① 

# t - test for single mean :

(i) If a random sample $x_i$ $(i = 1, 2 \ldots n)$ of size n has been drawn from a random population with specified mean say $\mu_0$, or.

(ii) If the sample mean differs significantly from the hypothetical value $\mu_0$ of the population mean.

Under the null hypothesis Ho :

(i) the sample has been drawn from the population with mean $\mu_0$ or

(ii) there is no significant differen b/w the sample mean $\bar{x}$ and the population mean $\mu_0$.

the statistic,

$$t = \frac{\bar{x} - \mu_0}{s/\sqrt{n}}.$$

where $\bar{x} = \frac{1}{n} \sum_{i=1}^{n} x_i$ and.

$$s^2 = \frac{1}{n-1} \sum_{i=1}^{n} (x_i - \bar{x})^2.$$

follows student's dist with $(n-1)$ d.f.

we now compare the calculated value of $t$, with the tabulated value at certain level of significance.

If calculated $|t| >$ tabulated $t$, null hypothesis is rejected and oof calcuted $|t| <$ tabulated $t$, Ho may be accepted at the level of significance ~~adop~~ accept.

Remark :1

on computation of $s^2$ for numerical problems:- If $\bar{x}$ comes out in integers, the formula can be convonentely used for computing $s^2$. However if $\bar{x}$ comes in fractions case, step deviation method given below, is quite useful.

If we take $di = xi - A$ where $A$ is any orbitary number, then.

$$S^2 = \frac{1}{n-1}\left[\sum[x_1 - \bar{x}]^2\right] = \frac{1}{n-1}\left[\sum x_i^2 - \left(\frac{\sum x_i^2}{n}\right)\right]$$

$$= \frac{1}{n-1}\left[\sum di^2 - \left(\frac{\sum di}{n}\right)^2\right] \text{ since variance}$$

is independent change of orgin.

Also in case $\bar{x} = A + \frac{\sum fdi}{N}$

2.) we know, the sample variance,

$$S^2 = \frac{1}{n}\sum_i^? (xi - \bar{x})^2.$$

$$\Rightarrow nS^2 = (n-1)s^2.$$

$$\frac{S^2}{n} = \frac{s^2}{(n-1)}.$$

Hence for numerical problems the test statistic on using becomes.

$$t = \frac{\bar{x} - \mu_0}{\sqrt{s^2/n}} = \frac{\bar{x} - \mu_0}{\sqrt{s^2/n-1}} \sim t_{n-1}.$$

paired t-test for difference of mean:

let us know consider the case when (i) the sample sizes are equal.

i.e., $n_1 = n_2 = n$ (say) and (ii) the two sample ④ are not independent but the sample observations are paired together.

i.e the paired of observations $(x_i, y_i)$ $(i = 1, 2 \ldots n)$ corresponds to the same (ith) sample unit. The problem is to test if the sample means differ significantly or not.

For example, suppose we want to test the efficacy of a particular drug, say the including sleep. Let $x$ and $y_i$ $(i = 1, 2 \ldots n)$ be the readings, in hours of sleep, on the individual before and after the drug is given respectively. Here instead of applying the difference of the means test discussed in, we apply the paired t-test given below.

Here we consider the increments

$$dp = x_i - y_i \quad (i = 1, 2 \ldots n).$$

under the null hypothesis, Ho that increments are due to fluctuations of sampling.

i.e the clag is not responsible for these increments , the statistic.

$$t = \frac{\bar{d}}{s/\sqrt{n}}.$$

where $\bar{d} = \frac{1}{n} \sum_{q=1}^{n} di$ and

$$s^2 = \frac{1}{n-1} \sum_{i=1}^{n} (di - \bar{d})^2.$$

follows students t-distribution with (n-1)df.

### T-test for diff means:

suppose we want to test if

the independent samples $xi (i=1,2...n)$

and $yi (i=1,2...n)$ of sizes $n_1$ and $n_2$

here been drawn from two normal

population with means $\mu_x$ and $\mu_r$ respectively.

Under the null hypothesis [Ho] that

the samples have been drawn from the

normal populations with means $\mu_x$ and

$\mu_r$ and under the assumption that

the population variance an equal

i.e., $(\sigma_x^2 = \sigma_r^2 = \sigma^2 (say))$. the statistic

$$t = \frac{(\bar{x} - \bar{y}) - (\mu_x - \mu_r)}{s \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}.$$

where $\bar{x} = \dfrac{1}{n_1} \sum\limits_{i=1}^{n_1} x_i$ & $\bar{y} = \dfrac{1}{n_2} \sum\limits_{i=1}^{n_2} y_i$

and $S^2 = \dfrac{1}{n_1+n_2-2} \left[ \sum\limits_i (x_i-\bar{x})^2 + \sum\limits_i (y_i-\bar{y})^2 \right]$.

is an unbiased estimate of the mean common population variance $\sigma^2$, following students distribution $(n_1+n_2-2)$ d.f.

proof :

distribution of $t$ defined in

$$\xi_i = \dfrac{(\bar{x}-\bar{y}) - E(\bar{x}-\bar{y})}{\sqrt{V(\bar{x}-\bar{y})}} \sim N(0,1).$$

put $E(\bar{x}-\bar{y}) = E(\bar{x}) - E(\bar{y}) = \mu_x - \mu_y$.

$V(\bar{x}-\bar{y}) = V(\bar{x}) + V(\bar{y}) = \dfrac{\sigma_x^2}{n_1} + \dfrac{\sigma_y^2}{n_2}$.

$$= \sigma^2 \left( \dfrac{1}{n_1} + \dfrac{1}{n_2} \right).$$

$$\xi_i \dfrac{(\bar{x}-\bar{y}) - (\mu_x - \mu_y)}{\sqrt{\sigma^2 \left( \dfrac{1}{n_1} + \dfrac{1}{n_2} \right)}} \sim N(0,1) \rightarrow ①$$

let $\chi^2 = \dfrac{1}{\sigma^2} \left[ \sum\limits_{i=1}^{n} (x_i-\bar{x})^2 + \sum\limits_{i=}^{n} (y_i-\bar{y})^2 \right]$.

$$= \left[ \sum\limits_i (x_i-\bar{x})^2 / \sigma^2 \right] + \left[ \sum\limits_i (y_i-\bar{y})^2 / \sigma^2 \right].$$

$$= \dfrac{n_1 s_x^2}{\sigma^2} + \dfrac{n_2 s_y^2}{\sigma^2} \rightarrow ②$$

Since $n_1 s_x^2/\sigma^2$ and $n_2 s_y^2/\sigma^2$ are independent ()
$x^2$ variates with $(n_1-1)$ and $(n_1-1)$ df respectively
& by the additive property of chi-square
distribution $x^2$ defined ① ② is a $x^2$ variate
with $(n_1-1)+(n_2-1)$.

i.e $n_1+n_2-2$ df further, since sample
variance are independently distributed $s$ and
$x^2$ are independent distribution variables.

Hence fisher's t-Statistic is given by

$$t = \frac{\bar{y}}{\sqrt{\frac{x^2}{n_1+n_2-2}}}$$

$$= \frac{(\bar{x}-\bar{y})-(\mu_x-\mu_y)}{\sqrt{\sigma^2\left(\frac{1}{n_1}+\frac{1}{n_2}\right)}} \times \frac{1}{n_1+n_2-2\{\sum_i c_i (x_i-\bar{x})^2 + \sum_i c_i (y_i-\bar{y})^2\}}$$

$$= \frac{(\bar{x}-\bar{y})-(\mu_x-\mu_y)}{s\sqrt{\frac{1}{n_1}+\frac{1}{n_2}}} \quad \text{where } s^2 = \frac{1}{n_1+n_2-2}$$

$$\left[\sum_i (x_i-\bar{x})^2 + \sum_i (y_i-\bar{y})^2\right].$$

and it follows students t-distribution
with $(n_1+n_2-2)$ df.

## Remark:

$s^2$ defined in is an unbiased estimate of the common population curve $\sigma^2$, since

$$E(s^2) = \frac{1}{n_1+n_2-2} \, E\left[\sum_1 (x_i-\bar{x})^2 + \sum_1 (y_i-\bar{y})^2\right].$$

$$= \frac{1}{n_1+n_2-2} \, E[(n_1-1)\,s_x^2 + (n_2-1)\,s_y^2].$$

$$= \frac{1}{n_1+n_2-2} \, [(n_1-1)\,E(s_x^2) + (n_2-1)\,E(s_y^2)].$$

$$= \frac{1}{n_1+n_2-2} \, [(n_1-1)\sigma^2 + (n_2-1)\sigma^2] = \sigma^2.$$

An important deduction which is of much practical utility is discussed below.

Suppose we want to test if : (a) two independent samples $x_i \, (i=1,2\ldots n_1)$ and $x_i \, (1,2\ldots n_2)$ have been drawn from the populations with same means or (b) the two samples mean $\bar{x}$ and $\bar{y}$ differ significantly or not.

t-test for testing the significance of an observed sample Correlation:

If r is the observed correlation co-efficient in a sample of n point of observations from a bivariate normal population, then proof, Fisher proved that under the null hypothesis $H_0: p=0$ i.e., population correlation co-eff is zero,

the statistic,

$$t = \frac{r}{\sqrt{(1-r^2)}} \sqrt{n-2}.$$

follows student's t-dist with (n-2)df. If the value of t comes out to be significant, we reject $H_0$ at the level of significance adopted and conclude that $p \neq 0$ i.e 'r' is significant or Correlation in the population.

If t comes out be non-significant then $H_0$ may be accepted and we conclude that variables may be regarded as un-correlated in the population.

## Application of the chi-square Test:

$\chi^2$ dist has a large number of applications in statistics, some of when enumerated below.

(i) to test if the hypothetical value of population variance is $\sigma^2 = \sigma_0^2$.

(ii) To test the "goodness of fit".

(iii) To test the independence of attributes.

(iv) to test the homogenaity of independent estimates of the population variance.

(v.) To combine various probablities obtained from independent experience give a sample test of significance.

(vi.) to test the homogeneity of independent estimates of the population Correlation Co-efficient.

# Goodness of fit test:

A very powerful test for testing the significance of the discrepancy b/w theory and experiment was given by Proof. Karl pearson 1900 and in known as " chisquare test of goodness of fit ". It enables us to find it the deviation of the experiment from theory is just by chance (or) is it really due to for independency of the theory to fit the observed data.

If $f_i$ (1,2,...n) is a set of observed (experimental) frequencies and i.e $(i=1,2...n)$ is the corresponding set of expected (theoretical or hypothetical) frequency then Kar pearson's chisquare given by

$$\psi^2 = \sum_{j=1}^{n} \left[ \frac{(f_i - e_i)^2}{e_i} \right] \left[ \sum_{j=1}^{n} f_i = \sum_{i=1}^{n} e_i \right]$$

follows chi-square dist with $(n-1)$ d.f.

⑫

Remark:

This is an approximate test for large values of $n$.

The goodness of fit test uses the chi-square test to determine if a hypothesis prob distribution for a population provides a good fit.

Decision rule:

Accept $H_0$ if $\phi^2 \leq \psi^2_{\alpha}(n-1)$

and reject $H_0$ if $\phi^2 > \psi^2_{\alpha}(n-1)$ where is the calculated value of chi-square obtain on using and $\psi^2_{\alpha}(n-1)$ is to tabulated value of chisquare of $(n-1)$ d.f and level of significance $\alpha$.

test of independence of attributes

Contingency tables:

Let us consider two attributes A and B, A divided to s classes $B_1, B_2 \dots B_s$ such a classification in which attributes are divided into where than two classes in known as monifold classification.

The various cell frequencies on be expressed in the following table known as $\partial \times s$ manifold contingency table where is the number of persons possessing the attribute $A_i$ ($i=1, 2 \dots \partial$) $(B_i)$ is the no. of persons possessing attribute $B_j (j=1, 2 \dots s)$. and $(A_i B_i)$ is the no. of persons possessing both the attributes $A_i$ and $B_j$ ($i=1, 2 \dots \partial$, $j=1, 2 \dots s$).

Also $\sum_{i=1}^{\partial} A_i = \sum_{j=1}^{s} B_j = N$.

( $\because$ where N is the total frequency).

| A / B | A1 | A2 ... | Ai ... | Aᵧ | total |
|---|---|---|---|---|---|
| B1 . | A1/B1 | A2B1 ... | AiBi | AᵧiB1 | B1 |
| B2 | A1/B2 | A2B1 ... | AiB2 | AᵧB2 | B2 |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ |
| Bj | A1 Bj | A2Bj ... | AiBj | AᵧBj | Bj |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ |
| Bs | A1,BS | A2Bs ... | AiBS | AᵧBS | Bs |
| total | A1 | A2 | Ai | Aᵧ | N. |

the problem is to test if the two attributes

A and B under Consider to independent or not.

under the null hypothesis that the attributes

are independent theory of frequency are

calculated follows.

$P[Ai]$ = prob that a person possesses

that attribute $Ai = \dfrac{(Ai)}{N}, i = 1, 2 \dots n.$

$P[Bj]$ = prob that a person possesses

the attribute $Bj = \dfrac{(Bj)}{N} \; j = 1, 2 \dots n.$

$P[A_i B_j]$ = prob that a person posseses the attributes $A_i$ and $B_j$ = $P(A_i) P(B_j)$.

(By compound prob theorem, since the attributes $A_i$ and $B_j$ are independent under the null hypothesis.)

$$P[A_i B_j] = \left(\frac{A_i}{N}\right)\left(\frac{B_j}{N}\right); \ell = 1,2\ldots\vartheta.$$

$(A_i B_j)_0$ = Expected no. of persons possesing both the attributes $A_i$ and $B_j$

$$= N \cdot P[A_i B_j] = \frac{A_i B_j}{N}.$$

$$\Rightarrow (A_i B_j)_0 = \frac{(A_i)(B_j)}{N}, \quad (i=1,2\ldots s).$$

By using this formula, we want find out expected frequencies for each of the frequencies $(A_i, B_j)$ $(i=1,2\ldots r, j=1,2\ldots s)$ under the null hypothesis at independence of attributes.

The exact test for the independence of attributes is very complicated but a fair degree of approximation is given, for large samples; (large N), by the $\psi^2$ test of goodness of fit, viz....

$$\phi^2 = \sum_{i=1}^{r} \sum_{j=1}^{s} \left[ \frac{\left[AiBj - AiBj'\right]_0^2}{AiBj_0} \right]$$ (16)

$$= \sum_i \sum_j \frac{(fij - eij)^2}{eij}.$$

where , $fij$ = observed frequency for contingency table category in column $i$ and row $j$.

$eij$ = ex. frequency for contigency, table category in column $i$ and row $j$, which is distributed as a $\psi^2$ variate with $(r-1)(s-1)$ df

[c.f Note b/w on d.f].

**Remark:**

$$\phi^2 = \psi^2/N$$ is known as mean square contingency.

Since the limits for $\psi^2$ and $\phi^2$ very in different cases, they cannot be used for establishing the closeness of the relationship b/w quantitative characters under study.

**proof:** Karl pearson suggested another measure, known as Co-eff of mean square contingency which is denoted by $c$ and is given by":

$$C = \sqrt{\frac{\psi^2}{\psi^2 + N}} = \sqrt{\frac{\phi^2}{1 + \phi^2}}.$$

Obviously c is always less than unity. ⑰

The max value of c depends on $r$ and s. the no. of. classes into which A and B divided. In a $r \times r$ Contingency table, the max value $c = \sqrt{\left(\frac{r-1}{r}\right)}$.

Since the max value of c differs for different classification. viz.. $r \times r (r = 2, 3, 4..)$ strictly, speaking, the values of c obtained from different types of classifications. are not comparable.