

Principles of Inference: a prelude to Hypothesis Testing

Goal :- Primary objectives of a statistical analysis is to use data from a sample to make inferences about the population from which the sample was drawn.

The sampling distributions discussed earlier in this lecture note (χ^2 , t, F) play a pivotal role in this aspect.

Statistical Inference

a statement (hypothesis)
about the value of
the parameter (that
we are interested
in estimating about
a population by using
a random sample)

a measure of the
reliability of that
statement in terms of
a probability (that the
decision is incorrect)

Hypothesis Testing

A hypothesis usually results from speculation concerning observed behavior, natural phenomena or established theory. If the hypothesis is stated in terms of population parameters such as the mean & variance, the hypothesis is called a statistical hypothesis. Data from a sample are used to test the validity of the hypothesis.

5-Step procedure for Hypothesis Testing

- ↳ Null & Alternate hypothesis -
 H_0 H_1
 - Decision Making / Rejection (Critical) region
 - Errors in Decision
- Types of error α , β
(Type I, Type II)

5-Step procedure for Hypothesis Testing.

Step ① :- Specify H_0 , H_1 & an acceptable level of α .

Step ② :- Define a sample based test statistic (e.g. \bar{X} , S^2 etc) & the rejection region for H_0 .

Step ③ :- Collect the sample data & calculate the test statistic.

Step ④ :- Make a decision to either reject or fail to reject H_0 .

Step ⑤ :- Interpret the result in the language of the problem
(i.e. provide Confidence Interval (C.I.)
Type & probability of error for the kind of decision being made).

Application of $Z \sim N(0, 1)$ D".

Example (Quality control in a packaging..)

A company that packages salted peanuts in 8kg jars is interested in maintaining control on the amount of peanuts put in jars by one of its machines. Control is defined as averaging 8kg per jar & not consistently over/underfilling the jars. To monitor this control, a sample of 16 jars is taken from the packaging line at random time intervals & their contents weighed. The mean wt. of peanuts in these 16 jars (test statistic) will be used to test the null hypothesis that the machine is indeed working properly. If it is found not to be doing so, a costly adjustment will be needed.

Step ①

Hypothesis $\begin{cases} H_0: \mu = 8 \\ H_1: \mu \neq 8 \end{cases}$; set α

$$\text{Test Statistic} \rightarrow \bar{X} = \frac{\sum_{i=1}^n X_i}{n}$$

Step ② & ③

Rejection Region \rightarrow eg. $\{ \bar{X} < 7.9 \text{ or } \bar{X} > 8.1 \}$

Step ④

Type of error

Type(I) :- When H_0 is actually true & our analysis based on our test statistic rejects H_0 .

Type(II) :- H_0 is actually not true & our inference based on test statistic fails to reject H_0 .

IN THE POPULATION

The Decision	H_0 is True	H_0 is NOT True
H_0 is NOT rejected	Decision is Correct	type II error (w/ probability β)
H_0 is rejected	Type I error (w/ probability α)	Decision is correct

Calculating α

$$\alpha = P(\bar{X} < 7.9 \text{ or } \bar{X} > 8.1 \text{ when } \mu = 8)$$

Assume (for now) that we somehow know

σ of the population (of jars) to be 0.2.

$n = 16$ (This will often not be the case & n may be \rightarrow for analysis)

$$\begin{aligned} P(\bar{X} < 7.9) &= P\left(\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} < \frac{7.9 - 8}{0.2/\sqrt{16}}\right) \\ &= P(Z < -2.0) = 0.0228 \end{aligned}$$

and

$$\begin{aligned} P(\bar{X} > 8.1) &= P\left(\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} > \frac{8.1 - 8}{0.2/\sqrt{16}}\right) \\ &= P(Z > 2.0) = 0.0228 \end{aligned}$$

$$\begin{aligned} \alpha &= P(\bar{X} < 7.9) + P(\bar{X} > 8.1) \\ &= 0.0228 + 0.0228 = 0.0456 \\ &= \text{Prob (Type I error)} \end{aligned}$$

Probability of adjusting the m/c when it does not need it is slightly less than 0.05 (5%).

Calculating β

This is not always straightforward.

A type I error is committed when H_0 is rejected if $\bar{X} \neq 8$ even though $\mu = 8$.

The sampling Dⁿ for the test statistic (\bar{X}) is different for different values of $\mu \neq 8$.

For practical purpose, several different values for β is obtained by considering a handful of μ values eg. $\{7.80, 7.90, 7.95,$
 $8.05, 8.10, 8.02\}$.
Each of this μ value will give a range of \bar{X} for a globally fixed α .

These β vs μ gives us an "Operating characteristic (OC)" curve.

Power of a test = $(1 - \beta)$; power vs μ is called power curve.

Reading assignment-

Importance of power of a test.

* There is always a trade-off between α & β . reducing β will inc. α & vice-versa
but generally α is more important to control/reduce at the expense of a reasonable β .

The significance level of a test is the maximum acceptable probability of rejecting a true NULL hypothesis (i.e. max^m allowable type I error probability).

→ This will enable us to define a rejection region disadv of significance level.

- I) One may not have a clear idea of what an appropriate max^m allowable type I error probability should be for a given test statistic.
- II) Significance level may also be sensitive to minor changes in sample statistic.

Alternatively, a method of reporting the results of a significance test w/o having to choose an exact level of significance, but instead leave that decision to the individual who will actually act on the conclusion of the test —

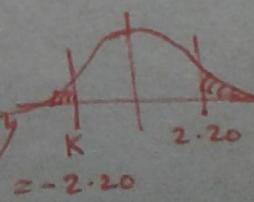
p-value report.

$$P\left(Z < \frac{7.89 - 8}{0.2/\sqrt{16}}\right) = P(Z < 1.6) = P(Z < -2.2)$$

e.g. if $\bar{X} = 7.89$ is obtained

$$P(|Z| > 2.2) \xrightarrow{\text{symmetry}} 2P(Z > 2.2) \\ \xrightarrow{N(0,1)} = 0.0278$$

∴ We can say that the probability of observing a test statistic at least this extreme if H_0 is true is 0.0278.



Truth (Actual Situation)

Decision

Do NOT
reject H_0

Reject H_0

		H_0 True	H_0 false
	H_0 True	Correct decision w/ prob $1 - \alpha$	Incorrect Decision β
	H_0 false	Incorrect Decision w/ prob α	Correct Decision $1 - \beta$
		Type I error	Type II error

$$\alpha = \text{Prob}(\text{Type I error})$$

$$\beta = \text{Prob}(\text{Type II error})$$

Application of t-Distribution in hypothesis testing

Two sample test for mean

$$X_1 \sim N(\mu_1, \sigma^2) \\ X_2 \sim N(\mu_2, \sigma^2)$$

Note σ^2 is same for both samples R.v.

$$X_{1i} \sim N(\mu_1, \sigma^2); i=1, 2, \dots, n_1$$

$$X_{2i} \sim N(\mu_2, \sigma^2); i=1, 2, \dots, n_2$$

Step 1:-

$$H_0: \mu_1 = \mu_2$$

$$H_1: \mu_1 \neq \mu_2 \text{ or } \mu_1 > \mu_2 \text{ or } \mu_1 < \mu_2$$

Also, select α level of significance.

Step ② & ③: Test Statistic

$$t = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{S \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

where $S^2 = \frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{(n_1+n_2-2)}$; s_i^2 is sample variance of i^{th} group.

Under H_0 :

$$t = \frac{\bar{X}_1 - \bar{X}_2}{S \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim t(n_1+n_2-2)$$

Step ④ Rejection (Critical region)

Reject H_0 in favor of $H_1 (\mu_1 \neq \mu_2)$ if
 $|t| \geq t(\alpha/2, n_1+n_2-2)$

Reject H_0 in favor of $H_1 (\mu_1 > \mu_2)$ if
 $t \geq t(\alpha, n_1+n_2-2)$

Reject H_0 in favor of $H_1 (\mu_1 < \mu_2)$ if
 $t \leq -t(\alpha, n_1+n_2-2)$

example

In comparison of two gasoline brands, a consumer survey reveals the following:

- A full tank of Brand (A) requires 4 cans and covers 546 kms w/ a std. deviation of 31 kms.
- A full tank of Brand (B) requires 4 cans & covers 492 kms w/ a std. deviation of 26 kms.

Assume that both populations (A) & (B) are sampled from Normal Dⁿs w/ equal variances. Test

$$H_0: \mu_1 = \mu_2$$

$$H_1: \mu_1 > \mu_2 \text{ at } 0.05 \text{ level of significance}$$

Soln :- $H_0: \mu_1 = \mu_2$

$$H_1: \mu_1 > \mu_2$$

$$\alpha = 0.05$$

$$\bar{x}_1 = 546, s_1 = 31, n_1 = 4$$

$$\bar{x}_2 = 492, s_2 = 26, n_2 = 4$$

$$S^2 = \frac{(4-1)31^2 + (4-1)26^2}{4+4-2}$$

$$\Rightarrow S = 28.609$$

$$\therefore t(\text{under } H_0) = \frac{546 - 492}{28.609 \sqrt{\frac{1}{4} + \frac{1}{4}}} = 2.67$$

from the table of t-Dⁿ:-

$$t(0.05, 6) = 1.943 \text{ (one-tail t-Dⁿ)}$$

$$\sim t(0.05, 4+2)$$

$$\therefore t = 2.67 > t(0.05, 6) \Rightarrow \text{Reject } H_0 \text{ in favor of } H_1.$$

$$= t(0.05, 6)$$

Analysis of Variance

30/6/18
Pg ①

The t - statistic test for 2 means cannot be generalized for multiple - groups/population means. For such a case, we require Anova.

Data :-

$y_{ij}, i=1, \dots, t, j=1, \dots, n_i$ | y_i i.e. t groups
 $\sum n_i$ observations

If $n_i = n \forall i \Rightarrow$ "Balanced" data.

Null hypothesis :- $H_0: \mu_1 = \mu_2 = \dots = \mu_t$
 $H_1:$ at least 1 equality is not satisfied

Assumption :- Each group is distributed $N(\mu_i, \sigma^2)$
Note σ^2 is same across population groups.

Organization of Data

Factor levels	Observations	Totals	Means	Sum
				sq's
1	$y_{11} y_{12} \dots y_{1n_1}$	$y_{1.}$	$\bar{y}_{1.}$	SS_1
2	$y_{21} y_{22} \dots y_{2n_2}$	$y_{2.}$	$\bar{y}_{2.}$	SS_2
\vdots	\vdots	\vdots	\bar{y}_{\vdots}	\vdots
i	$y_{i1} y_{i2} \dots y_{in_i}$	$y_{i.}$	$\bar{y}_{i.}$	SS_i
t	$y_{t1} y_{t2} \dots y_{tn_t}$	$y_{t.}$	$\bar{y}_{t.}$	SS_t
Overall		$y_{..}$	$\bar{y}_{..}$	SS_p

Factor level totals

$$Y_{ij} = \sum_j Y_{ij}$$

Factor level means

$$\bar{y}_{ij} = \frac{Y_{ij}}{n_i}$$

Overall total $Y_{..} = \sum_{i,j} Y_{ij}$

Overall mean $\bar{y}_{..} = \frac{Y_{..}}{\sum_i n_i}$

The computation of estimated variance follows through the calculations below

$$SS_i = \sum_j (Y_{ij} - \bar{y}_{ij})^2 ; i=1, 2, \dots, t$$

$$\text{or} \quad \sum_j Y_{ij}^2 - \frac{(\bar{y}_{..})^2}{n_i}$$

Pooled sum of squares

$$SS_p = \sum_i SS_i$$

Pooled degrees of freedom $\sum n_i - t$

$$s_p^2 = \frac{SS_p}{\sum n_i - t} = \frac{\sum_i SS_i}{\sum n_i - t}$$

If the individual variances are available s_i^2 ,

$$\text{then } s_p^2 = \frac{\sum_i (n_i - 1) s_i^2}{\sum n_i - t}$$

$$S_{\text{means}}^2 = \frac{\sum_i (\bar{y}_{..} - \bar{y}_{..})^2}{(t - 1)}$$

is the variance estimate of the factor level means.

Under the null hypothesis &

By the 1st thm. of Sampling Dⁿ; the factor level means have Dⁿ w/ mean μ & variance = σ^2/n

varied that variance of samples at each factor level is the same, σ^2 . Pg 9

$$S_{\text{means}}^2 = \frac{\sigma^2}{n} \Rightarrow n S_{\text{means}}^2 = \sigma^2 \text{ w/ } (t-1) \text{ d.o.f}$$

An alternate estimate of σ^2 is S_p^2 w/ $t(n-1)$ d.o.f
(here $n_i = n + i$)

Recall from the def'n. of F, D^{*}; the F - D^{*}
Represents the ratio of 2 independent estimates of
a common variance.

$$\therefore F_{\text{calculated}} = \frac{n S_{\text{means}}^2}{S_p^2} > F_{\alpha}(t-1, t(n-1))$$

Then reject H_0 .

example:- An experiment to compare the yield of 4 varieties of rice was conducted. Each of 16 plots on a test-farm where soil fertility was fairly homogeneous was treated alike relative to H₂O & fertilizer. Four plots were randomly assigned each of the 4 varieties of rice. The yield in kg/acre was recorded for each plot for this randomized experiment. Do the data presented in the following table indicate a difference in the mean yield bet'n the 4 varieties?

Variety	Yields				N _i	\bar{Y}_i	SS _i
	1	2	3	4			
1	934	1041	1028	935	3	984.50	10085.00
2	882	963	924	946	3	928.25	3868.75
3	987	951	976	840	3	938.50	13617.00
4	992	1143	1140	1191	3	1116.50	22305.00
Overall					15	991.94	49875.75

Soln :-

Hypothesis $H_0: \mu_1 = \mu_2 = \mu_3 = \mu_4$

$H_1:$ not all varieties have the same mean

Where μ_i is the mean yield per acre for variety i .

$$S_{\text{mean}}^2 = n \sum_i (\bar{y}_{i\cdot} - \bar{y}_{\cdot\cdot})^2 / (t-1) ; \begin{matrix} n=4 \\ t=4 \end{matrix}$$
$$= 4 \left\{ (984.5 - 991.94)^2 + \dots + (1116.50 - 991.94)^2 \right\} / (4-1)$$
$$= 29977.06$$

$$S_p^2 = \sum_i Ss_i^2 / t(n-1) = \frac{(10,085 + \dots + 22,305)}{4 \times 3} = \frac{49875.75}{12}$$
$$= 4156.31$$

$$F_{\text{cal}} = \frac{29977.06}{4156.31} = 7.21$$

$\alpha = 0.01$ (^{says} set opinion by statistician)

$$F_{0.01}(3, 12) = 5.95$$

$\uparrow \quad \uparrow$

d.o.f in numerator
 $= (t-1)$
 $= 3$

d.o.f in denominator
 $t(n-1)$
 $= 12$

$$\therefore F_{\text{cal}} > F_{0.01}(3, 12)$$

Reject H_0 !

i.e. Difference may exist in the yields of 4 varieties.

Computational convenience for calculating F - Statistic

Pg③

Between group sum of squares,

$$SSB = \sum_i \frac{y_{i.}^2}{n_i} - \frac{y_{..}^2}{\sum n_i} \quad w/ df_B = (t-1)$$

W/in group sum of squares,

$$SSW = \sum_{i,j} y_{ij}^2 - \sum_i \frac{y_{i.}^2}{n_i} \quad w/ df_W = \sum n_i - t$$

Total sum of squares,

$$TSS = SSB + SSW$$

$$\begin{aligned} w/ df_T &= df_B + df_W \\ &= \sum n_i - 1 \end{aligned}$$

1-way Analysis of Variance table

Source	df	ss	$MS = \frac{ss}{df}$	F
Bet'n groups	t-1	SSB	MSB	$\frac{MSB}{MSW}$
W/in groups	$\sum n_i - t$	SSW	MSW	
Total	$\sum_i n_i - 1$	TSS		

ANOVA table for our earlier example

Source	df	ss	MS	F
Bet'n varieties	3	89,931.19	29977.06	7.21
W/in varieties	12	49,875.75	4156.31	
Total	15	139,806.94		

Same concl
as before