P(\chi_{9, 0.975}^2 < 5.83), P(\chi_{9, 0.975}^2 > 5.83)) = 0.486.$$

# Extensions

Let  $X_1, \dots, X_m$  and  $Y_1, \dots, Y_n$  be two iid random samples with  $X_i \sim N(\mu_x, \sigma_x^2)$ ,  $Y_j \sim N(\mu_y, \sigma_y^2)$  and  $X_i$  is independent from  $Y_j$ , all parameters unknown:

- (i) test for  $H_0 : \mu_x = \mu_y$  against  $H_1 : \mu_x \neq \mu_y$ ; it's known as the two-sample  $t$ -test, single-tailed versions are also possible.
- (ii) test for  $H_0 : \mu_x = \mu_y$  against  $H_1 : \mu_x \neq \mu_y$  when  $X_i, Y_j$  are dependent and  $n = m$ ; known as paired  $t$ -test.
- (iii) test for  $H_0 : \sigma_x^2 = \sigma_y^2$  against  $H_1 : \sigma_x^2 \neq \sigma_y^2$ ; easier if restated as ratio of variances.
- (iv) three samples, say  $W_1, \dots, W_r$  iid and  $W_h \sim N(\mu_w, \sigma_w^2)$  indep., it may be of interest to test  $H_0 : \mu_x = \mu_y = \mu_w$ ; this is called ANOVA (analysis of variance) test.

All may be derived by means of the likelihood ratio test. The LRT test in (i) gives an exact test under the assumption  $\sigma_x^2 = \sigma_y^2$ .

Indeed ...

### Example 14 (Two-sample t-test)

By a reasoning similar to that of Example 12 and assuming  $\sigma_x^2 = \sigma_y^2$ , it is possible to show that the LRT statistic for the two-sample problem is equivalent to the statistic

$$T_n = \frac{\bar{X} - \bar{Y}}{\sqrt{S_p^2 \left( \frac{1}{m} + \frac{1}{n} \right)}}, \Rightarrow \text{distr. } t_{n+m-2}$$

*gradi di libert *

where  $S_{pool}^2$  is the pooled variance estimator

$$S_p^2 = \frac{(n-1)S_y^2 + (m-1)S_x^2}{n+m-2}$$

The test of size  $\alpha$  is to reject  $H_0 : \mu_x = \mu_y$  if  $|T_n^{obs}| \geq t_{n+m-2, \alpha/2}$

The  $p$ -value is computed by  $P(|t_{n+m-2}| \geq |T_n^{obs}|)$ .



## Note to two-sample t-test

The variance homogeneity assumption, i.e.  $\sigma_x^2 = \sigma_y^2$ , is often not realistic. With this assumption relaxed, we can use the test statistic

$$T'_n = \frac{\bar{X} - \bar{Y}}{\sqrt{S_x^2/m + S_y^2/n}},$$

and reject  $H_0$  for large values of  $|T'_n|$ .

There is no exact test for this problem (aka Behrens–Fisher problem), but an approximate test of size  $\alpha$  that's extremely accurate is: Reject  $H_0$  if  $|T'_n| > t_{\nu, 1-\alpha/2}$ , where

$$\nu = \frac{(S_x^2/m + S_y^2/n)^2}{S_x^4/(m^2(m-1)) + S_y^4/(n^2(n-1))}.$$

All one and two-sample version of the t-test are implemented in the `t.test` function of R.

## Example 15

paired t-test Let  $(X_1, Y_1), \dots, (X_n, Y_n)$  be an iid random sample where  $(X_i, Y_i)$  follows a bivariate distribution with mean  $\mu = (\mu_x, \mu_y)$  and covariance matrix  $\Sigma$ .

Often we assume that the bivariate distribution is a bivariate normal, and we are interested in testing

$$H_0 : \mu_x = \mu_y \text{ against } H_1 : \mu_x \neq \mu_y.$$

In principle we could apply the LRT, but the easiest solution is to work with the differences. Indeed, if  $D_i = X_i - Y_i$ , then

$$\mu_D = E(D_i) = \mu_x - \mu_y, \text{ and } \text{var}(D_i) = \text{var}(X_i) + \text{var}(Y_i) - 2\text{cov}(X_i, Y_i).$$

Then we can apply a t-test (Example 12) for  $H_0 : \mu_D = 0$  vs  $H_1 : \mu_D \neq 0$  using the sample  $D_1, \dots, D_n$ .