

PRE1: APPLIED STATISTICS

Hypothesis Testing

Marie-Anne Poursat

AI Master
Université Paris-Saclay

October 15, 2020

Statistical testing uses data to decide whether a statement (= *null hypothesis*) is true.

- **EX1** : *European regulations require that any presence of genetically modified organisms (GMOs) in food be labelled as soon as the level of GMOs exceeds the 0.9% threshold.*

↔ Is the level of GMOs τ below the 0.9% threshold regulations ?

Data : 20 packets of cereals, $\hat{\tau} = 0.91$ is the observed percentage of OGMs

$$H_0 : \tau \leq 0.9 \text{ versus } H_1 : \tau > 0.9$$

Is the observed difference a **real difference** likely to appear in the larger population ? Or is it observed in the sample **by chance** ?

- EX2 : Is the coin toss fair ?

Data : $X = 18$ heads out of 30 tosses

Model : $X \sim \text{Bin}(30, p)$,

$H_0 : p = 0.5$ versus $H_1 : p \neq 0.5$

- EX3 : Is a given gene *differentially expressed* between two cell types (normal and tumor cells) ?

Data : 2 Gaussian samples, $\bar{x}_1 - \bar{x}_2 = 1.25$,

$H_0 : \mu_1 - \mu_2 = 0$ versus $H_1 : \mu_1 - \mu_2 \neq 0$

The first step in the testing procedure : find H_0 and H_1

H_0 and H_1 are written according to the model parameters.

- " H_0 is **accepted**" = not rejected ! There is no evidence from the data that H_0 is wrong (but H_0 is not necessarily true)
- " H_0 is **rejected** in favor of H_1 " = H_1 explains the data significantly better than H_0 so we decide H_1
- A hypothesis can be **simple** or **composite** : a *simple* hypothesis specifies the value of the unknown parameters
 $H_0 : p = 0.5$ and $H_1 : p \neq 0.5$, H_0 is *simple* and H_1 is *composite*, *two-sided*
 $H_0 : \tau \leq 0.9$ and $H_1 : \tau > 0.9$: H_0 and H_1 are *composite*, *one-sided*

2 types of errors

The decision to accept or reject H_0 is based on data observed from a random process

↔ the decision is random and may be incorrect

2 types of errors :

- 1 to reject H_0 when it is true : Type I error
- 2 to retain H_0 when it is false : Type II error

It is not possible to ensure that the probabilities of making a Type I error and a Type II error are both arbitrarily small

↔ Classical testing paradigm : focus on the Type I error

- the probability of a type I error is kept below α , the level of the test ($\alpha = 0.05, 0,01$)
- the probability of a type II error is **not** controlled

The second step in the testing procedure : determine a test statistic T
This is the quantity calculated from the data whose numerical value leads to acceptance or rejection of H_0 .

- **EX1** : suppose that X_1, \dots, X_{20} are i.i.d $\mathcal{N}(\tau, \sigma^2)$
A reasonable choice is to reject H_0 if $(\bar{X} - 0.9) > c$.
 $T = \bar{X} - 0.9$ is a **test statistic** and c is a **critical value**
- **EX2** : $X \sim \text{Bin}(30, p)$.
Take $T = |X - 15|$ and reject H_0 if $T > c$.

In some cases there are several possible choices for T (corresponding to different statistical tests); in more complicated cases, the choice of T is not straightforward.

Rejection regions

The third step in the testing procedure : determine the rejection region

The rejection region is the set of observed values of the test statistic that lead to reject H_0 .

- EX1 : $\mathcal{R} = \{(X_1, \dots, X_{20}) : T(X_1, \dots, X_{20}) > c\}$
- EX2 : $\mathcal{R} = \{X : T(X) > c\}$

Usually, the rejection region is of the form

$$\mathcal{R} = \{(X_1, \dots, X_n) : T(X_1, \dots, X_n) > c\}$$

$\Leftrightarrow c$ The value of c is determined by

$$P(\text{reject } H_0 \text{ when } H_0 \text{ is true}) = P_{H_0}(T > c) = \alpha$$

Example 1 (GMOs'rate)

Example : one-sample t -test

EX1 : X_1, \dots, X_{20} i.i.d $\mathcal{N}(\tau, \sigma^2)$

① $H_0 : \tau \leq 0.9$ and $H_1 : \tau > 0.9$

② $T = (\bar{X} - 0.9)$

③ reject H_0 if $T > c$

$$P_{H_0}(T > c) = P\left(\frac{\bar{X} - 0.9}{S/\sqrt{20}} > \frac{c}{S/\sqrt{20}}\right) = \alpha$$

Thus, $c = t_{1-\alpha}(19)S/\sqrt{20}$ where $t_{1-\alpha}(19)$ is the quantile of the t -distribution with 19 degrees of freedom.

The test using the t -quantile is called the *one-sample t -test*.

④ observed data : $\bar{X}^{obs} = 0.91$, $S^{obs} = .06$; $t_{.95}(19) = 1.73$

Then $T^{obs} = 0.01$ and $c = .023$.

⑤ Decision : we do not reject H_0 ; according to these data, there is no evidence that the product does not respect the european regulations.

Rather than specifying α and computing c , we calculate the P-value of the test

Definition : The P-value for a sample is defined as the smallest value of α for which the null hypothesis is rejected.

\hookrightarrow To perform the test, find the p-value of the sample and then H_0 is rejected if we decide to use α larger than the p-value :

$$\text{reject } H_0 \iff \text{p-value} < \alpha$$

Interpretation :

- a small p-value is evidence *against* H_0
- a large p-value shows that the *data are consistent* with H_0
- the p-value tells us whether the decision to reject or accept H_0 is close to α

Proposition

Suppose that the rejection region is of the form $T > c$. Then,

$$\text{p-value} = P(T \geq T^{obs})$$

where T is the test statistic and T^{obs} is the observed numerical value of T on the data.

Statistical softwares calculate p-values.

Typically, the software output uses the evidence scale :

- p-value $< .001$ very strong evidence against H_0
- p-value $< .01$ strong evidence against H_0
- p-value $< .1$ weak evidence against H_0
- p-value $> .1$ little or no evidence against H_0

EX1 (GMOs' rate) : p-value

$$\text{p-value} = P(\bar{X} - 0.9 \geq T^{obs})$$

Answer : $P(T \geq 0.037) = 0.485$ where T is a $t(19)$ variable.

There is no evidence against H_0 : H_0 is not rejected (but we don't know the type II error).

Testing for a mean

Normal model : X_1, \dots, X_n i.i.d. $\mathcal{N}(\mu, \sigma^2)$

- $H_0 : \mu = \mu_0$ $H_1 : \mu > \mu_0$

$$T = \frac{\bar{X} - \mu_0}{S/\sqrt{n}} \sim t(n-1) \quad (\text{exact Student distribution})$$

\Leftrightarrow p-value = $P(T \geq T^{obs}) = 1 - \text{stats.t.cdf}(T^{obs}, df=n-1)$

- if $H_1 : \mu < \mu_0$,
p-value = $P(T \leq T^{obs}) = \text{stats.t.cdf}(T^{obs}, df=n-1)$
- if $H_1 : \mu \neq \mu_0$,
p-value = $P(|T| \geq |T^{obs}|) = 2 * (1 - \text{stats.t.cdf}(T^{obs}, df=n-1))$

Testing for a mean

Non-normal model X_1, \dots, X_n i.i.d. $\mu = E(X_i)$, $\sigma^2 = \text{Var}(X_i)$

- $H_0 : \mu = \mu_0$ $H_1 : \mu > \mu_0$

$$T = \frac{\bar{X} - \mu_0}{\hat{\sigma}/\sqrt{n}} \sim \mathcal{N}(0, 1) \quad (\text{approximate distribution via CLT})$$

\Leftrightarrow p-value = $P(T \geq T^{obs}) = 1 - \text{stats.norm.cdf}(T^{obs})$

- if $H_1 : \mu \neq \mu_0$,
p-value = $P(|T| \geq |T^{obs}|) = 2 * (1 - \text{stats.norm.cdf}(|T^{obs}|))$
- if $H_1 : \mu < \mu_0$,
p-value = $P(T \leq T^{obs}) = \text{stats.norm.cdf}(T^{obs})$

Testing for one parameter

X_1, \dots, X_n i.i.d. with distribution depending on θ , $\hat{\theta}$ MLE, $\widehat{\text{s.e.}}$ estimated standard error of $\hat{\theta}$

- $H_0 : \theta = \theta_0$ $H_1 : \theta > \theta_0$

$$T = \frac{\hat{\theta} - \theta_0}{\widehat{\text{s.e.}}} \sim \mathcal{N}(0, 1) \quad (\text{approximate distribution for the MLE})$$

\Leftrightarrow p-value = $P(T \geq T^{obs}) = 1 - \text{stats.norm.cdf}(T^{obs})$

- if $H_1 : \theta \neq \theta_0$,
p-value = $P(|T| \geq |T^{obs}|) = 2 * (1 - \text{stats.norm.cdf}(|T^{obs}|))$
- if $H_1 : \theta < \theta_0$,
p-value = $P(T \leq T^{obs}) = \text{stats.norm.cdf}(T^{obs})$

There is a relationship between the rejection regions of tests of level α and the $1 - \alpha$ confidence intervals :

The test of $H_0 : \theta = \theta_0$ versus $H_1 : \theta \neq \theta_0$ with test statistic $T = \frac{\hat{\theta} - \theta_0}{\widehat{\text{s.e.}}}$ and level α is the same test as the test that rejects H_0 if and only if $\theta_0 \notin [\hat{\theta} - q_{1-\alpha/2}\widehat{\text{s.e.}}; \hat{\theta} - q_{\alpha/2}\widehat{\text{s.e.}}]$

q_α is the α -quantile of the distribution of T .

Testing $\theta = \theta_0$ is equivalent to checking whether θ_0 is in the confidence interval.

The testing procedure

- 1 Specify the model and the hypotheses H_0 and H_1
- 2 Find an appropriate test statistic T :
 - T is calculable (does not depend on unknown parameters)
 - the distribution of T under H_0 is known
- 3 Find the form of the rejection region (*look at H_1 !*)
- 4 Calculate the p-value and make your decision
 - if you reject H_0 , the risk of an incorrect decision is less than α (*type I error*)
 - if you don't reject H_0 , the risk of an incorrect decision is usually unknown (*type II error*)

The **power of the test** is defined as the probability of rejecting H_0 when H_1 is true = 1 - type II error.

The power depends on the actual value of the parameter under H_1 and is not calculable. In practice, choosing the test (of given level) that maximizes the power is difficult.

The two-sample t -test

Comparison of two means

EX3 : Is a given gene *differentially expressed* between two cell types ?

- 2 independent normal samples X_{11}, \dots, X_{1n_1} and X_{21}, \dots, X_{2n_2} with means μ_1 and μ_2 and identical variances σ^2 .

$$H_0 : \mu_1 = \mu_2 \quad \text{versus} \quad H_1 : \mu_1 \neq \mu_2$$

- test statistic $T = \frac{\bar{X}_1 - \bar{X}_2}{S \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$, S^2 is the *pooled estimator* of σ^2

$$S^2 = \frac{\sum_i (X_{1i} - \bar{X}_1)^2 + \sum_i (X_{2i} - \bar{X}_2)^2}{n_1 + n_2 - 2}$$

The distribution of T under H_0 is the t -distribution with $n_1 + n_2 - 2$ df.

- Rejection region $\mathcal{R} = \{T > c\}$
- P-value = $P_{H_0}(|T| \geq |T^{obs}|)$

The Likelihood Ratio Test (LRT)

X_1, \dots, X_n i.i.d. with distribution depending on $\theta \in \mathbb{R}^p$, $\hat{\theta}$ MLE

$$H_0 : \theta = \theta_0 \quad H_1 : \theta \neq \theta_0$$

- The Likelihood ratio test rejects H_0 if

$$\text{LRT} = 2 \left[\log L(\hat{\theta}) - \log L(\theta_0) \right] \geq c$$

- In simple models, the exact distribution of a transformation of LRT can be found to calculate an exact p-value
- In more complex models, the approximate distribution of LRT under H_0 is the **chi-squared distribution** with p df.
- This test can be generalized to test two *nested* models and is very popular.

The degrees of freedom of the χ^2 distribution is the difference of the number of parameters under H_1 and under H_0 .

Exercise 1

Dataset *Quine* of Lab 2. Can we say that the mean number of days absent from school is greater than 15?

```
print(np.mean(df['Days']))  
16.46
```

Exercise 2

An election opposes two candidates A and B. Denote by p_A the proportion of voters who vote for candidate A in the total population of size $N = 60$ millions. In a poll, $n = 1000$ voters are questioned : 52% of these voters announce that they will vote for the candidate A. Can we conclude that candidate A will win (at the 5% level) ?