# STAT 231 - Statistics

Cameron Roopnarine

Last updated: March 31, 2020

# Contents

## 0.1 2020-03-06

---
<div align="center">2020-03-13</div>

---

<u>Roadmap:</u>

(i) Recap and the relationship between Confidence and Hypothesis

(ii) Example: Bias Testing

(iii) Testing for variance (Normal)

(iv) What if we don't know how to construct a Test-Statistic?

---

**EXAMPLE 0.1.1.** $Y_1, \ldots Y_n$ iid $N(\mu, \sigma^2)$
- $\sigma^2 = $ known
- $\mu = $ unknown
- Sample: $\{y_1, \ldots, y_n\}$
- $\overline{y} = $ sample mean
- $H_0$: $\mu = \mu_0$
- $H_1$: $\mu \neq \mu_0$

$$D = \left| \frac{\overline{Y} - \mu_0}{\frac{\sigma}{\sqrt{n}}} \right| \qquad \rightarrow \quad \text{Test-Statistic (r.v.)}$$

$$d = \left| \frac{\overline{y} - \mu_0}{\frac{\sigma}{\sqrt{n}}} \right| \qquad \rightarrow \quad \text{Value of the Test-Statistic}$$

$$p\text{-value} = P(D \geqslant d) \qquad \text{assuming } H_0 \text{ is true}$$
$$= P(|Z| \geqslant d) \qquad Z \sim N(0,1)$$

---

<u>Question:</u> Suppose the $p$-value for the test $> 0.05$ if and only if $\mu_0$ belongs in the $95\%$ confidence interval for $\mu$?

YES.

Suppose $\mu_0$ is in the $95\%$ confidence interval for $\mu$, i.e.

$$\overline{y} \pm 1.96 \frac{\sigma}{\sqrt{n}}$$

$$\mu_0 \leqslant \overline{y} + 1.96 \frac{\sigma}{\sqrt{n}}$$
$$\mu_0 \geqslant \overline{y} - 1.96 \frac{\sigma}{\sqrt{n}}$$

These two equations yield

$$d = \left| \frac{\overline{y} - \mu_0}{\frac{\sigma}{\sqrt{n}}} \right| \leqslant 1.96$$

$$P(|Z| \geqslant d) > 0.05$$

<u>General result</u> (assuming same pivot)

$p$-value of a test $H_0$: $\theta = \theta_0$ vs $H_1$: $\theta \neq \theta_0$ is more than $q\%$, then $\theta_0$ belongs to the $100(1 - q)\%$ confidence interval and vice versa.

**EXAMPLE 0.1.2** (Bias). A 10 kg weighted 20 times $(y_1, \ldots, y_n)$
- $H_0$: The scale is unbiased
- $H_1$: The scale is biased

If the scale was unbiased,

$$Y_1, \ldots, Y_n \sim N(10, \sigma^2)$$

If the scale was biased,

$$Y_1, \ldots, Y_n \sim N(10 + \delta, \sigma^2)$$

- $H_0$: $\delta = 0$ (unbiased)
- $H_1$: $\delta \neq 0$ (biased)

is equivalent to
- $H_0$: $\mu = 10$
- $H_1$: $\mu \neq 10$

Test-statistic:

$$D = \left| \frac{\overline{Y} - 10}{\frac{s}{\sqrt{n}}} \right|$$

Compute $d$.

$$d = \left| \frac{\overline{y} - 10}{\frac{s}{\sqrt{n}}} \right|$$

$$p\text{-value} = P(D \geqslant d)$$
$$= P(|T_{19}| \geqslant d)$$

**EXAMPLE 0.1.3** (Draw Conclusions). $Y_1, \ldots, Y_n = $ co-op salaries. $Y_1, \ldots, Y_n \sim N(\mu, \sigma^2)$
- $H_0$: $\mu = 3000$
- $H_1$: $\mu < 3000$ ($\mu \neq 3000$)

$$D = \left| \frac{\overline{Y} - \mu_0}{\frac{s}{\sqrt{n}}} \right|$$

$$D = \begin{cases} 0 & \overline{Y} > \mu_0 \\ \frac{\overline{Y} - \mu_0}{\frac{s}{\sqrt{n}}} & \overline{Y} < \mu_0 \end{cases}$$

If $n$ is large, then

$$Y_1, \ldots, Y_n \sim f(y_i; \theta)$$

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

- $H_1$: $\theta \neq \theta_0$

$$\Lambda(\theta) = -2 \ln \left[ \frac{L(\theta_0)}{L(\tilde{\theta})} \right]$$

where $\Lambda$ satisfies all the properties of $D$. Also,

$$\lambda(\theta) = -2 \ln \left[ \frac{L(\theta_0)}{L(\hat{\theta})} \right]$$

and

$$p\text{-value} = P(\Lambda \geqslant \lambda) = P(Z^2 \geqslant \lambda)$$

Roadmap:

   (i)   General info

  (ii)   Testing for variance for Normal

 (iii)   An example

The general problem: $Y_1, \ldots, Y_n \sim N(\mu, \sigma^2)$ iid where $\mu$ and $\sigma^2$ are both unknown. $H_0$: $\sigma^2 = \sigma_0^2$ vs two sided alternative.

   (i)   Test statistic? Problem

  (ii)   Convention?

The pivot is:

$$U = \frac{(n-1)s^2}{\sigma_0^2} \sim \chi_{n-1}^2$$

can we use this as our test statistic?

---

**EXAMPLE 0.1.4.**
 - Normal population: $\{y_1, \ldots, y_n\}$
 - $n = 20$, $\sum y_i = 888.1$, $\sum y_i^2 = 39545.03$
 - $H_0$: $\sigma = 2$
 - $H_1$: $\sigma \neq 2$

What is the $p$-value? We know

$$s^2 = \frac{1}{n-1}\left[\sum y_i^2 - n\bar{y}^2\right] = 5.7342$$

Compute $U$:

$$U = \frac{(n-1)s^2}{\sigma_0^2} = 27.24$$

$\chi_{19}^2$

$$
\begin{aligned}
p\text{-value} &= 2P(U \geqslant 27.24) \\
&= 2P(\chi_{19}^2 \geqslant 27.24) \\
&= 10\% \text{ and } 20\%
\end{aligned}
$$

so, $p > 0.1$ means there is no evidence against null-hypothesis.

---

## 2020-03-18

Roadmap:

   (i)   5 min recap

  (ii)   LTRS for large $n$

 (iii)   An example

$Y_1, \ldots, Y_n$ iid $\sim N(\mu, \sigma^2)$

 - $H_0$: $\sigma^2 = \sigma_0^2$
 - $U = \frac{(n-1)s^2}{\sigma_0^2} \sim \chi_{n-1}^2$

We calculated the $p$-value:

$$U = \frac{(n-1)s^2}{\sigma_0^2}$$

If

- $U > \text{median } \chi_{n-1}^2 \implies p\text{-value} = 2P(U \geqslant u)$
- $U < \text{median } \chi_{n-1}^2 \implies p\text{-value} = 2P(U \leqslant u)$

<u>Exercise</u>: Construct the $95\%$ confidence interval for $\sigma^2$. Then, check if $\sigma_0^2(4) \in 95\%$ confidence interval.

- $H_0$: $\sigma^2 = 4$ (more than $10\%$, so it is in the $95\%$ confidence interval)

<u>Likelihood Ratio Test Statistic</u> (one parameter)

$Y_1, \ldots, Y_n$ iid $f(y_i; \theta)$ with $n$ large.

- Sample: $\{y_1, \ldots, y_n\}$
- $\theta = $ unknown parameter
- $H_0$: $\theta = \theta_0$
- $H_1$: $\theta = \theta_0$

<u>Step 1</u>: Test statistic:

$$\Lambda = -2\ln\left[\frac{L(\theta)}{L(\tilde{\theta})}\right]$$

If $H_0$ is true:

$$\Lambda = -2\ln\left[\frac{L(\theta)}{L(\tilde{\theta})}\right] \sim \chi_1^2$$

<u>Step 2</u>: Calculate $\lambda$

$$\lambda = -2\ln\left[\frac{L(\theta_0)}{L(\hat{\theta})}\right] = -2\ln\left[R(\theta_0)\right]$$

$$\begin{aligned} p\text{-value} &= P(\Lambda \geqslant \lambda) \\ &= P(Z^2 \geqslant \lambda) \\ &= 1 - P(|Z| \leqslant \lambda) \end{aligned}$$

**EXAMPLE 0.1.5.** Suppose $Y_1, \ldots, Y_n \sim f(y_i; \theta)$ iid. where

$$f(y, \theta) = \frac{2y}{\theta} e^{-y^2/\theta}$$

Data: $n = 20$, $\sum y_i^2 = 72$
We want to test $H_0$: $\theta = 5$ (two sided alternative).
- $\hat{\theta} = \frac{1}{n}\sum y_i^2 = 3.6$
- $R(\theta_0) = \frac{\hat{\theta}}{\theta_0} e^{(1-\hat{\theta}/\theta_0)^n}$
- $\lambda(\theta_0) = \cdots$

We know $\lambda = -2\ln\left[R(\theta_0)\right] = 1.9402$ and so

$$R(\theta_0) = \frac{L(\theta_0)}{L(\hat{\theta})} = 0.3791$$

also $\theta_0 = 5$. Lastly, calculate the $p$-value.

$$p\text{-value} = P(\Lambda \geqslant \lambda)$$
$$= P(Z^2 \geqslant 1.9402)$$
$$\approx 16.5\%$$

Thus, no evidence against null-hypothesis ($H_0$).

A few final points:

(i) Careful about the previous example.

(ii) $\lambda$ and the relationship with $R$

(iii) Next video

- $n = 20$ is not large
- $\lambda = -2\ln[R(\theta_0)]$: high values of $\lambda \implies$ low values of $R(\theta_0)$

---

## 2020-03-20

Roadmap:

(a) Housekeeping

       Modified Syllabus + Incentives

       Extra materials

       Dropbox link + Mathsoc

(b) Gaussian Response Model: An introduction

Gaussian Response Models

Assumption: $Y_1, \ldots, Y_n \sim$ Normal

Before: $Y_1, \ldots, Y_n \sim N(\mu, \sigma^2)$ iid. with $\mu, \sigma^2 =$ unknown.

$$Y_i = \mu + R_i$$

where $R_i \sim N(0, \sigma^2)$ and $R_i$'s independent for each $i \in [1, n]$. We call:

- $Y_i$ **response** variable
- $\mu$ **systematic part**
- $R$ **random part**

Now:

- $x =$ explanatory variable
- $\mu = \mu(x)$
- $\sigma^2 = \sigma^2(x)$

For example,

$$Y_i \sim N(\mu(x), \sigma^2(x_i))$$

Simple Linear Regression: $\mu = \alpha + \beta x$ and $\sigma^2 =$ constant.

> **EXAMPLE 0.1.6.**
>   - Response: $Y_i = $ STAT 231 score of student $i$
>   - Explanatory (Covariate): $x_i = $ STAT 230 score of student $i$ (given)
>
> Can $Y$ be explained by $x$?
>
> <u>Simple Linear Regression Model</u>
>
> $$Y_i \sim N(\alpha + \beta x_i, \sigma^2)$$
>
> for each $i \in [1, n]$ independent.
>
> Our assumptions are:
>   - $E(Y) = \mu(x) = \alpha + \beta x$
>   - $Y \sim$ Normal
>   - $\sigma^2 = $ constant (independent of $x$)
>   - independent
>
> We want to estimate $\alpha$ and $\beta$.

---

## 2020-03-23

<u>Roadmap:</u>

(i) 5 min recap

(ii) MLE for $\alpha$, $\beta$, $\sigma$

(iii) Least Squares

(iv) Example

<u>Recap:</u>

General: $Y \sim N(\mu(x), r(x))$

<u>Assumptions for the Simple Linear Regression Model</u> (Gauss Markov Assumptions)

(i) One covariate (for the time being)

(ii) Normality: $Y_i$'s are Normal

(iii) Linearity: $E(Y) = \alpha + \beta x$

(iv) Independence: $Y_i$'s are all independent

(v) Homoscedasticity: $\sigma^2 = \sigma^2(x) = \sigma^2$ for all $x$

We call it a Simple since $x$ is the only explanatory variate. If we used more than one explanatory variate, we call it a multi-variable regression (not covered in this course).

<u>MLE Calculation</u>

$$Y_i \sim N(\alpha + \beta x_i, \sigma^2)$$

for each $i \in [1, n]$ independent. We can also write

$$Y_i = (\alpha + \beta x_i) + R_i$$

where $R_i \sim N(0, \sigma^2)$ and $R_i$'s independent.

$$f(y_i) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2\sigma^2}(y_i - (\alpha + \beta x_i))^2}$$

$$L(\alpha, \beta, \sigma) = \frac{1}{(2\pi)^{n/2}\sigma^n} e^{-\frac{1}{2\sigma^2}\sum [y_i - (\alpha + \beta x_i)]^2}$$

so,

$$\ell(\alpha, \beta, \sigma) = -\frac{n}{2}\ln(2\pi) - n\ln(\sigma) - \frac{1}{2\sigma^2}\sum[y_i - (\alpha + \beta x_i)]^2$$

$$\frac{\partial\ell}{\partial\alpha} = 0 \implies \hat{\alpha} = \overline{y} - \hat{B}\overline{x}$$

$$\frac{\partial\ell}{\partial\beta} = 0 \implies \hat{\beta} = \frac{S_{xy}}{S_{xx}} = \frac{\sum(x_i - \overline{x})(y_i - \overline{y})}{\sum(x_i - \overline{x})^2}$$

$$\frac{\partial\ell}{\partial\sigma} = 0 \implies \hat{\sigma}^2 = \frac{1}{n}\sum\left[y_i - (\hat{\alpha} + \hat{\beta}x_i)\right]^2$$