## MS = HYPOTHESIS TESTING

In this chapter, we introduce parametric hypothesis testing. Non-parametric techniques are discussed in the next chapter.

## I. SIMPLE HYPOTHESIS TESTING

### I.1. The general principle.

x <u>Introductive example</u>: The crop that we wish to grow is known to give the best yield in soils with a pH of 7 (i.e. neutral). The pH of the soil was measured at various locations, giving observations

$$\{6.0 \quad 5.7 \quad 6.2 \quad 6.3 \quad 6.5 \quad 6.4 \quad 6.9 \quad 6.6$$
$$6.8 \quad 6.7 \quad 6.8 \quad 7.1 \quad 6.8 \quad 7.1 \quad 7.1 \quad 7.5$$
$$7.0\} = \{x_1, \dots, x_n\}, \quad n = 17.$$

The sample mean pH of the soil is $\bar{x} := \frac{1}{n}\sum_{i=1}^{n} x_i = 6.676$. Should we add some chemicals to change the pH of the soil? In other words, is the population mean pH of the soil different from 7?

↳ Thus assuming that our observations are a realization of a random sample $\mathcal{L}_n = \{X_1, \dots, X_n\}$, where the $X_i$ are iid, with mean $\mathbb{E}X$.

↳ We first ask another question: what is the probability of observing a sample mean as small as 6.676, if the population mean is 7?

To answer this question, we need to specify a distribution for the soil pH. Assume that the $X_1, \dots, X_n$ are $\mathcal{N}(\mu, \sigma = 0.5)$. Then $\bar{X} \sim \mathcal{N}(\mu, \frac{0.5}{\sqrt{17}})$, and we find that $\mathbb{P}(\bar{X} \leqslant 6.676) = 0.004$, with $\mu = 7$.

This probability is very small. So...

(i) Either $\mu = 7$, and we have observed something very unlikely, or

(ii) Our assumption about the mean pH is wrong, or

(iii) Our assumption about the model is wrong.

Here, we would most likely conclude that $\mu < 7$, and add some chemicals.

• In simple hypothesis testing, we compare two hypotheses:

↳ The **NULL HYPOTHESIS** (denoted $H_0$) = the statement whose validity is to be tested. Often the null hypothesis can be expressed in terms of parameters of a model (e.g. $H_0: \mu = 7$).

↖ An example of a <u>simple</u> hypothesis: the parameter of the distribution is specified. As opposed to a <u>composite</u> hypothesis, for which the parameter of the distribution is not completely specified (e.g: $\mu < 7$)

The null hypothesis often expresses an absence of effect, a reference, a "status quo".

Ex: x amount of savings of customers in a bank is equivalent to the amount of savings of customers of another bank.

x prices of a major sports brand are identical to the prices of its main competitor

↘ The **ALTERNATIVE HYPOTHESIS** (denoted $H_1$) = it specifies what happens if $H_0$ is false. $H_1$ often specifies what we hope, or expect, to be true.

The alternative hypothesis often expresses a difference, a departure from a reference, the presence of an effect.

Ex: $(H_1 : \mu = \mu_1)$, with $\mu_1 \neq 0$ for $(H_0 : \mu = 0.)$

Back to our example, we may consider

$$H_0 : \mu = 7$$
$$H_1 : \mu < 7$$

or

$$H_0 : \mu = 7$$
$$H_1 : \mu \neq 7 \quad \text{"two-sided"}$$

↑ if we suspect that the soil is acidic

↑ if we are unable (or unwilling) to specify the direction in which the mean pH may differ from 7.

"one-sided"

× **Classical approach to Hypothesis Testing (HT):**

- **Rejecting the null** : Rejecting $H_0$ when it is true is called a **type I error**. Its probability $\alpha$ is called the **SIGNIFICANCE LEVEL**. Denoting $H_0 : \theta = \theta_0$, we have
$$\alpha = \mathbb{P}_{\theta_0}(\text{reject } H_0)$$
↖ the probability is computed under the null distribution

(We have not specified yet a rule for rejecting $H_0$)

- **Failing to reject the null**: Failing to reject $H_0$ when it is false is called a **type II error**. Its probability is usually denoted $\beta$. The quantity $1-\beta$ is referred to as the **POWER** of the test, and corresponds to the probability of rejecting $H_0$ when it is indeed false. To compute $\beta$, the

---

alternative hypothesis must be exactly specified:
e.g. $H_1 : \theta = \theta_1$, so that $\beta = \mathbb{P}_{\theta_1}(\text{fail to reject } H_0)$.

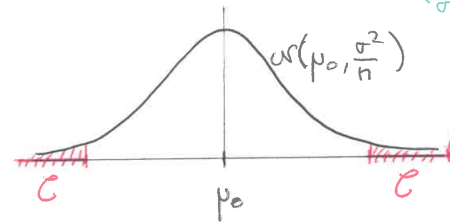If $\theta_1$ is not specified, $\beta$ cannot be calculated.

- The decision of a test (reject or not) is often based on determining a **critical region** $\mathcal{C}$. A statistic $T(X_1, \dots, X_n) \in \mathcal{C}$ is thought to be unlikely to have occured if $H_0$ is true:
$$\text{reject } H_0 \iff T(X_1, \dots, X_n) \in \mathcal{C}$$

× **Example:** Testing for $H_0 : \mu = \mu_0$, under the assumption that $X_1, \dots, X_n$ are $\mathcal{N}(\mu, \sigma^2)$ iid. The sample mean $\overline{X} = \frac{1}{n} \sum_{i=1}^{n} X_i = T(X_1, \dots, X_n)$ is $\mathcal{N}(\mu, \frac{\sigma^2}{n})$ distributed.

⇒ Reject $H_0$ if the observed value $\overline{x}$ of $\overline{X}$ falls in the tails of the $\mathcal{N}(\mu_0, \sigma^2)$ distribution. For example,

↖ $\sigma$ is assumed to be known.



$\mathcal{N}(\mu_0, \frac{\sigma^2}{n})$

observing $\overline{x}$ here is highly unlikely if the data $X_1, \dots, X_n$ arised from the $\mathcal{N}(\mu_0, \sigma^2)$ distribution.
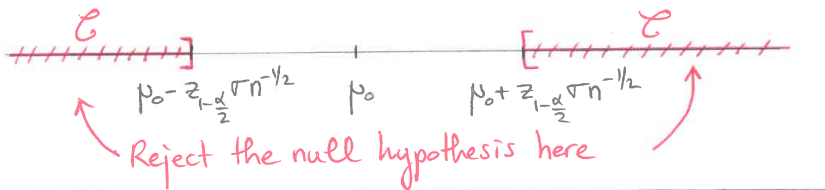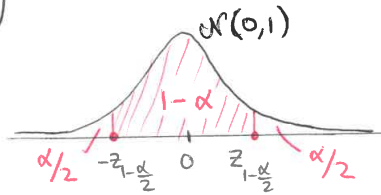⇒ Reject $H_0$

The significance level $\alpha$ of the test corresponds precisely to the probability that $T(X_1, \dots, X_n) \in \mathcal{C}$. It is usually fixed in advance by the practitioner. [Common values are $\alpha = 0.1$ or $\alpha = 0.05$.]
↳ Let $z_{1-\alpha/2}$ be the $(1 - \frac{\alpha}{2})$-quantile of the standard normal distribution, we

obtain $\quad \mathscr{E} = \left(-\infty, \ \mu_0 - z_{1-\frac{\alpha}{2}}\sigma n^{-1/2}\right] \cup$

$$\left[\mu_0 + z_{1-\frac{\alpha}{2}}\sigma n^{-1/2}, \ +\infty\right),$$

since

$$\mathbb{P}_{\mu_0}\left(\ \left|\frac{\overline{X}-\mu_0}{\sigma n^{-1/2}}\right| > z_{1-\frac{\alpha}{2}}\right) = \alpha$$

"$Z \sim \mathcal{N}(0,1)$"



$\mathcal{N}(0,1)$

$1-\alpha$

$\frac{\alpha}{2}$ $\quad -z_{1-\frac{\alpha}{2}}$ $\quad 0 \quad$ $z_{1-\frac{\alpha}{2}}$ $\quad \frac{\alpha}{2}$

$\mathscr{E}$ $\qquad \mathscr{E}$

$\mu_0 - z_{1-\frac{\alpha}{2}}\sigma n^{-1/2}$ $\quad \mu_0 \quad$ $\mu_0 + z_{1-\frac{\alpha}{2}}\sigma n^{-1/2}$

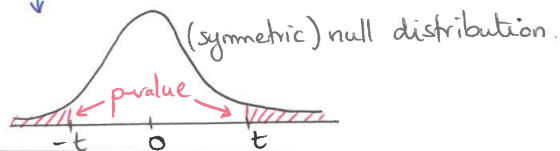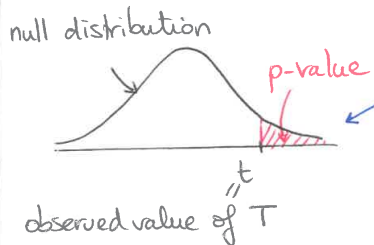Reject the null hypothesis here

Instead of specifying $\alpha$, one often reports the p-VALUE = the probability of observing a value of the test statistic $T(X_1,..,X_n)$ as or more <u>extreme</u> than the one actually observed, assuming $H_0$ is true.

there is some flexibility in what is meant by 'more extreme':
$$\mathbb{P}_{H_0}(T \geq t) \text{ or } \mathbb{P}_{H_0}(T \leq t) \text{ for}$$
a tail event; or
$$2\min\left(\mathbb{P}_{H_0}(T \geq t), \ \mathbb{P}_{H_0}(T \leq t)\right) \text{ for}$$
a double tail event.

null distribution



p-value

observed value of $T$

(symmetric) null distribution.

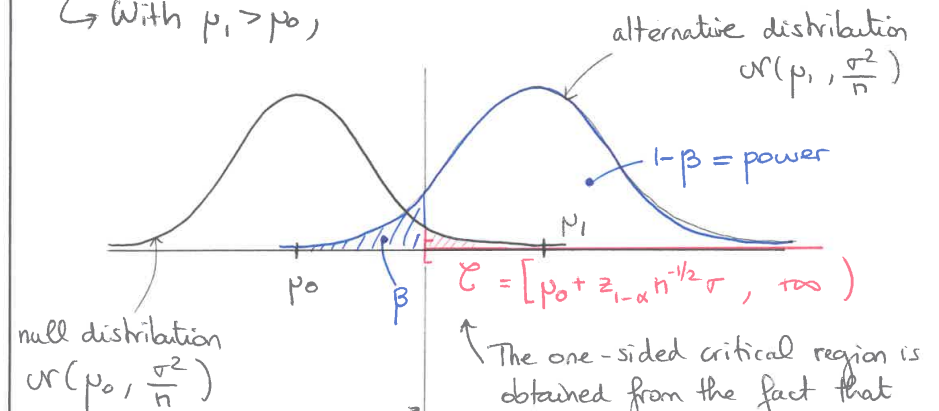p-value

$-t \quad 0 \quad t$

---

$H_1$ determines what is meant by "more extreme":

Tail events correspond to one-sided alternatives, such as $H_1 : \mu < \mu_0$ or $H_1 : \mu > \mu_0$; while double tail events correspond to two-sided alternatives; $H_1 : \mu \neq \mu_0$.

Assuming $H_1 : \mu = \mu_1$, we can compute the power $(1-\beta)$ of the test, since
$$\beta = \mathbb{P}_{\mu_1}(\text{fail to reject } H_0)$$
$$= \mathbb{P}_{\mu_1}(T(X_1,..,X_n) \in \mathscr{E})$$

$\hookrightarrow$ With $\mu_1 > \mu_0$,



alternative distribution $\mathcal{N}\left(\mu_1, \frac{\sigma^2}{n}\right)$

$1-\beta = $ power

$\mu_1$

$\mathscr{E} = \left[\mu_0 + z_{1-\alpha}n^{-1/2}\sigma, \ +\infty\right)$

$\mu_0$

$\beta$

null distribution $\mathcal{N}\left(\mu_0, \frac{\sigma^2}{n}\right)$

the threshold is fixed by the critical region constructed from $H_0$.

The one-sided critical region is obtained from the fact that under $H_0$,
$$\mathbb{P}_{\mu_0}\left(\frac{\overline{X}-\mu_0}{\sigma n^{-1/2}} > z_{1-\alpha}\right) = \alpha$$
$\sim \mathcal{N}(0,1)$

<u>In summary</u>, when conducting a HT, one aims at
$\searrow$ Minimizing the type I error (aka $\alpha$)
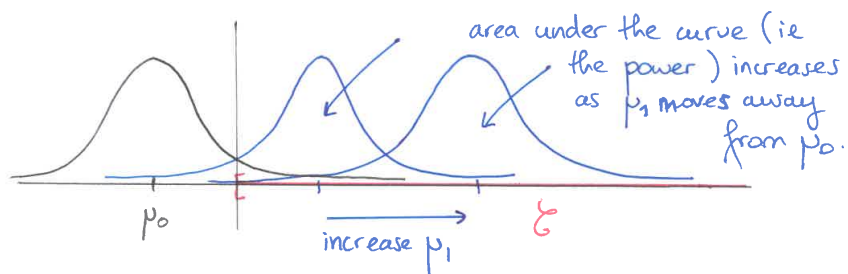$\searrow$ Maximizing the power (aka $1-\beta$).

The power of a statistical test is calculated for a particular alternative hypothesis $H_1 : \mu = \mu_1$.

$\Rightarrow$ Compute the value of $(1-\beta)$ for a range of simple alternative hypothesis.
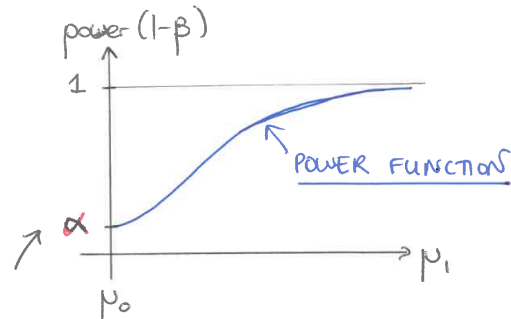
<u>Ex</u>: Compute the value of $(1-\beta)$ under $H_1 : \mu = \mu_1$, for all $\mu_1 > \mu_0$.

$\hookrightarrow$ This leads us to the concept of <u>POWER FUNCTIONS</u>.

Note that as $\mu_1$ increases, the power of the test increases as the alternative distribution is shifted to the right.



area under the curve (ie the power) increases as $\mu_1$ moves away from $\mu_0$.

$\mu_0$    increase $\mu_1$    $\xi$

$\Rightarrow$ Plot the value of the power as a function of $\mu_1$, for $\mu_1 \in [\mu_0, \infty)$ :
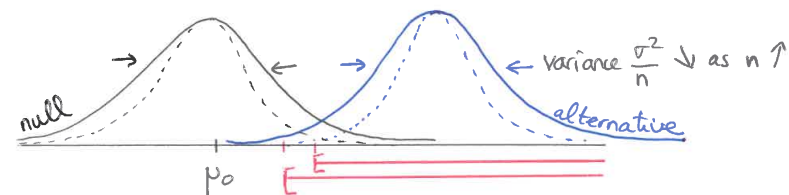


power $(1-\beta)$

POWER FUNCTION

At $\mu_1 = \mu_0$, the power of the test is equal to $\alpha$, as the null & alternative distributions coincide

---

The sample size has an effet on the power of a test:

For a given $\alpha$, as $n$ increases, the power increases. Indeed, in our previous example, the boundary is shifted to the left as $n \uparrow$, while the variance $\frac{\sigma^2}{n}$ of the normal distributions decreases:



null    $\leftarrow$ variance $\frac{\sigma^2}{n} \downarrow$ as $n \uparrow$    alternative

$\mu_0$

As $n \uparrow$, the end point $\mu_0 + z_{1-\alpha} n^{-1/2} \sigma$ decreases towards $\mu_0$

$\Rightarrow$ One can calculate the minimum sample size required to achieve a pre-specified power, for a given alternative hypothesis.

<u>Ex</u>: Looking back at the picture on page 6, the power of the test considered is

$$1 - \beta = \mathbb{P}_{\mu_1} \left( Y \geqslant \mu_0 + z_{1-\alpha} n^{-1/2} \sigma \right), \text{ where } Y \sim \mathcal{N}\left(\mu_1, \frac{\sigma^2}{n}\right)$$

$$= \mathbb{P}_{\mu_1} \left( \underbrace{\frac{Y - \mu_1}{\sigma n^{-1/2}}}_{\sim \mathcal{N}(0,1)} \geqslant \frac{\mu_0 - \mu_1}{\sigma n^{-1/2}} + z_{1-\alpha} \right)$$
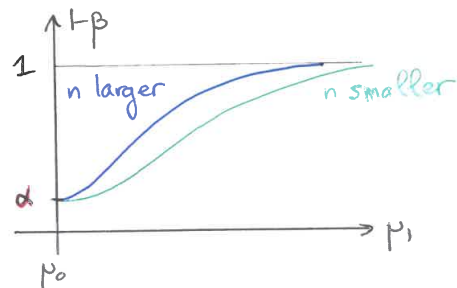
Since $\mu_0 - \mu_1 < 0$, as $n \uparrow$, the power of the test tends to $1$.

Denoting $z_\beta$ the $\beta$-quantile of the standard normal

distribution, we need to solve $\quad \dfrac{\mu_0 - \mu_1}{\sigma n^{-1/2}} + z_{1-\alpha} = z_\beta$

$\Rightarrow$ Choose $n \geqslant \left( \dfrac{\mu_0 - \mu_1}{\sigma (z_\beta - z_{1-\alpha})} \right)^2$ to ensure a power of at least $1-\beta$.



Remark: The test can be conducted without computing its power. However, you have no theoretical guarantee that you are doing something meaningful.

Summary: When designing a HT, you need to

(i) Decide on the null & alternative hypothesis.

(ii) Fix a desired level of type I error (aka $\alpha$)

(iii) Construct a test statistic $T(X_1, .., X_n)$ from the data $X_1, .., X_n$, and such that the distribution of $T$ is known under $H_0$.

(iv) Construct a rejection region based on your choices (i) (ii), (iii).

(v) Optional (but recommended): the sample size should be chosen to ensure that the type II error remains small (aka $\beta$) / the power $(1-\beta)$ is high.

(vi) Conclude: Reject $H_0$, or not ($\equiv$ presence of an effect, or not).

---

Remark: Hypothesis testing was introduced by Ronald Fisher in 1925. His approach however differs from the one presented in this section, in that he did not consider an alternative hypothesis to the null. Instead, given $H_0$, and a test statistic $T(X_1, .., X_n)$, Fisher suggested calculating the p-value, without the need to fix a desired level of significance $\alpha$. A result with a low p-value is taken as statistical evidence against the null.

On the top of page 6, we interpreted tail-events and double tail events in terms of a one-sided or two-sided alternative. This association is however superfluous. For Fisher, alternatives to $H_0$ are implicitely "all what is not $H_0$". Mathematically, this would translate as $H_1 : \mu \neq \mu_0$, if $H_0 : \mu = \mu_0$.

Type I errors and type II errors were introduced later by Neyman & Pearson, with the concept of an alternative hypothesis: given a fixed level $\alpha$, they advocate selecting the test that has the most power.

And indeed, we will see in section I. that to a given problem, you can construct a serie of tests with the same significance level $\alpha$. The test that should be retained, according to Neyman & Pearson, is the one with maximum power.

# I.2. Further examples.

x **Example 1**: Testing for the mean of a normal population, when the variance is unknown.

Testing $\left(H_0 : \mu = \mu_0\right)$ for $X_1, \cdots, X_n$ iid $\mathcal{N}(\mu, \sigma^2)$, $\sigma$ unknown.

Then we know from the result on page 25/26 in <span style="color:green">MS = PARAMETRIC INFERENCE</span> that $\overline{X} := \frac{1}{n}\sum_{i=1}^{n} X_i$ and $S^2 := \frac{1}{n-1}\sum_{i=1}^{n}(X_i - \overline{X})^2$

are independent, and that $\overline{X} \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$

$$\& \quad \frac{(n-1)S^2}{\sigma^2} \sim \chi^2(n-1).$$

It follows that under $H_0$,

$$T := \frac{(\overline{X} - \mu_0)/\sigma\sqrt{n}}{\sqrt{\frac{(n-1)S^2}{\sigma^2}/(n-1)}} = \frac{n^{1/2}(\overline{X} - \mu_0)}{S} \sim t(n-1).$$

$\Rightarrow$ Consider the test statistic $T(X_1, \cdots, X_n) = \dfrac{n^{1/2}(\overline{X} - \mu_0)}{S}$,
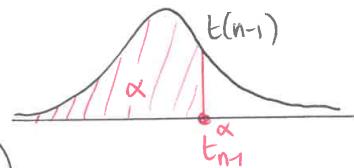
whose distribution under $H_0$ is $t(n-1)$

The rejection region $C$ is constructed from the quantiles of the $t(n-1)$ distribution: denoting $t_{n-1}^\alpha$ the $\alpha$-quantile, we obtain

$$C = \left(-\infty, \mu_0 - t_{n-1}^{1-\alpha/2} S n^{-1/2}\right]$$
$$\cup \left[\mu_0 + t_{n-1}^{1-\alpha/2} S n^{-1/2}, +\infty\right)$$

<span style="color:blue">for a two-sided alternative. The case of a one-sided alternative is treated similarly.</span>

x **Example 2**: Testing for a proportion.

Testing $\left(H_0 : p = p_0\right)$ for $X_1, \cdots, X_n$ iid $B(p)$.

The estimator $\hat{p} = \overline{X} = \frac{1}{n}\sum_{i=1}^{n} X_i$ is a consistent estimator of $p$.

Moreover, under $H_0$,

$$Z := \frac{\hat{p} - p_0}{\sqrt{np_0(1-p_0)}} \xrightarrow{d} \mathcal{N}(0,1) \quad \text{as } n \to +\infty.$$

Therefore, assuming that $n$ and $p_0$ are such that the normal approximation holds ( <span style="color:blue">a rule of thumb : $np_0 \geqslant 5$, and $n(1-p_0) \geqslant 5$</span> ) we can use $Z$ as our test statistic, and construct a rejection region based on the standard normal distribution. As usual, denoting $z_\alpha$ the $\alpha$-quantile of $\mathcal{N}(0,1)$, we obtain the (two-sided) critical region

$$C = \left(-\infty, p_0 - z_{1-\frac{\alpha}{2}} n^{-1/2}(p_0(1-p_0))^{1/2}\right]$$
$$\cup \left[p_0 + z_{1-\frac{\alpha}{2}} n^{-1/2}(p_0(1-p_0))^{1/2}, +\infty\right),$$

with nominal level $\alpha$.
<span style="color:blue">(and similarly for a one sided region )</span>

$\rightarrow$ We discuss next common approaches for constructing test statistics :
— Wald tests
— likelihood ratio tests

## I.3. The Wald test.

Consider a random sample $X_1, \ldots, X_n$ iid $\sim P_\theta$, for $\theta \in \Theta \in \mathbb{R}^d$.

Test $(H_0 : \theta = \theta_0)$ for some fixed $\theta_0 \in \mathbb{R}^d$.

Let $\hat{\theta}_{ML}$ = maximum likelihood estimate of $\theta$

Under some technical conditions, the MLE is consistent and asymptotically normally distributed, so that we can write:

$$n^{1/2} I_d(\hat{\theta}_{ML})^{1/2} (\hat{\theta}_{ML} - \theta_0) \xrightarrow{d} \mathcal{N}(0, I_d).$$

The multivariate version of the result on pages 15/16 in MS = MAX. LIK. EST.

$I_d(\hat{\theta}_{ML})$ is the $(d \times d)$ Fisher matrix, evaluated at the MLE.

$$I_d(\theta) = (I_{ij}(\theta)), \text{ with } I_{ij}(\theta) = \mathbb{E}\left\{\frac{\partial \log f(\underline{X};\theta)}{\partial \theta_i} \frac{\partial \log f(\underline{X};\theta)}{\partial \theta_j}\right\}$$

$$f(\underline{X};\theta) = \prod_{i=1}^{n} f(X_i;\theta) = \text{joint density of } \underline{X} = (X_1, \ldots, X_n).$$

Also, $A^{1/2}$ denotes the square root of the matrix $A$.

It follows that

$$\underbrace{n (\hat{\theta}_{ML} - \theta_0)^t I_d(\hat{\theta}_{ML}) (\hat{\theta}_{ML} - \theta_0)}_{T(X_1, \ldots, X_n) = \underline{\text{Wald test statistic.}}} \xrightarrow{d} \chi^2(d)$$

Test has significance level $\alpha \in (0,1)$ for the rejection region $[q^d_{1-\alpha}, +\infty)$, where $q^d_\alpha = \alpha$-quantile of the $\chi^2(d)$ distribution.

---

## I.4. Likelihood ratio tests.

Consider a random sample $\underline{X} = (X_1, \ldots, X_n)$, where the $X_i$ are iid and $P_\theta$ distributed, for some $\theta \in \Theta$. The probability measure $P_\theta$ is assumed to have density $f_\theta(x) = f(x;\theta)$

( likelihood ratio tests are presented in the context of AC RVs, but the results hold true for discrete RVs as well ).

Put
$$\bullet \quad L(\underline{x};\theta) := \prod_{i=1}^{n} f(x_i;\theta)$$

$$\bullet \quad \ell(\underline{x};\theta) := \log L(\underline{x};\theta) = \sum_{i=1}^{n} \log f(x_i;\theta)$$

$$\underline{x} = (x_1, \ldots, x_n)$$

Let $\Theta_0$ be a subset of $\Theta$, and consider the MLEs $\hat{\theta}_0$ and $\hat{\theta}$, computed over the parameter spaces $\Theta_0$ and $\Theta$, respectively:

$$\bullet \quad \hat{\theta}_0 = \operatorname*{argmax}_{\theta \in \Theta_0} L(\underline{x};\theta)$$

$$\bullet \quad \hat{\theta} = \operatorname*{argmax}_{\theta \in \Theta} L(\underline{x};\theta).$$

The ratio $\left| \Lambda(\underline{x}) := \frac{L(\underline{x};\hat{\theta}_0)}{L(\underline{x};\hat{\theta})} \right.$ is called the

LIKELIHOOD RATIO STATISTIC.

A few observations:

↘ $0 \leq \Lambda(\underline{x}) \leq 1$

↘ If $\hat{\theta}_0$ is far from $\hat{\theta}$, then expect $\Lambda(\underline{x})$ to be small.

↘ If the true parameter is in $\Theta_0$, then $\Lambda(\underline{x}) \to 1$ as $n \to +\infty$.

$\Rightarrow$ The likelihood ratio statistic may be used as a test statistic to test for $H_0 : \theta \in \Theta_0$, by comparing $\Lambda(x)$ to some threshold $c$.

It is common to consider as well the quantity

$$\lambda(x) := -2 \log \Lambda(x)$$
$$= 2 \left( \ell(\underline{x}; \hat{\theta}) - \ell(\underline{x}; \hat{\theta}_0) \right)$$

x **Example:** Let $X_1, \ldots, X_n$ iid $\sim B(p)$, for $p \in \Theta = \{p \mid p_0 \leq p < 1\}$,

for some $p_0 \in (0,1)$.

Consider $(H_0 : p = p_0)$.

↖ so that $\Theta_0 = \{p_0\}$

We have :

- $\hat{p}_0 := \underset{p \in \Theta_0}{\text{argmax}} \ L(\underline{x}; p) = p_0$

  since $\Theta_0$ contains a single element

- $\hat{p} := \underset{p \in \Theta}{\text{argmax}} \ L(\underline{x}; p)$

$$= \begin{cases} p_0 & \text{if} \quad \bar{x} \leq p_0 \\ \bar{x} := \frac{1}{n} \sum_{i=1}^{n} x_i & \text{if} \quad \bar{x} > p_0 \end{cases}$$

$\Rightarrow \Lambda(\underline{x}) = \begin{cases} 1 & \text{if} \quad \frac{x}{n} \leq p_0 \\ \dfrac{p_0^x (1-p_0)^x}{\left(\frac{x}{n}\right)^x \left(1-\frac{x}{n}\right)^{n-x}} & \text{if} \quad \frac{x}{n} > p_0 \end{cases}$

$X := \sum_{i=1}^{n} X_i$

↑ The LR may have a complicated form …



1

$x$

0         $n$

We see that $\Lambda(\underline{x})$ is a decreasing function of $x = \sum_{i=1}^{n} x_i$.

$\Rightarrow$ Reject $H_0$ if $\Lambda(\underline{x})$ is smaller to some $c$ $\iff$ $x$ is larger to some value. Alternatively, we may compute the p-value $p := \mathbb{P}_{p_0} \left( \Lambda(\underline{X}) \leq \Lambda(\underline{x}) \right)$

$$= \mathbb{P}_{p_0} \left( X \geq x \right)$$

↖ where $X = \sum_{i=1}^{n} X_i$, $x = \sum_{i=1}^{n} x_i$.

x **Example:** $X_1, \ldots, X_n$ iid $\sim \mathcal{N}(\mu, \sigma^2)$, where $\mu \in \Theta = \mathbb{R}$.

$>0$, assumed known.

Consider $(H_0 : \mu = \mu_0)$, for some $\mu_0 \in \mathbb{R}$.

↖ So that $\Theta_0 = \{\mu_0\}$.

Then

- $\hat{\mu}_0 := \underset{\mu \in \Theta_0}{\text{argmax}} \ L(\underline{x}; \mu) = \mu_0$

- $\hat{\mu} := \underset{\mu \in \Theta}{\text{argmax}} \ L(\underline{x}; \mu) = \bar{x}$

Thus $\lambda(\underline{x}) = 2 \left( \ell(\underline{x}; \hat{\mu}) - \ell(\underline{x}; \hat{\mu}_0) \right)$

$$= \frac{n}{\sigma^2} (\bar{x} - \mu_0)^2$$

Note that $\lambda(\underline{X})$ has a $\chi^2(1)$ distribution, since

$\lambda(\underline{X}) = Z^2$; with $Z := \dfrac{\bar{X} - \mu_0}{\sigma/\sqrt{n}} \sim \mathcal{N}(0,1)$.

$= \dfrac{\bar{x} - \mu_0}{\sigma/\sqrt{n}}$

Moreover, since $\lambda(\underline{X}) \geq \lambda(\underline{x}) \iff |Z| \geq |z|$,

the p-value is $p := \mathbb{P}_{\mu_0} \left( \lambda(\underline{X}) \geq \lambda(\underline{x}) \right)$

$$= \mathbb{P}_{\mu_0} \left( |Z| \geq |z| \right) = 2 \left( 1 - \Phi(|z|) \right)$$

**Theorem:** Let $\theta_0$ be an interior point of $\theta \in \mathbb{R}^d$. Under some technical assumptions (the same as the ones needed to compute the asymptotic properties of the MLE), if the null ($H_0: \theta = \theta_0$) is true, then the log-likelihood ratio statistic $\lambda(\underline{X}) := -2\log \Lambda(\underline{X})$ has a limiting $\chi^2$ distribution, with $d$ degrees of freedom.

The limiting $\chi^2(d)$ distribution may be used to decide on a threshold for example.

**proof:** We have $\lambda(\underline{x}) = 2\left( \ell(\underline{x}; \hat{\theta}) - \ell(\underline{x}; \theta_0) \right)$.

Expanding $\ell(\underline{x}; \theta_0)$ around $\hat{\theta}$, we get

$$\ell(\underline{x}; \theta_0) = \ell(\underline{x}; \hat{\theta}) + (\theta_0 - \hat{\theta})^t \underbrace{\ell'(\underline{x}; \hat{\theta})}_{=0}$$

$$+ \frac{1}{2}(\theta_0 - \hat{\theta})^t \underbrace{\ell''(\underline{x}; \tilde{\theta})}_{(d \times d)} (\theta_0 - \hat{\theta}).$$

where $\tilde{\theta}$ is such that
$$\| \tilde{\theta} - \theta_0 \| \leq \| \hat{\theta} - \theta_0 \|$$

$$\Rightarrow \lambda(\underline{X}) = \sqrt{n}(\theta_0 - \hat{\theta})^t \left\{ -\frac{1}{n}\ell''(\underline{X}; \tilde{\theta}) \right\} \sqrt{n}(\theta_0 - \hat{\theta})^t,$$

$$\downarrow \qquad\qquad \downarrow \qquad \text{FISHER MATRIX}$$

$$\mathcal{N}(0, I_d^{-1}(\theta_0)) \qquad I_d(\theta_0) = (I_{ij}(\theta_0))_{k \times k}$$

where
$$I_{jk} = \mathbb{E}\left\{ \frac{\partial \log f(\underline{X}; \theta)}{\partial \theta_i} \frac{\partial \log f(\underline{X}; \theta)}{\partial \theta_j} \right\}$$

Multivariate version of the asymptotic properties of the MLE established in the previous chapter

and $\lambda(\underline{X}) \xrightarrow{d} \chi^2(d)$ as required ∎

---

To test ($H_0: \theta = \theta_0$), we may consider several test statistics. Consider for example $T(X_1, \ldots, X_n)$ and $S(X_1, \ldots, X_n)$, both constructed from the same dataset $\underline{X} = (X_1, \ldots, X_n)$. There is no guarantee that the two test agree on rejecting or not the null.

⇒ Which test should we consider?

Adopting the paradigm introduced by Neyman & Pearson, we may introduce an alternative simple hypothesis ($H_1: \theta = \theta_1$) and at a given significance level $\alpha$, keep the test that has highest power. Neyman & Pearson proved that when testing ($H_0: \theta = \theta_0$) against ($H_1: \theta = \theta_1$), the likelihood ratio test is the most powerful. This result is known as the Neyman-Pearson lemma:

## NEYMAN-PEARSON LEMMA

Let $X_1, \ldots, X_n$ iid with density $f(x; \theta)$, $\theta \in \Theta$.
We test ($H_0: \theta = \theta_0$) against ($H_1: \theta = \theta_1$).
The likelihood ratio is $\Lambda(\underline{x}) = \dfrac{f(\underline{x}; \theta_1)}{f(\underline{x}; \theta_0)}$, and consider

$$\phi(\underline{x}) := \mathbb{1}(\Lambda(\underline{x}) \geq k), \text{ for some } k \geq 0.$$

This test is the most powerful test amongst the class of tests with significance level equal to $\alpha := \mathbb{P}_{\theta_0}(\phi(\underline{X}) = 1)$.
In other words, for any other test $\phi'(\underline{X}) \in \{0, 1\}$ such that $\mathbb{P}_{\theta_0}(\phi'(\underline{X}) = 1) \leq \alpha$, we have that

$$\mathbb{P}_{\theta_1}(\phi'(\underline{X}) = 1) \leq \mathbb{P}_{\theta_1}(\phi(\underline{X}) = 1)$$

proof = let $C := \{ \underline{x} \mid \phi(\underline{x}) = 1 \} \rightarrow \phi(\underline{x}) = \mathbb{1}(\underline{x} \in C)$

$C' := \{ \underline{x} \mid \phi'(\underline{x}) = 1 \} \rightarrow \phi'(\underline{x}) = \mathbb{1}(\underline{x} \in C')$.

Critical regions for each test.

• Take $\underline{x} \in C$. Then

$\downarrow \phi'(\underline{x}) - \phi(\underline{x}) = \mathbb{1}(\underline{x} \in C') - \underbrace{\mathbb{1}(\underline{x} \in C)}_{=1} \leq 0$ ——————— (*)

$\downarrow f(\underline{x}, \theta_1) - k f(\underline{x}, \theta_0) \geq 0$ ———————————— (**)

by definition of $\phi$

Multiplying (*) & (**) together yields

$$\left( \phi'(\underline{x}) - \phi(\underline{x}) \right)\left( f(\underline{x}, \theta_1) - k f(\underline{x}, \theta_0) \right) \leq 0$$

This expression holds for $\underline{x} \in C'$ as well, since the signs in (*) & (**) are reversed to $\geq 0$ and $< 0$.

$\Rightarrow$ The expression holds $\forall \underline{x}$. Integrating it with respect to $x$ gives

$$0 \geq \int \left( \phi'(\underline{x}) - \phi(\underline{x}) \right)\left( f(\underline{x}, \theta_1) - k f(\underline{x}, \theta_0) \right) d\underline{x}$$

$$= \int \phi'(\underline{x}) f(\underline{x}, \theta_1) d\underline{x} - \int \phi(\underline{x}) f(\underline{x}, \theta_1) d\underline{x}$$

$$+ k \left( \int \phi(\underline{x}) f(\underline{x}, \theta_0) d\underline{x} - \int \phi'(\underline{x}) f(\underline{x}, \theta_0) d\underline{x} \right)$$

$$= \mathbb{P}_{\theta_1}(\phi'(\underline{X}) = 1) - \mathbb{P}_{\theta_1}(\phi(\underline{X}) = 1)$$

$$+ k \left( \underbrace{\mathbb{P}_{\theta_0}(\phi(\underline{X}) = 1)}_{\overset{\|}{\alpha}} - \underbrace{\mathbb{P}_{\theta_0}(\phi'(\underline{X}) = 1)}_{\leq \alpha} \right)$$

$$\underbrace{\phantom{xxxxxxxxxxxxxxxxxxxxxx}}_{\geq 0}$$

$$\geq \mathbb{P}_{\theta_1}(\phi'(\underline{X}) = 1) - \mathbb{P}_{\theta_1}(\phi(\underline{X}) = 1), \text{ as required.} \blacksquare$$

---

# II. TWO-SAMPLE TESTS

## II.1. Comparing means.

Let $X_1, \dots, X_{n_1}$ iid $\mathcal{N}(\mu_1, \sigma_1^2)$
$Y_1, \dots, Y_{n_2}$ iid $\mathcal{N}(\mu_2, \sigma_2^2)$  $\delta := \mu_1 - \mu_2$

We test for equality of the means: ($H_0 : \delta = 0$).

• $\sigma_1, \sigma_2$ known.

Then $\bar{X} \sim \mathcal{N}\left(\mu_1, \dfrac{\sigma_1^2}{n_1}\right)$ and $\bar{Y} \sim \mathcal{N}\left(\mu_2, \dfrac{\sigma_2^2}{n_2}\right)$
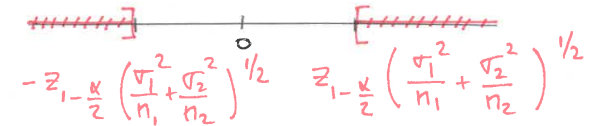
$\Rightarrow \bar{X} - \bar{Y} \sim \mathcal{N}\left(\mu_1 - \mu_2, \dfrac{\sigma_1^2}{n_1} + \dfrac{\sigma_2^2}{n_2}\right)$,

So that $\dfrac{\bar{X} - \bar{Y} - (\mu_1 - \mu_2)}{\sqrt{\dfrac{\sigma_1^2}{n_1} + \dfrac{\sigma_2^2}{n_2}}} \sim \mathcal{N}(0,1)$
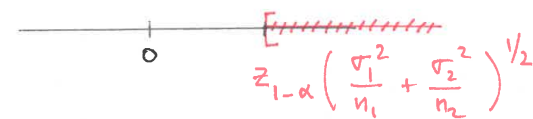
Under $H_0$, $\mu_1 = \mu_2$, so $T(\underline{X}, \underline{Y}) = \dfrac{\bar{X} - \bar{Y}}{\sqrt{\dfrac{\sigma_1^2}{n_1} + \dfrac{\sigma_2^2}{n_2}}} \sim \mathcal{N}(0,1)$

Our test statistic

$\hookrightarrow$ Critical region is



$-z_{1-\frac{\alpha}{2}}\left(\dfrac{\sigma_1^2}{n_1} + \dfrac{\sigma_2^2}{n_2}\right)^{1/2}$   $z_{1-\frac{\alpha}{2}}\left(\dfrac{\sigma_1^2}{n_1} + \dfrac{\sigma_2^2}{n_2}\right)^{1/2}$

For a one-sided test $H_0 : \delta \leq 0$, we get



$z_{1-\alpha}\left(\dfrac{\sigma_1^2}{n_1} + \dfrac{\sigma_2^2}{n_2}\right)^{1/2}$

$\llcorner$ **Power** of the one-sided test; under $(H_1 : \delta = \delta_0 > 0)$

$$1 - \beta = \mathbb{P}_{H_1}\left( W \geq z_{1-\alpha}\left( \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2} \right)^{1/2} \right),$$

where $W \sim \mathcal{N}\left( \delta_0, \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2} \right)$.

$$\Leftrightarrow \beta = \mathbb{P}_{H_1}\left( Z \leq -\frac{\delta_0}{\left( \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2} \right)^{1/2}} + z_{1-\alpha} \right),$$

where $Z \sim \mathcal{N}(0,1)$.

Assuming $n_1 = n_2 = n$ ( for sample size calculations, we are looking for the minimum value of $n$ to achieve a given power $1-\beta$ ).

Then $z_\beta = -\frac{n^{1/2} \delta_0}{\sqrt{\sigma_1^2 + \sigma_2^2}} + z_{1-\alpha}$

Taking $\sigma_1^2 = \sigma_2^2 = \sigma^2 \Rightarrow$
$$\boxed{ n \geq 2 (z_{1-\alpha} - z_\beta)^2 \left( \frac{\sigma}{\delta_0} \right)^2 } \qquad \text{(\textcolor{red}{*})}$$

As we expect,
the larger $\sigma$, or the smaller $\delta_0$,
and the larger $n$ should be.

A function of $\alpha, \beta, \sigma$ and $\delta_0$.

• $\sigma_1, \sigma_2$ unknown, assumed equal : $\sigma_1 = \sigma_2 = \sigma$.

Then $\bar{X} - \bar{Y} \sim \mathcal{N}\left( \mu_1 - \mu_2, \sigma^2 \left( \frac{1}{n_1} + \frac{1}{n_2} \right) \right)$

Moreover, $\frac{(n_1 - 1) S_1^2}{\sigma^2} \sim \chi^2(n_1 - 1)$ and $\frac{(n_2 - 1) S_2^2}{\sigma^2} \sim \chi^2(n_2 - 1)$

independent where $S_i^2 = \frac{1}{n_i - 1} \sum_{i=1}^{n} (X_i - \bar{X})^2$

$\Rightarrow \frac{(n_1 - 1) S_1^2}{\sigma^2} + \frac{(n_2 - 1) S_2^2}{\sigma^2} \sim \chi^2(n_1 + n_2 - 2)$

---

Thus:

our test statistic

$$\frac{\dfrac{\bar{X} - \bar{Y} - (\mu_1 - \mu_2)}{\sigma \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}}{\sqrt{\dfrac{(n_1-1) S_1^2 + (n_2-1) S_2^2}{\sigma^2} \Big/ (n_1 + n_2 - 2)}} = \frac{\bar{X} - \bar{Y}}{S_p^2 \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim t(n_1 + n_2 - 2)$$

$\sim \mathcal{N}(0,1)$

indpt

$\sim \chi^2(n_1 + n_2 - 2)$

Under $H_0$

where $S_p^2 = \frac{(n_1 - 1) S_1^2 + (n_2 - 1) S_2^2}{n_1 + n_2 - 2}$ is known as the "pooled"

estimator of $\sigma^2 \rightarrow$ usually, use \textcolor{red}{(*)} page 21 for sample size calculations, plugging in $S_p^2$ for $\sigma^2$.

• $\sigma_1, \sigma_2$ unknown.

If we do not assume $\sigma_1$ and $\sigma_2$ equal, we cannot pool the estimates, and one need to consider estimates of $\sigma_1$ and $\sigma_2$ separately. We have:

$$\bar{X} - \bar{Y} \sim \mathcal{N}\left( \mu_1 - \mu_2, \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2} \right)$$

Inspired by the fact that $\frac{(n_i - 1) S_i^2}{\sigma_i^2} \sim \chi^2(n_i - 1)$,

we are looking for a value of $v$ such that the quantity

$$U := \frac{v \left( \frac{S_1^2}{n_1} + \frac{S_2^2}{n_2} \right)}{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \quad \text{is approximately } \chi^2(v)$$

Once $v$ is computed, this will allow us to consider the statistic $T(X_1, \ldots, X_n, Y_1, \ldots, Y_n)$, defined by

$$T(\underline{X},\underline{Y}) = \frac{\frac{\overline{X}-\overline{Y}-(\mu_1-\mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1}+\frac{\sigma_2^2}{n_2}}}}{\sqrt{\frac{\left(\frac{S_1^2}{n_1}+\frac{S_2^2}{n_2}\right)\nu}{\frac{\sigma_1^2}{n_1}+\frac{\sigma_2^2}{n_2}} \Big/ \nu}} = \frac{\overline{X}-\overline{Y}}{\sqrt{\frac{S_1^2}{n_1}+\frac{S_2^2}{n_2}}} \approx t(\nu),$$

$\uparrow$ Under $H_0$

whose distribution can be approximated under $H_0$.

To find $\nu$, we use the method of moments. Let $V \sim \chi^2(\nu)$.

Then $\mathbb{E}V = \nu$, and $\operatorname{Var} V = 2\nu$.

Since $\mathbb{E}S_i^2 = \sigma_i^2$ (unbiased estimator of $\sigma_i$), we see that the first moments of $U$ and $V$ coincide.

Next, $\dfrac{(n_i-1)^2 \operatorname{Var} S_i^2}{\sigma_i^4} = 2(n_i-1) = $ Variance of a $\chi^2(n_i-1)$ RV

Thus, $\operatorname{Var}(S_i^2) = \dfrac{2\sigma_i^4}{n_i-1}$, and

$$\operatorname{Var} U = \frac{\nu^2}{\left(\frac{\sigma_1^2}{n_1}+\frac{\sigma_2^2}{n_2}\right)^2}\left(\frac{\operatorname{Var} S_1^2}{n_1^2}+\frac{\operatorname{Var} S_2^2}{n_2^2}\right)$$

$$= \frac{\nu^2}{\left(\frac{\sigma_1^2}{n_1}+\frac{\sigma_2^2}{n_2}\right)^2}\left(\frac{2\sigma_1^4}{n_1^2(n_1-1)}+\frac{2\sigma_2^4}{n_2^2(n_2-1)}\right) = \operatorname{Var} V = 2\nu$$

Solving for $\nu$, we obtain

$$\nu = \frac{\left(\frac{\sigma_1^2}{n_1}+\frac{\sigma_2^2}{n_2}\right)^2}{\frac{\sigma_1^4}{n_1^2(n_1-1)}+\frac{\sigma_2^4}{n_2^2(n_2-1)}}$$

$\longrightarrow$ Plug-in $S_1^2$ and $S_2^2$ in practice, for an estimate of $\nu$.

$\leftarrow$ the so-called WELCH-SATTERTHWAITE FORMULA

## II.2. Comparing two proportions.

Let $X_1, \ldots, X_{n_1} \sim B(p_1)$
$\quad\quad Y_1, \ldots, Y_{n_2} \sim B(p_2)$ $\quad\quad \delta = p_1 - p_2$

We test for equality of the two proportions $(H_0 : p_1 = p_2)$.

Put $\hat{p}_1 = \dfrac{1}{n_1}\sum_{i=1}^{n_1} X_i$, and $\hat{p}_2 = \dfrac{1}{n_2}\sum_{i=1}^{n_2} Y_i$.

The CLT ensures that $\hat{p}_1 - \hat{p}_2 \approx \mathcal{N}\left(p_1-p_2, \dfrac{p_1(1-p_1)}{n_1}+\dfrac{p_2(1-p_2)}{n_2}\right)$.

Under $H_0$, $p_1 = p_2 = p$, (equivalently $p = \dfrac{p_1+p_2}{2}$), and

$$\frac{\hat{p}_1 - \hat{p}_2}{\sqrt{p(1-p)\left(\frac{1}{n_1}+\frac{1}{n_2}\right)}} \approx \mathcal{N}(0,1).$$

$\nwarrow$ Replace $p$ by a pooled estimate $\hat{p} = \dfrac{\sum(X_i+Y_i)}{n_1+n_2}$, since $X_i$ and $Y_i$ are iid $B(p)$

Take $T(\underline{X},\underline{Y}) := \dfrac{\hat{p}_1 - \hat{p}_2}{\sqrt{\hat{p}(1-\hat{p})\left(\frac{1}{n_1}+\frac{1}{n_2}\right)}}$ as your test statistic,

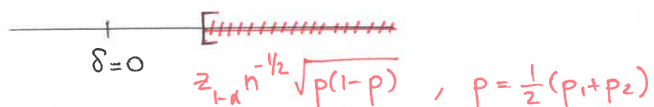and compare its value with the quantiles of the standard normal distribution.

- Remark: Alternatively, you may consider estimates of $p_1$ and $p_2$ separately, and use $\dfrac{\hat{p}_1 - \hat{p}_2}{\sqrt{\dfrac{\hat{p}_1(1-\hat{p}_1)}{n_1}+\dfrac{\hat{p}_2(1-\hat{p}_2)}{n_2}}}$.

## • Power & sample size calculations

Assuming equal sample sizes $n_1 = n_2 = n$, the one-sided rejection region is



$$z_{1-\alpha}\, n^{-1/2}\sqrt{p(1-p)} \quad,\quad p = \frac{1}{2}(p_1 + p_2)$$

Under $(H_1 : \delta = \delta_0 > 0)$, the power of the test is

$$1-\beta = \mathbb{P}_{H_1}\left( W \geq z_{1-\alpha}\, n^{-1/2}\sqrt{p(1-p)} \right),$$

where $W \approx \mathcal{N}\left( \delta_0, \dfrac{p_1(1-p_1) + p_2(1-p_2)}{n} \right)$, $p_1 - p_2 = \delta_0 > 0$

$$\Leftrightarrow \quad \beta = \mathbb{P}_{H_1}\left( Z \leq \frac{-\delta_0 n^{1/2}}{\sqrt{p_1(1-p_1) + p_2(1-p_2)}} + z_{1-\alpha}\sqrt{\frac{p(1-p)}{p_1(1-p_1) + p_2(1-p_2)}} \right)$$

where $Z \approx \mathcal{N}(0,1)$

Given $\beta$, the minimum value of $n$ achieving the required power must satisfy

$$z_\beta = \frac{-\delta_0 n^{1/2}}{\sqrt{p_1(1-p_1) + p_2(1-p_2)}} + z_{1-\alpha}\sqrt{\frac{p(1-p)}{p_1(1-p_1) + p_2(1-p_2)}},$$

— so that —

$$n \geq \frac{1}{\delta_0^2}\left( z_{1-\alpha}\sqrt{p(1-p)} - z_\beta\sqrt{p_1(1-p_1) + p_2(1-p_2)} \right)^2$$

where $p_1 - p_2 = \delta_0 > 0$
& $p = \frac{1}{2}(p_1 + p_2)$

In practice, plug in the sample estimates $\hat{p}_1, \hat{p}_2$.

---

## III - MULTIPLE TESTING

Suppose we want to test multiple null hypothesis $H_{0,1}, \ldots, H_{0,d}$.
For each $H_{0,i}$ $(i=1,\ldots,d)$, we have a p-value $p_i$ associated with the test statistic $T_i$. It is defined as

$$p_i := \mathbb{P}_{0,i}(T_i \geq t_i), \text{ where } t_i \text{ is the observed value.}$$

↑ Probability computed under the null $H_{0,i}$.

$$= 1 - F_i(t_i), \text{ denoting } F_i := \text{distribution of } T_i \text{ under } H_{0,i}.$$

Then, under $H_{0,i}$, the random variable $P_i := 1 - F_i(T_i)$ has a $\mathcal{U}(0,1)$ distribution. Indeed,

$$\mathbb{P}_{0,i}(P_i \leq p) = \mathbb{P}_{0,i}(1 - F_i(T_i) \leq p)$$
$$= \mathbb{P}_{0,i}(F_i(T_i) \geq 1-p)$$
$$= \mathbb{P}_{0,i}(T_i \geq F_i^-(1-p))$$

↑ the generalized inverse of $F_i$

$$= 1 - F_i(F_i^-(1-p))$$
$$= 1 - (1-p)$$
$$= p.$$

Now, assume that we decide to reject $H_{0,i}$ if $p_i < \alpha$; so that the probability of a type I error for a single test is $\alpha$. For $d$ tests, the probability of a type I error associated with the global test $H_0 : \bigcap_{i=1,\ldots,d} H_{0,i}$ is $1 - (1-\alpha)^d$

↑ Since we reject $H_0$ as soon as one of the $H_{0,i}$ is rejected.

For $\alpha = 0.05$ and $d = 0$, we get $1 - (1-\alpha)^d \approx 0.40$, quite far from the significance level $0.05 \Rightarrow$ Need for correction.

## III.1 Bonferroni Correction.

For $H_{0,1}$ , ... , $H_{0,d}$ = a set of $d$ null hypothesis.

Test for $H_0 := \bigcap_{i=1,-,d} H_{0,i}$

Reject each individual test at significance level $\frac{\alpha}{d}$.

$\Rightarrow$ Bonferroni rejects $H_0$ if $\min_{1 \leq i \leq d} P_i < \frac{\alpha}{d}$, where

$P_i$ = p-value associated with $H_{0,i}$.

$\hookrightarrow$ PROCEDURE

The overall significance level is

$$\mathbb{P}\left( \min_{1 \leq i \leq d} P_i < \frac{\alpha}{d} \right) = \mathbb{P}\left( \bigcup_{1 \leq i \leq d} P_i < \frac{\alpha}{d} \right)$$

Under $H_0$.    $\}$ sub-additivité

$$\leq \sum_{i=1}^{d} \mathbb{P}\left( P_i < \frac{\alpha}{d} \right) \quad\quad (*)$$

$\}$ Since $P_i \sim \mathcal{U}(0,1)$

$$= \sum_{i=1}^{d} \frac{\alpha}{d} = \alpha$$

The bound may be crude, but if the p-values are independent, then $\mathbb{P}\left( \min P_i < \frac{\alpha}{d} \right) = 1 - \mathbb{P}\left( \min P_i \geq \frac{\alpha}{d} \right)$

$$= 1 - \prod_{i=1}^{d} \mathbb{P}\left( P_i \geq \frac{\alpha}{d} \right)$$

$$= 1 - \left( 1 - \frac{\alpha}{d} \right)^d \quad \} \text{ if } d \text{ is large}$$

$$\approx 1 - e^{-\alpha} \quad\quad \} \text{ if } \alpha \text{ is small}$$

$$\approx \alpha$$

$\Rightarrow$ For independent tests, the bound $(*)$ is reasonable.
However, Bonferroni correction is known to be too conservative in many practical applications, due to the conservative individual thresholds $\frac{\alpha}{d}$, which can be smaller than needed.

---

## III.2. Fisher's combination test.

Bonferroni is looking at the smallest of the p-values to test for $H_0$. Alternatively, Fisher suggests to aggregate all the p-values, and considers the quantity

$$X^2 := -2 \sum_{i=1}^{d} \log P_i$$

$\uparrow$ when the $P_i$ are small, $X^2$ is large.
$\Rightarrow$ Reject the null $H_0 := \bigcap_{i=1,-,d} H_{0,i}$ if $X^2$ is large. How large?

Under the assumption that $P_1, ..., P_d$ are independent, $X^2 \sim \chi^2(2d)$. The distribution of $X^2$ may be used to decide on a threshold.

indeed, if $U \sim \mathcal{U}(0,1)$, then $-2 \log U \sim \chi^2(2)$, since
$$\mathbb{P}(-2 \log U \leq x) = \mathbb{P}\left( U \geq \exp\left(-\frac{x}{2}\right) \right) = 1 - e^{-x/2},$$
with density $\frac{1}{2} e^{-\frac{x}{2}}$.
& the sum of $d$ independent $\chi^2(2)$ RVs is $\chi^2(2d)$.

### III.3. Controlling the False Discovery Rate (FDR).

Instead of controlling the intersection $\bigcap H_{0,i}$, we may look at the test separately:

| | Not Reject | Reject |
|---|---|---|
| $H_{0,i}$ true | | FP |
| $H_{0,i}$ false | | TP |

FP $\leftarrow$ total number of false rejections
TP $\leftarrow$ total # of true rejections.

The <u>F</u>amily-<u>W</u>ise <u>E</u>rror <u>R</u>ate (FWER) is the probability of making at least one false rejection:

$\text{FWER} = \mathbb{P}(\text{FP} \geqslant 1)$. We say that the FWER is controlled at level $\alpha$ if $\text{FWER} \leqslant \alpha$.

$\searrow$ Bonferroni's correction controls the FWER at level $\alpha$ since
$$\mathbb{P}(\text{FP} \geqslant 1) = \mathbb{P}\left( \bigcup_{i \in \mathcal{I}_0} \text{rejecting } H_{0,i} \right)$$

$\mathcal{I}_0$ = set of indices in $\{1,..,d\}$ corresponding to the true nulls $H_{0,i}$

Let $d_0 = |\mathcal{I}_0|$
$\quad$ = # elements in $\mathcal{I}_0$

$$\leqslant \sum_{i \in \mathcal{I}_0} \mathbb{P}(\text{rejecting } H_{0,i})$$
$$= \alpha \frac{d_0}{d} \leqslant \alpha.$$

Instead of controlling the FWER, we may wish to control the mean proportion $\left( \dfrac{FP}{FP+TP} \right)$ of false positives, referred to as the <u>F</u>alse <u>D</u>iscovery <u>R</u>ate (FDR) in the literature:

$$\boxed{\text{FDR} = \mathbb{E}\left\{ \frac{FP}{FP+TP} \, \mathbb{1}(FP+TP \geqslant 1) \right\}.}$$

<u>Goal</u>: Design a procedure to control the FDR; and make sure it remains below a certain level (say $\alpha$).

The first method we discuss is due to Benjamini & Hochberg (1995) <u>Controlling the False Discovery Rate: A Practical & Powerful Approach to Multiple Testing</u>. JRSS B. Vol 57, No 1, p. 289-300.
Also, see Giraud (2015), <u>Introduction to high dim statistics</u>.

---

Consider $H_{0,1}, ..., H_{0,d} = d$-null hypotheses, with p-values $P_1, ..., P_d$.

Let $P_{(1)}, ..., P_{(d)}$ be the ordered p-values, and $H_{0,(i)}$ be the null hypothesis associated with $P_{(i)}$.

Let $k$ be the largest $i$ for which $P_{(i)} \leqslant \dfrac{i\alpha}{d}$.

Then, reject all $H_{(i)}$ for $i = 1, -, k$.

<u>PROCEDURE</u> $\qquad$ (<u>BENJAMINI & HOCHBERG</u>)
$\qquad\qquad\qquad\qquad\qquad$ (1995)

<u>Result</u>: If the $P_1, ..., P_d$ are independent, then the Benjamini & Hochberg procedure ensures that FDR $\leqslant \alpha$.

<u>proof</u>: let $\mathcal{I}_0 = \{ 1 \leqslant i \leqslant d \mid H_{0,i} \text{ is true} \}$.
$\qquad K$ = (random) number of rejected hypothesis.
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ (= FP+TP)

Then
$$\text{FDR} = \mathbb{E}\left\{ \frac{FP}{FP+TP} \, \mathbb{1}(FP+TP \geqslant 1) \right\}$$
$$= \mathbb{E}\left\{ \frac{|\{ i \in \mathcal{I}_0 : P_i \leqslant \frac{\alpha K}{d} \}|}{K} \, \mathbb{1}(K \geqslant 1) \right\}$$

Since all the rejected hypothesis have a p-value smaller than $\frac{\alpha K}{d}$

$$= \sum_{i \in \mathcal{I}_0} \mathbb{E}\left\{ \mathbb{1}\left( P_i \leqslant \frac{\alpha K}{d} \right) \frac{\mathbb{1}(K \geqslant 1)}{K} \right\}.$$
$$= \sum_{i \in \mathcal{I}_0} \mathbb{E}\, \mathbb{E}\left[ \{ \cdots \} \mid K \right],$$

where $\mathbb{E}\left[ \{ \cdots \} \mid K = k \right] = \dfrac{\mathbb{1}(k \geqslant 1)}{k} \, \mathbb{P}\left( P_i \leqslant \frac{\alpha K}{d} \mid K = k \right).$

Thus

$$\mathbb{E}\,\mathbb{E}\big[\{\cdots\}\mid K\big] = \sum_{k=1}^{d} \frac{1}{k}\,\mathbb{P}\Big(P_i \le \frac{\alpha k}{d}\mid K=k\Big)\mathbb{P}(K=k),$$

and

$$FDR = \sum_{i\in I_0}\sum_{k=1}^{d} \frac{1}{k}\,\mathbb{P}\Big(K=k,\ P_i \le \frac{\alpha k}{d}\Big)$$

$$= \sum_{i\in I_0}\sum_{k=k_i}^{d} \frac{1}{k}\,\mathbb{P}\Big(K=k\mid P_i \le \frac{\alpha k}{d}\Big)\mathbb{P}\Big(P_i \le \frac{\alpha k}{d}\Big),$$

where $k_i$ is defined to be the smallest integer $\ge 1$ such that $\mathbb{P}\big(P_i \le \frac{\alpha k_i}{d}\big) > 0$.

Now, we have $P_i \sim \mathcal{U}(0,1)$, so that $\mathbb{P}\big(P_i \le \frac{\alpha k}{d}\big) = \frac{\alpha k}{d}$, and

$$FDR = \frac{\alpha}{d}\sum_{i\in I_0}\sum_{k=k_i}^{d} \mathbb{P}\Big(K=k\mid P_i \le \frac{\alpha k}{d}\Big)$$

$$= \frac{\alpha}{d}\sum_{i\in I_0}\sum_{k=k_i}^{d}\Big\{\mathbb{P}\Big(K\le k\mid P_i \le \frac{\alpha k}{d}\Big) - \mathbb{P}\Big(K\le k-1\mid P_i \le \frac{\alpha k}{d}\Big)\Big\}$$

We claim that $\forall k \ge k_i$,

$$\mathbb{P}\Big(K\le k\mid P_i \le \frac{\alpha k}{d}\Big) \le \mathbb{P}\Big(K\le k\mid P_i \le \frac{\alpha(k+1)}{d}\Big) \quad\text{---}\quad (**)$$

Assuming that $(**)$ holds, the last written sum is a telescoping sum, and

$$FDR \le \frac{\alpha}{d}\sum_{i\in I_0}\mathbb{1}(k_i \le d)\,\mathbb{P}\Big(K\le d\mid P_i \le \frac{\alpha(d+1)}{d}\Big) \le \frac{d\alpha}{d}$$

$$\le \alpha .$$

It remains to show $(**)$ under the independence assumption of the p-values.

---

Note that $\mathbb{1}(K\le k) = \mathbb{1}\Big(\max\{i\mid P_i \le \frac{\alpha i}{d}\} \le k\Big)$

$$= \text{a function of } P_1,\ldots,P_d$$

$$=: g(P_1,\ldots,P_d).$$

The function $g : [0,1]^d \to \{0,1\}$
$$(p_1,\ldots,p_d) \mapsto g(p_1,\ldots,p_d)$$
is a nondecreasing function of $p_1,\ldots,p_d$. (why?)

To get $(*)$, we are studying the function

$$u \mapsto \mathbb{E}\{\mathbb{1}(K\le k)\mid P_i \le u\},$$

and we show that it a non decreasing function of $u \in [0,1]$, provided the $P_1,\ldots,P_d$ are independent. This function can be rewritten

$$u \mapsto \mathbb{E}\{g(P_1,\ldots,P_d)\mid P_i \le u\}.$$

More generally, a set of distributions satisfying this property $\forall$ positive & bounded $g$ is said to fulfill the Weak Positive Regression Dependency (WPRD) property. As we now show, it holds under the independence assumption, but the WPRD property hold under other assumptions as well, so the result on page 28 is more general than as stated.

We have $\qquad\qquad\qquad$ consider wlog $i=1$.

$$\mathbb{E}\{g(P_1,\ldots,P_d)\mid P_1 \le u\}$$

$$= \int_{(x_2,\ldots,x_d)\in[0,1]^{d-1}} \mathbb{E}\{g(P_1,x_2,\ldots,x_d)\mid P_1\le u\}$$
$$\times \mathbb{P}(P_2\in dx_2,\ldots,P_d\in dx_d)$$

Under the independence assumption of $P_1,\ldots,P_d$.

$\Rightarrow$ We only need to show that $\forall x_2, .., x_d$, the function

$$u \mapsto \mathbb{E}\{ g(P_1, x_2, .., x_d) \mid P_1 \leq u \}$$

is non-decreasing with $u$.

Since $g$ is nondecreasing, the function $g_1 : x_1 \mapsto g(x_1, .., x_d)$ is also nondecreasing. Thus

$$\mathbb{E}\{ g(P_1, x_2, .., x_d) \mid P_1 \leq u \} = \mathbb{E}\{ \underbrace{g_1(P_1)}_{\text{a non-neg RV}} \mid P_1 \leq u \}$$

$$= \int_0^{+\infty} \mathbb{P}( g_1(P_1) \geq x \mid P_1 \leq u) \, dx$$

$$= \int_0^{+\infty} \mathbb{P}( P_1 \geq g_1^-(x) \mid P_1 \leq u) \, dx$$

generalized inverse of $g_1$:
$$g_1^-(x) = \inf_{u \in [0,1]} \{ g_1(u) \geq x \}$$

Since $\mathbb{P}(P_1 \geq g_1^-(x) \mid P_1 \leq u) = \left( 1 - \dfrac{\mathbb{P}(P_1 < g_1^-(x))}{\mathbb{P}(P_1 \leq u)} \right)_+$

is a nondecreasing function of $u$ for all $x$, we obtain (✱). ∎

The Benjamini–Hochberg procedure is a powerful procedure, but theoretical guarantees are obtained under distributional assumptions of the p-values (such as independence). We present next an alternative procedure, introduced by Benjamini & Yekutieli (2001): The Control of the false discovery rate in Multiple Testing under Dependency. Annals of Statistics - Vol 29 - p. 1165-1188.

The idea is to replace the procedure on page 28 by "Let $k$ be the largest $i$ for which $P_{(i)} \leq \dfrac{\beta(i)\alpha}{d}$", for some appropriate function $\beta : \{1, .., d\} \to \mathbb{R}_+$, nondecreasing. The Benjamini–Hochberg procedure uses $\beta(i) = i$. The FDR is

$$FDR = \sum_{i \in I_0} \mathbb{E}\left\{ \mathbb{1}\left( P_i \leq \dfrac{\alpha \beta(K)}{d} \right) \dfrac{\mathbb{1}(K \geq 1)}{K} \right\}.$$

See the derivation on page 28.

Noting that on the event $\{K \geq 1\}$, $\left| \dfrac{1}{K} = \sum_{j \geq 1} \dfrac{\mathbb{1}(j \geq K)}{j(j+1)} \right.$,

we have

$$FDR = \sum_{i \in I_0} \sum_{j \geq 1} \dfrac{1}{j(j+1)} \mathbb{E}\left\{ \mathbb{1}\left( P_i \leq \dfrac{\alpha \beta(K)}{d} \right) \mathbb{1}(j \geq K) \mathbb{1}(K \geq 1) \right\}.$$

For $j \geq K$, we have that $\beta(K) \leq \beta(j \wedge d)$, so that

$$\mathbb{E}\left\{ \mathbb{1}\left( P_i \leq \dfrac{\alpha \beta(K)}{d} \right) \mathbb{1}(j \geq K) \mathbb{1}(K \geq 1) \right\} \leq \mathbb{P}\left( P_i \leq \dfrac{\alpha \beta(j \wedge d)}{d} \right)$$

$$= \dfrac{\alpha \beta(j \wedge d)}{d},$$

and we obtain $\left| FDR \leq \alpha \dfrac{d_0}{d} \sum_{j \geq 1} \dfrac{\beta(j \wedge d)}{j(j+1)} \right.$

We conclude that as long as $\sum_{j \geq 1} \dfrac{\beta(j \wedge d)}{j(j+1)} \leq 1$, the FDR of the procedure is less than $\alpha$. The choice $\beta(i) = i$ yields $\sum_{j \geq 1} \dfrac{\beta(j \wedge d)}{j(j+1)} = 1 + \dfrac{1}{2} + .. + \dfrac{1}{d} =: H_d > 1$, and

therefore does not guarantee a FDR $\leq \alpha$ (note (35)
that Benjamini - Hochberg ensure that FDR $\leq \alpha$ under
additional assumptions, such as independence. The bound
derived previously is crude, but does not require any
distributional assumptions on the p-values $\Rightarrow$ more
general, but more conservative).

A popular choice for $\beta$ ensuring that $\sum_{j \geq 1} \frac{\beta(j \wedge d)}{j(j+1)} \leq 1$

is $\beta(i) = \frac{i}{H_j}$, where $H_i = 1 + \frac{1}{2} + \cdots + \frac{1}{i}$ grows as $\log i$.

This procedure is known as the Benjamini - Yekutieli procedure.

Consider $H_{0,1}, \ldots, H_{0,d} = d$-null hypotheses, with

p-values $p_1, \ldots, p_d$.

Let $p_{(1)}, \ldots, p_{(d)}$ be the ordered p-values, and

$H_{0,(i)}$ the null hypothesis associated with $p_{(i)}$.

Let $k$ be the largest $i$ for which $p_{(i)} \leq \frac{i\alpha}{d H_d}$.

Then, reject all $H_{(i)}$ for $i = 1, -, k$.

PROCEDURE      (BENJAMINI & YEKUTIELI)
                        (2001)

The B & Y procedure ensures that FDR $\leq \alpha$.

Conclusion:
• B & H (p.28) more powerful, but theoretical
  guarantees depend on distributional properties of
  the p-values
• B & Y (p.33) more conservative, but more general,
  and theoretical guarantees under dependency.

---

• <u>Summary</u> = Repeated sampling inflates the type-I error rate. (36)
Bonferroni's correction controls the FWER $\mathbb{P}(FP \geq 1)$
at level $\alpha$ by adjusting the individual threshold to $\frac{\alpha}{d}$,
where $d$ = number of experiments, since

$$FWER = 1 - \left(1 - \frac{\alpha}{d}\right)^d \approx 1 - e^{-\alpha} \approx \alpha$$
  (d large)           ($\alpha$ small)

Alternatively, we may look at different quantities, such as the

$$FDR = \mathbb{E}\left(\frac{FP}{FP + TP} \mathbb{1}(FP + TP \geq 1)\right), \text{ and control using}$$

the Benjamini & Hochberg or the Benjamini & Yekutieli procedure.
(Bayesian techniques may be even more appropriate, see
  MS = BAYESIAN STATISTICS)

⁎ <u>Ex</u>: Run 100 tests with $\alpha = 0.05$
(i) 1 test in 10 is truely effective, $\beta = 0.8$
        ↳ detect 80% of them; TP = 8    } FDR $\approx 0.36$
        ↳ (100 - 10) × 0.05 = 4.5̄ = FP
(ii) 1 test in 20 is truely effective, $\beta = 0.8$ [LOW BASE RATE]
        ↳ TP = 4, FP = 4.75 $\Rightarrow$ FDR $\approx 0.54$ ↗
(iii) 1 test in 10 is truely effective, $\beta = 0.3$ [UNDERPOWERED]
        ↳ TP = 3, FP = 4.75 $\Rightarrow$ FDR $\approx 0.62$ ↗
(iv) Tests are all ineffective $\Rightarrow$ FDR = 1
    + "regression to the mean" effect: all tests flagged as
    improvements will be likely to not reproduce the results
    in a repeated experiment (cf Kahneman:
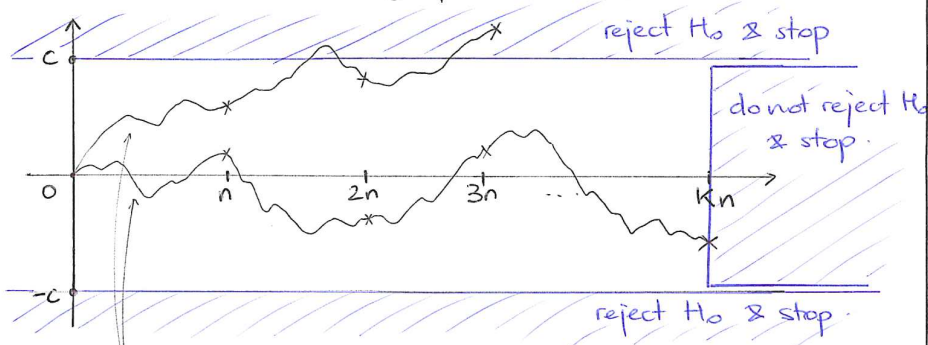    performance = a bit of talent + a lot of chance).

↳ Repeated sampling        low-base Rates
  increase the FWER    +   &               increase the FDR.
                           Underpowered tests

# IV. GROUP SEQUENTIAL TESTING

We investigate how the usual testing procedure of Neyman behaves when data are received sequentially, and when an "optimal stopping" of the data collection is performed. The procedure, also known as data peeking, can be described as follows:

(i) Collect $n$ observations $X_1, \ldots, X_n$ and compute the test statistic $T(X_1, \ldots, X_n)$.

Reject $H_0$ and stop the experiment if $|T(X_1, \ldots, X_n)| > c$. If not, do not reject and continue.

(ii) Collect another $n$ observations $X_{n+1}, \ldots, X_{2n}$ and compute $T(X_1, X_2, \ldots, X_n, X_{n+1}, \ldots, X_{2n})$.

Reject $H_0$ & stop if $|T(X_1, \ldots, X_{2n})| > c$. If not, continue.

(iii) $\ldots / \ldots$

(iv) Until a pre-specified number $K$ of maximum number of iterations is reach, or significance is achieved.



reject $H_0$ & stop

do not reject $H_0$ & stop.

reject $H_0$ & stop.

Two possible trajectories. A simple simulation study would show that repeating this many times under the case where $H_0$ is true would lead to many paths hitting the $\pm c$ boundary before reaching "do not reject $H_0$ & stop".

In other words, the type I error $\alpha$ is underlined{inflated}. We artificially detect effects when there is not, at a rate that can well exceed $\alpha$.

The inflation is know to be not as high as in the case of multiple testing, but can nevertheless easily reach $0.10 / 0.20$ with $\alpha = 0.05$.

To control for the type-I error in sequential testing, the boundaries $\pm c$ need some adjustments. We derive those in the simplest case where observations are iid and normally distribution, leading to Pocock (1977) and O'Brien-Fleming (1979) boundaries. First, we review the procedure with no peeking at the data.

**Case 1: no peeking.**

$X_1, \ldots, X_{n_0} \sim \mathcal{N}(\mu_X, \sigma^2)$ — Collect $n_0$ observations in each group, and test for $H_0 : \mu_X = \mu_Y$ vs $H_1 : \mu_X \neq \mu_Y$.

$Y_1, \ldots, Y_{n_0} \sim \mathcal{N}(\mu_Y, \sigma^2)$

Assume $\sigma^2$ known.

Then $\bar{X} = \frac{1}{n_0} \sum_{i=1}^{n_0} X_i \sim \mathcal{N}(\mu_X, \frac{\sigma^2}{n_0})$, $\bar{Y} \sim \mathcal{N}(\mu_Y, \frac{\sigma^2}{n_0})$

and $\bar{X} - \bar{Y} \sim \mathcal{N}(\mu_X - \mu_Y, \frac{2\sigma^2}{n_0})$.

Put $I_0^{-1}(1) := \frac{2\sigma^2}{n_0} =$ Fisher information.

Then $\sqrt{I_0(1)}\, (\bar{X} - \bar{Y}) \sim \mathcal{N}(\sqrt{I_0(1)}\, \delta_0, 1)$, $\delta_0 := \mu_X - \mu_Y$.

When testing only once, under $H_0$, $\sqrt{I_0(1)}\, (\bar{X} - \bar{Y}) \sim \mathcal{N}(0,1)$, and $T(1) := \sqrt{I_0(1)}\, (\bar{X} - \bar{Y})$ can be used to construct a rejection region at significance level $\alpha$: $\mathbb{P}_{H_0}(|T(1)| \geq z_{1-\frac{\alpha}{2}}) = \alpha$.

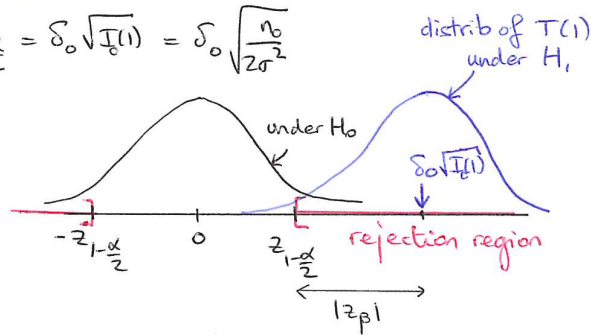Under $H_1 : \mu_X - \mu_Y = \delta_0 > 0$, the required sample size can be

calculated ensuring enough power:

$$\mathbb{P}_{H_1}\left( \sqrt{I_0(1)}\, (\bar{X} - \bar{Y}) \geq z_{1-\frac{\alpha}{2}} \right) = 1 - \beta \quad , \quad \text{where under}$$

neglecting the left region, which has little mass under $H_1$

$$H_1, \quad \sqrt{I_0(1)}\,(\bar{X} - \bar{Y}) \sim \mathcal{N}\left(\sqrt{I_0(1)}\,\delta_0, 1\right)$$

$$\Rightarrow \text{Take} \quad -z_\beta + z_{1-\frac{\alpha}{2}} = \delta_0 \sqrt{I_0(1)} = \delta_0 \sqrt{\frac{n_0}{2\sigma^2}}$$

distrib of $T(1)$ under $H_1$

$$n_0 \geq 2\sigma^2 \left( \frac{z_{1-\beta} + z_{1-\frac{\alpha}{2}}}{\delta_0} \right)^2$$

under $H_0$

$\delta_0 \sqrt{I(1)}$

rejection region

$-z_{1-\frac{\alpha}{2}}$   0   $z_{1-\frac{\alpha}{2}}$

$|z_\beta|$

• **Case 2: Optimal stopping.**

Suppose now that we allow ourselves to peek at regularly spaced intervals $t_1, .., t_K$, where $K$ is fixed in advance. For simplicity, assume that $n$ observations are collected in each group between two successive times $t_k$ and $t_{k+1}$, so that at time $t_k$, each group has collected $kn$ observations.

Reject $H_0$ at time $t_k$ & stop the experiment if $|T(k)| \geq b(k)$ where $T(k) := \sqrt{I(k)}\,(\bar{X} - \bar{Y}) \sim \mathcal{N}(0,1)$ under $H_0$,

$I^{-1}(k) = \frac{2\sigma^2}{kn}$, $k = 1, .., K$, for some carefully chosen $b(1), .., b(K)$.

$$\text{Reject } H_0 \iff \bigcup_{k=1}^{K} \left\{ |T(k)| \geq b(k) \right\}$$

$$\text{Inflated type-I error} = \mathbb{P}_{H_0}\left( \bigcup_{k=1}^{K} \left\{ |T(k)| \geq b(k) \right\} \right) > 0.05$$
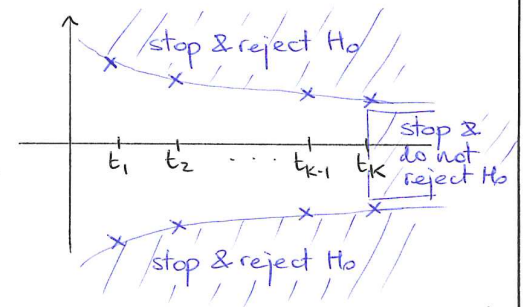
if $K \geq 2$.

The goal is to select boundaries $b(1), .., b(K)$ ensuring a

type-I error of $\alpha$, and a sample size $n$ ensuring a power of $1-\beta$ for all $\delta_0 = \mu_X - \mu_Y$ larger than some value.

$$\mathbb{P}_{H_0}\left( |T(1)| < b(1), ..., |T(K)| < b(K) \right) = 1 - \alpha$$

Reject at time $t_k$ if
$$\left\{ |T(1)| < b(1), \right.$$
$$..., \quad |T(k-1)| < b(k-1),$$
$$\left. |T(k)| \geq b(k) \right\}$$

stop & reject $H_0$

stop & do not reject $H_0$

$t_1$  $t_2$  $\cdots$  $t_{k-1}$  $t_K$

stop & reject $H_0$

$\Rightarrow$ We need to derive the joint distribution of the vector $\begin{pmatrix} T(1) \\ \vdots \\ T(K) \end{pmatrix}$.

Recall that $T(k) = \sqrt{I(k)}\,(\bar{X} - \bar{Y})$, where $I(k) = \frac{kn}{2\sigma^2}$.

Put $W(k) := \sqrt{I(k)}\, T(k)$
$$= I(k)\,(\bar{X} - \bar{Y})$$
$$= \frac{kn}{2\sigma^2}\left( \frac{X_1 + \cdots + X_{kn}}{kn} - \frac{Y_1 + \cdots + Y_{kn}}{kn} \right)$$
$$= \frac{1}{2\sigma^2}\left( X_1 + \cdots + X_{kn} - Y_1 - \cdots - Y_{kn} \right)$$

$\left\{ W(k) \right\}_{k=1..K}$ defines an independent increment process since

$$\begin{cases} W(1) = W(1) \\ W(2) = W(1) + [W(2) - W(1)] \\ W(3) = W(1) + [W(2) - W(1)] + [W(3) - W(2)] \\ .../... \end{cases}$$

independent

$\Rightarrow \text{var } W(k) = I(k)$

and $\text{cov}\left(W(k), W(\ell)\right) = \text{cov}\left(W(k), \left[W(\ell) - W(k)\right] + W(k)\right)$    (41)

$(k < \ell)$

$$= \underbrace{\text{cov}\left(W(k), W(\ell) - W(k)\right)}_{=0} + \text{var } W(k)$$

$$= I(k).$$

Thus, since $T(k) = I^{-1/2}(k) W(k)$, $\text{var } T(k) = 1$, and

$$\text{cov}\left(T(k), T(\ell)\right) = \text{cov}\left(I^{-1/2}(k) W(k), I^{-1/2}(\ell) W(\ell)\right)$$

$(k < \ell)$

$$= I^{-1/2}(k) I^{-1/2}(\ell) \underbrace{\text{cov}\left(W(k), W(\ell)\right)}_{= I(k)}$$

$$= \sqrt{\frac{I(k)}{I(\ell)}} = \sqrt{\frac{k}{\ell}}$$

↖ holds as well for $k = \ell$

* Summary: Under $H_0$, the vector $(T(1), \ldots, T(K))^t$ is multivariate normal with mean vector $0$ and covariance matrix $\Sigma = (\Sigma_{k\ell})$, where $\Sigma_{k\ell} = \sqrt{\frac{k}{\ell}}$, $k \leq \ell$.

The conditional distributions of a multivariate normal vector with known mean and covariance matrix can be computed using a recursive numerical integration algorithm, see e.g. Armitage, McPherson & Rowe (1969), which can be used to choose $b(1), \ldots, b(K)$.

↳ There are $\infty$-many choices of the $b(1), \ldots, b(K)$ that fulfil the requirement

$$\mathbb{P}_{H_0}\left(|T(1)| < b(1), \ldots, |T(k)| < b(k)\right) = 1 - \alpha.$$

Following Wang & Tsiatis (1987), take the parametric form $b(k) = c \, k^{\gamma - 1/2}$, where $\gamma \in [0, 1/2]$ is referred to as the shape parameter.

---

$\gamma = 1/2$ leads to Pocock (1977) boundaries.

$\gamma = 0$ leads to O'Brien - Fleming (1979) boundaries.

We are looking for a value of $c$ such that

$$\mathbb{P}_{H_0}\left(\bigcap_{k=1}^{K}\left\{|T(k)| < c \, k^{\gamma - \frac{1}{2}}\right\}\right) = 1 - \alpha$$

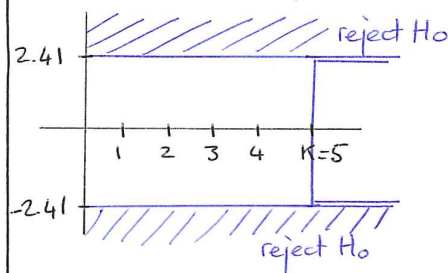↳ denote the solution $c(\alpha, K, \gamma)$.

* Ex: take $K = 5$, $\alpha = 0.05$.

Then Pocock: $c(\alpha, K, \frac{1}{2}) = 2.4135$
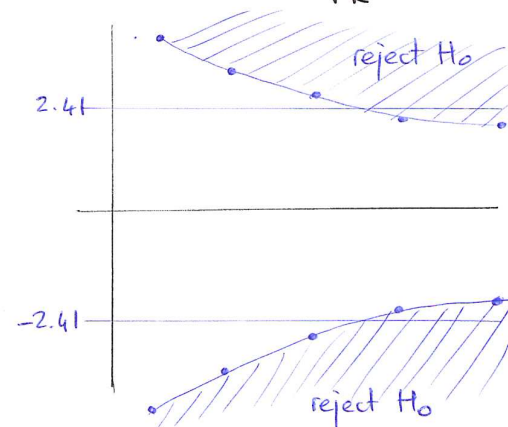
O'Brien-Fleming: $c(\alpha, K, 0) = 4.5618$

Pocock leads to constant boundaries:

Reject at time $t_k$ if $|T(k)| \geq 2.41$.



O'Brien-Fleming rejects $H_0$ at time $t_k$ if

$$|T(k)| \geq \frac{4.5618}{\sqrt{k}}.$$

We see that the Pocock procedure is likely to stop earlier

x <u>Remark</u> = generalizations.

The key step in the derivation above is the independence of the increments of the process $\{W(k)\}$. In the case where $\sigma^2$ is unknown and must be estimated from the data, this property does not necessarily hold. Another problematic example is when testing for the equality of two binomial proportions $p_X = p_Y$, where the variance $\sigma^2$ of the difference of the sample means depends on the unknown parameters themselves. However, although the process does not have independent increments, independence holds asymptotically, as noted in Scharfstein, Tsiatis & Robins (1997).

"Any efficient based test or estimator for [ the difference $\mu_X - \mu_Y$ ], when computed sequentially over time, has, asymptotically, a normal independent increment process whose distribution depends only on the [ difference $\mu_X - \mu_Y$ ] and the statistical information."

$\uparrow$ $I(1), .., I(K)$.

$\uparrow$ asymptotic : as $n \to \infty$ in between to times $t_k$ and $t_{k+1}$.

— The covariance $\text{cov}(T(k), T(\ell)) = \sqrt{\dfrac{I(k)}{I(\ell)}}$ — asymptotic

• <u>Power & sample size.</u>

Under $H_1$ $T(k) = \sqrt{I(k)} \, (\bar{X} - \bar{Y}) \sim \mathcal{N}(\sqrt{I(k)} \, \delta_0 , 1)$,

where $I(k) = \dfrac{kn}{2\sigma^2}$. The power under $H_1 : \mu_X - \mu_Y = \delta_0 > 0$ is

$$1 - \mathbb{P}_{H_1}\Big( \underbrace{|T(1)| < b(1), .., |T(K)| < b(K)}_{\text{fail to reject } H_0} \Big)$$

Information at time $t_K$ is $I(K) = \dfrac{Kn}{2\sigma^2}$.

Put $I(k) = \dfrac{kn}{2\sigma^2} = \dfrac{k}{K} I(K)$.

$\Rightarrow$ Under $H_1$, $T(k) \sim \mathcal{N}\Big( \delta_0 \sqrt{\dfrac{k}{K} I(K)} , 1 \Big)$,

so that $(T(1), .., T(K))^t$ is multivariate normal, with mean vector $\Big( \delta_0 \sqrt{\dfrac{1}{K} I(K)} , \delta_0 \sqrt{\dfrac{2}{K} I(K)} , .., \delta_0 \sqrt{I(K)} \Big)^t$,

and covariance matrix $\Sigma = (\Sigma_{k\ell})$ ; $\Sigma_{k\ell} = \sqrt{\dfrac{k}{\ell}}$ , $k \leq \ell$.

Power is

$$1 - \beta = 1 - \mathbb{P}\Big( \bigcap_{k=1}^{K} \{ |T(k)| < c(\alpha, K, \gamma) \, k^{\gamma - \frac{1}{2}} \} \Big),$$

where $\mathbb{P}$ = distribution of $(T(1), .., T(K))^t$ given above.

For fixed $\alpha$, $K$ and $\gamma$, the power is an increasing function of $\delta_0$. It can be computed numerically using recursive integration.

<u>sample size calculations.</u>

The mean vector can be rewritten

$$\delta_0 \sqrt{I(K)} \Big( \sqrt{\dfrac{1}{K}} , \sqrt{\dfrac{2}{K}} , .., 1 \Big)^t ,$$

where $I(K) = \dfrac{Kn}{2\sigma^2}$ depends on the sample size $n$ between two times $t_k$ and $t_{k+1}$.

$\searrow$ The power is an increasing function of $\delta_0 \sqrt{I(K)}$, and we can solve for $\delta_0 \sqrt{I(K)}$ that gives power $1 - \beta$. Denote this solution $\delta(\alpha, \beta, K, \gamma)$.

$\delta_0 \sqrt{I(K)}$ plays the role of the non-centrality parameter

as in the case of no peeking (page 39, it is given by $\delta_0 \sqrt{I_0(1)}$ )

no peeking solves for

$$\delta_0 \sqrt{I_0(1)} = z_{1-\beta} + z_{1-\frac{\alpha}{2}}$$

sequential testing solves for

$$\delta_0 \sqrt{I(K)} = \delta(\alpha, \beta, K, \gamma)$$

Take $K=1$, $\gamma = 1/2$, and $\delta(\alpha, \beta, 1, 1/2) = z_{1-\beta} + z_{1-\frac{\alpha}{2}}$

$$\delta_0^2 \boxed{\frac{n_0}{2\sigma^2}} = (z_{1-\beta} + z_{1-\frac{\alpha}{2}})^2$$

$$\underset{\parallel}{I_0(1)} = \frac{(z_{1-\beta} + z_{1-\frac{\alpha}{2}})^2}{\delta_0^2}$$

$$\delta_0^2 \boxed{\frac{nK}{2\sigma^2}} = \delta^2(\alpha, \beta, K, \gamma)$$

$$\underset{\parallel}{I(K)} = \frac{\delta^2(\alpha, \beta, K, \gamma)}{\delta_0^2}$$

$$\Rightarrow \quad I(K) = I_0(1) \times \left( \frac{\delta(\alpha, \beta, K, \gamma)}{z_{1-\beta} + z_{1-\alpha/2}} \right)^2$$

Information after K successive peaks at the data = Information required if testing only once × Inflation Factor (IF)

"The relative increase of information necessary for a group-sequential test to have the same power as a single fixed sample test. It depends on $\alpha$, $\beta$, and on the group sequential design parameters $K$, $\gamma$.

$nK$ = total number of obs in each group.
To derive this value (ensuring power $1-\beta$), multiply $n_0$ from a single analysis by the IF to get the total sample size $nK$.

Ex: For $\alpha = 0.05$, $1-\beta = 0.8$, $K = 5$,
      IF = 1.23 (Pocock) and IF = 1.03 (O'Brien)

For $\alpha = 0.05$, $1-\beta = 0.9$, $K = 5$,
      IF = 1.21 (Pocock) and IF = 1.03 (O'Brien)

In practice, compute the number $n_0$ of observations required to achieve some power using a usual fixed sample design, and multiply this number by IF to have the same power in a group-sequential test. Interim analyses would be conducted after every $\frac{n_0}{K}$ IF observations in each group.

Which of Pocock or O'Brien should be used?

IF indicate that Pocock require larger sample sizes ($\equiv$ larger information) than O'Brien to do the sequential test, at the same significance level and power.
But Pocock tests have a better chance to stop early because of the shape of the boundary.

Let's formalize this.

We compute the average information required to stop a sequential test when $H_1$ is true, and the true effect is $\delta_0 > 0$. To ensure a given power of $1-\beta$, recall that

$$I(K) = \frac{\delta^2(\alpha, \beta, K, \gamma)}{\delta_0^2}, \quad \leftarrow I(K) = I_0(1) \times IF$$

while the fixed sample design is such that $I_0(1) = \frac{(z_{1-\beta} + z_{1-\frac{\alpha}{2}})^2}{\delta_0^2}$.

Let $V$ = # interim analyses before a study is stopped.

$\Rightarrow$ The average information before a study is stopped is then

$$\frac{I(K)}{K} \mathbb{E} V = I_0(1) \frac{IF}{K} \mathbb{E} V \quad \text{computed numerically under } H_1 : \delta_0 > 0$$

For $K=5$, $\alpha = 0.05$ and $1-\beta = 0.9$, we obtain

$\quad$ EV $= 2.83$ (Pocock)

$\quad$ EV $= 3.65$ (O'Brien),

and the average information is $\quad$ $0.68 \, I_0(1)$ $\quad$ (Pocock)

$\qquad\qquad\qquad\qquad\qquad\qquad$ $0.75 \, I_0(1)$ $\quad$ (O'Brien)

Compare with the maximum info $\quad$ $1.21 \, I_0(1)$ $\quad$ (Pocock)

$\qquad\qquad\qquad\qquad\qquad\quad$ $1.03 \, I_0(1)$ $\quad$ (O'Brien)

$\Downarrow$

$\searrow$ If we want a design which stops on average with less information when there truly is a treatment difference, then Pocock is preferable over O'Brien.

$\searrow$ However, a design with better stopping properties under $H_1$ need greater maximum information & thus more observations will be needed if $H_0$ is true.

[REF] B. Tsiatis Lecture Notes
$\qquad\quad$ ST 520, Statistical Principles of Clinical Trials (2017).