

# Data Analytics

Course Taught at IIFT

*Day 6: Analysis of Variance (ANOVA)*

**Dr. Tanujit Chakraborty**

*Centre for Data Sciences*

IIIT Bangalore

# Today's Topics.....

- Hypothesis Tests and P-values
- What is Analysis of variance?
- Why ANOVA?
- How to do ANOVA?
  - *One – way ANOVA*
  - *Two–way ANOVA*



# What is Hypothesis Tests?

# Hypothesis Tests...

- A hypothesis test is a form of statistical inference where we attempt to answer a specific question about the distribution of some measurement in the population.
- Basically, we want to know “Yes or No, does the following statement hold true for the distribution of our measurement in our population?” (Most often, our hypothesis will be about one or two population means or one or two population proportions).
- When we make such a decision, we must keep in mind that we could make an incorrect decision. Hence, we want to design our decision procedures so as to minimize the chances of errors. We will formally see how this is done in the context of statistical decisions below.
- We will restrict our attention to the case of two alternatives. The first will be called the **and** and will be represented by the symbol  $H_0$  (**null hypothesis**). The other will be called the **alternative hypothesis** and will be represented by the symbol  $H_a$ .

# Hypothesis Tests...

- The alternative hypothesis is usually some positive statement about the experiment being analyzed and is often called the “**research hypothesis**” because it is what the researcher is trying to show in the experiment.
- The null hypothesis is often a representation of the status quo or other representation of lack of positive results in the experiment. Hence the aspect of “**null**” about the null hypothesis.
- **Example:** In an analgesic drug experiment the alternate hypothesis might be that the new drug relieves pain faster than the old drug. The null hypothesis would then be that the new drug does not relieve pain faster.
- **Example:** In a teaching experiment, the alternative hypothesis might be that the new method results in higher scores than the old method. The null hypothesis would be that this is not so.



What is P-value?

# P-values...

- P-value is a measure of consistency between the null hypothesis and the observed data.
- Hence, we will be more inclined to believe the alternative hypothesis when the P-value is small and less likely to believe the alternative hypothesis when the P-value is not so small.
- We make a rule for statistical decisions as follows. We establish a cut-off value called the **level of significance**.
- Then we reject the null hypothesis in favor of the alternative hypothesis when the P-value is less than the level of significance.
- If the P-value is not less than the level of significance, then we fail to reject the null hypothesis in favor of the alternative hypothesis.

# P-values...

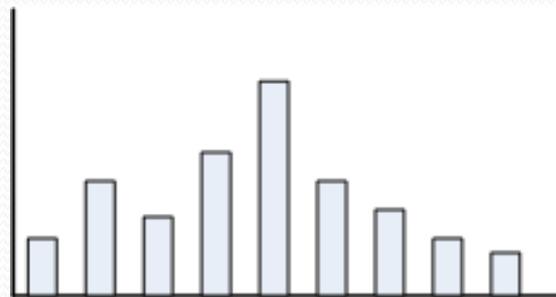
- If the P-value is small, the data is inconsistent with  $H_0$ , so we reject  $H_0$  in favor of  $H_a$ .
- If the P-value is not small, the data is not inconsistent with  $H_0$ , so we fail to reject  $H_0$  in favour of  $H_a$ .
- The natural question at this point is “How small does the P-value have to be to be small?” (or, “What amount of evidence is beyond reasonable doubt?”)
- If the P-value is less than  $\alpha$ , then we reject  $H_0$  in favor of  $H_a$ .
- **Remember :** P-value is NOT the probability that the null hypothesis is true.

| $\alpha$ | <b>Interpretation</b> (if p-value is less than the given $\alpha$ ) | <b>Strength of Evidence</b> (if p-value is less than the given $\alpha$ ) |
|----------|---------------------------------------------------------------------|---------------------------------------------------------------------------|
| 0.10     | Approaching statistical significance                                | Potential Evidence                                                        |
| 0.05     | Statistically significant                                           | Sufficient Evidence                                                       |
| 0.01     | Highly statistically significant                                    | Strong Evidence                                                           |
| 0.001    | Extremely statistically significant                                 | Very Strong Evidence                                                      |

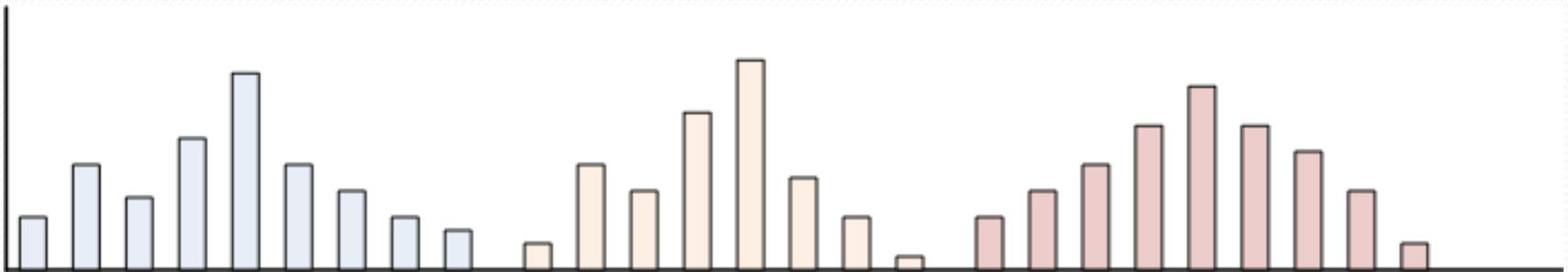


# What is Analysis of Variance?

# What is analysis of variation?



Single population



Multiple population

# Example : Single vs. Multiple population



# What is the issue?

- Are the statistical inference valid?

$\mu$

$\sigma$

# Example 1: The issue in Statistical Testing

A recent study claims that using music in a class enhances the concentration and consequently helps students absorb more information.

- What if it affected the results of the students in a negative way?
- or
- What kind of music would be a good choice for this?

We should have some proof that it actually works or not.

# Design of Experiment

- The teacher decided to implement it on a smaller group of randomly selected students from **three different** classes.

Three different groups of **ten randomly selected students** from three different classrooms were taken.

Each classroom was provided with **three different environments** for students to study.

- Classroom A had **constant music** being played in the background
  - Classroom B had **variable music** being played in the background
  - Classroom C was a regular class **with no music playing**
- A test was conducted after one month for all the three groups and their test scores were collected.

# Test Result

|                          | Test scores of students (out of 10) |   |   |   |   |   |   |              |   |     | Mean |
|--------------------------|-------------------------------------|---|---|---|---|---|---|--------------|---|-----|------|
| Class A (constant music) | 7                                   | 9 | 5 | 8 | 6 | 8 | 6 | 10           | 7 | 4   | 7    |
| Class B (variable music) | 4                                   | 3 | 6 | 2 | 7 | 5 | 5 | 4            | 1 | 3   | 4    |
| Class C (no music)       | 6                                   | 1 | 3 | 5 | 3 | 4 | 6 | 5            | 7 | 3   | 4.3  |
|                          |                                     |   |   |   |   |   |   | Grand Mean : |   | 5.1 |      |

# Observations from the results

- It is noticed that the mean score of students from **Group A** is definitely greater than the other two groups, so the treatment must be helpful.
- Maybe it's true, but there is also a slight chance that we happened to select the best students from class A, which resulted in better test scores (remember, the selection was done at random).
- This leads to a few questions:
  1. How do we decide that these three groups performed differently because of the different situations and **not merely by chance**?
  2. In a statistical sense, how **different are these three samples** from each other?

# Analysis of Variance (ANOVA)

## Definition :

- Analysis of Variance (ANOVA) is derived from a partitioning of total variability into its component parts.
  - ANOVA is a statistical technique that is used to check **if the means of two or more groups are significantly different from each other.**
    -
  - ANOVA checks the impact of one or more factors by comparing the means of different samples.
- 
- This technique was invented by **Sir Ronald Aylmer Fisher** (1921) and is often referred to as Fisher's ANOVA.



# Why ANOVA?

# Statistical Inferences

- ANOVA is a statistical technique
  - It is similar in application to techniques such as t-test, Z-test and  $\chi^2$ -test in that it is used to compare means and the relative variance between them.
- Why not use t-test, Z-test and  $\chi^2$ -test ?
- Why analysis of variance for comparing means?

t-test is used to:

- To infer **mean of a single population**
- t-test can be used to compare two populations

*However, t-test is not useful to compare mean of more than two populations.*



# Extending the two population procedure

- Construct pairwise comparison on all means.
- For 5 populations  $\rightarrow$  10 possible pairs.
- Considering  $\alpha = 0.05$ , probability of correctly failing to reject the null hypothesis for all 10 tests is  $(0.95)^{10}$ , assuming that the tests are independent
- Thus the true value of  $\alpha$  for this set of comparison is 0.4, instead of .05
- It inflates the Type I error.

# Extending the two population procedure

- Statistical Inference I
  - A car magazine wishes to compare the average petrol consumption of **THREE** models for car and has available **SIX** vehicles of each model.

| Model 1 | Model 2 | Model 3 |
|---------|---------|---------|
|         |         |         |
|         |         |         |
|         |         |         |
|         |         |         |
|         |         |         |
|         |         |         |

- There are **THREE** populations
- There are samples each of size six from each population

# Extending the two population procedure

- Statistical Inference II
  - A teacher is interested in a comparison of the average percentage marks obtained in the examinations of five different subjects and has available the marks of eight students who all completed each examination.

| Subject 1 | Subject 2 | Subject 3 | Subject 4 | Subject 5 |
|-----------|-----------|-----------|-----------|-----------|
|           |           |           |           |           |
|           |           |           |           |           |
|           |           |           |           |           |
|           |           |           |           |           |
|           |           |           |           |           |
|           |           |           |           |           |

- What is the number of populations?
- How many samples? What are their sizes?? Are each samples independent to each other?

# Example 2 : Why ANOVA?

Consider the two sets of contrived data as shown below:

| Set 1           |                  |                  | Set 2           |                  |                  |
|-----------------|------------------|------------------|-----------------|------------------|------------------|
| Sample 1        | Sample 2         | Sample 3         | Sample 1        | Sample 2         | Sample 3         |
| 5.7             | 9.4              | 14.2             | 3.0             | 5.0              | 11.0             |
| 5.9             | 9.8              | 14.4             | 4.0             | 7.0              | 13.0             |
| 6.0             | 10.0             | 15.0             | 6.0             | 10.0             | 16.0             |
| 6.1             | 10.2             | 15.6             | 8.0             | 13.0             | 17.0             |
| 6.3             | 10.6             | 15.8             | 9.0             | 15.0             | 18.0             |
| $\bar{y} = 6.0$ | $\bar{y} = 10.0$ | $\bar{y} = 15.0$ | $\bar{y} = 6.0$ | $\bar{y} = 10.0$ | $\bar{y} = 15.0$ |

## Observations:

- Looking only at the means, we can see that they are identical for the three populations in both the sets.
- Using the means alone, we would state that there is no difference between the two sets.

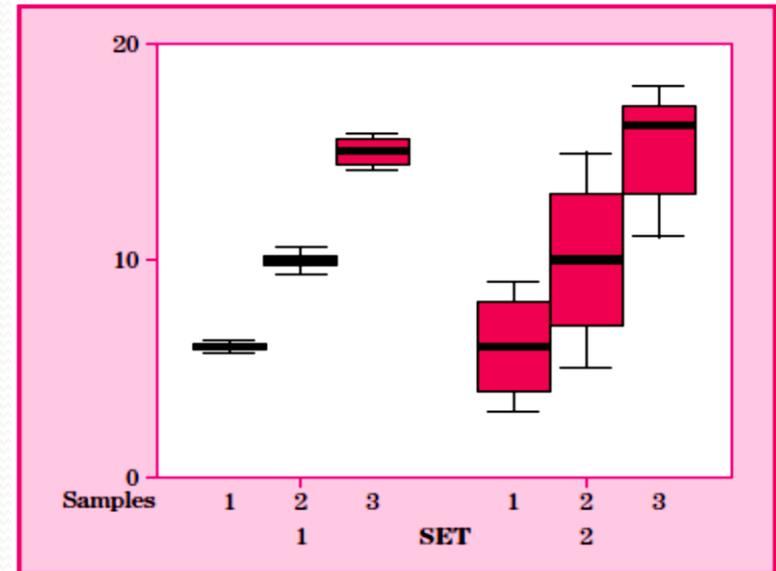
# Box plots of the two experiments

## Observation from Box plots

- It appears that there is stronger evidence of differences among means in Set 1 than among means in Set 2.
- The observations *within* the samples are more closely bunched in Set 1 than they are in Set 2,
- We know that **sample means from populations with smaller variances** will also be less variable.

### (Central Limit Theorem)

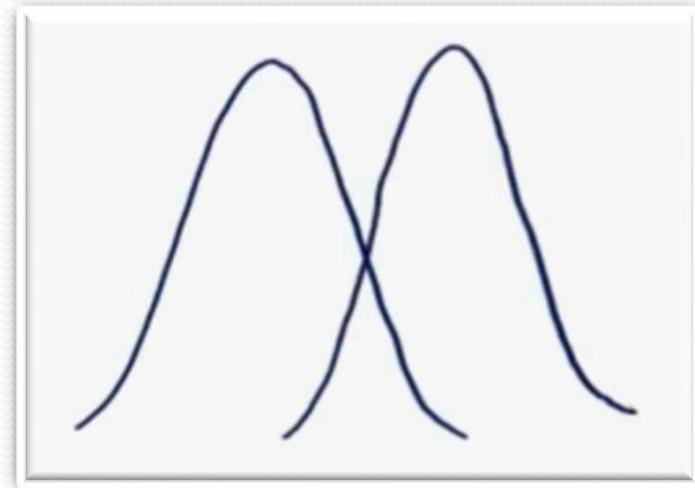
- Thus, although the variances among the means for the two sets are identical, the variance among the observations within the individual samples is smaller for Set 1 and is the reason for the apparently stronger evidence of different means.
- This observation is the basis for using the analysis of variance for making inferences about differences among means
- The analysis of variance is based on the **comparison of the variance among the means of the populations to the variance among sample observations within the individual populations.**



# Between Group Variability

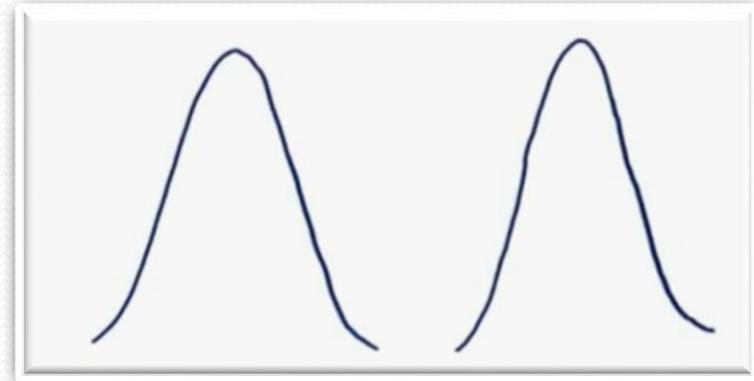
Variance among the means of the populations

- Consider the distributions of the below two samples.
- As these samples overlap, their individual means won't differ by a great margin.
- Hence, the difference between their individual means and grand mean won't be significant enough.
- Mean is a simple or arithmetic average of a range of values. There are two kinds of means that we use in ANOVA calculations, which are separate sample means ( $\mu_1$  and  $\mu_2$ ) and the grand mean  $\mu$
- The grand mean is the mean of sample means or the mean of all observations combined, irrespective of the sample.



# Between Group Variability

Now consider these two sample distributions. As the samples differ from each other by a big margin, their individual means would also differ. The difference between the individual means and grand mean would therefore also be significant.

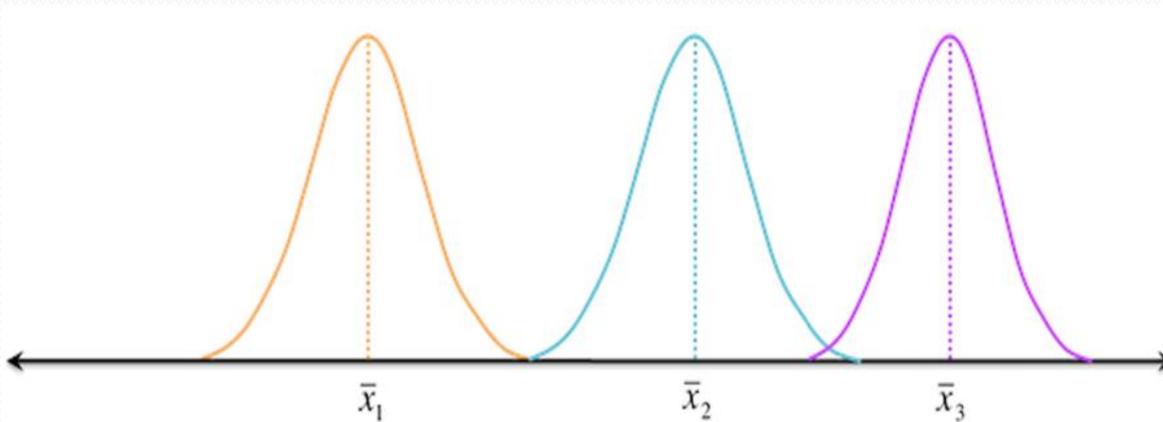
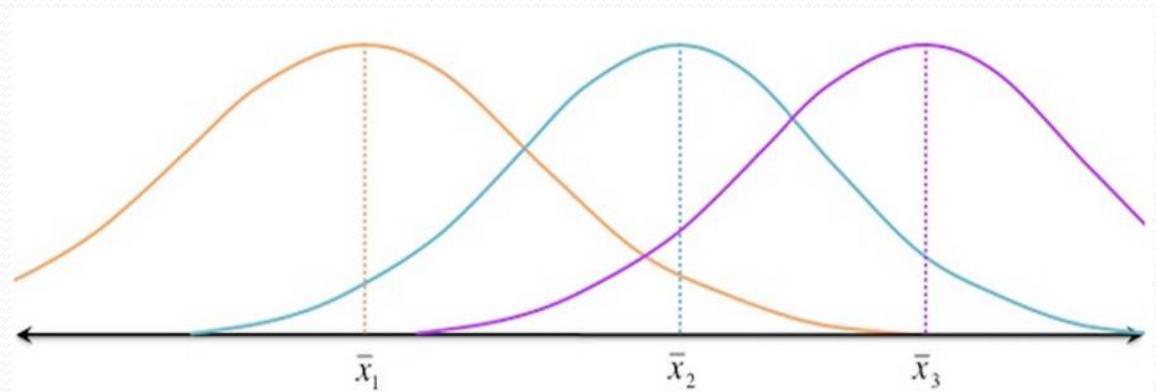


- Such variability between the distributions called *Between-group variability or variance among the means of the populations*.
- Each sample is looked at and the difference between its mean and grand mean is calculated to calculate the variability.
- If the distributions overlap or are close, the grand mean will be similar to the individual means, whereas if the distributions are far apart, difference between means and grand mean would be large.

# Within Group Variability

Variance among sample observations

- Consider the given distributions of three samples. As the spread (variability) of each sample is increased, their distributions overlap and they become part of a big population.



- Now consider another distribution of the same three samples but with less variability. Although the means of samples are similar to the samples in the above image, they seem to belong to different populations.



# How to do ANOVA?

# Some Terminologies

- Factor
  - A characteristic under consideration, thought to influence the measured observations
- Level (also called treatment)
  - A value of the factor

Typical data for a **Single-Factor** Experiment

| Level | Observations |          |     |            | Total | Mean |
|-------|--------------|----------|-----|------------|-------|------|
| 1     | $Y_{11}$     | $Y_{12}$ | ... | $Y_{1n_1}$ |       |      |
| 2     | $Y_{21}$     | $Y_{22}$ | ... | $Y_{2n_2}$ |       |      |
| ...   | ...          | ...      | ... | ...        |       |      |
| ...   | ...          | ...      | ... | ...        |       |      |
| ...   | ...          | ...      | ... | ...        |       |      |
| k     | $Y_{k1}$     | $Y_{k2}$ | ... | $Y_{kn_k}$ |       |      |

# Example 3: Single-Factor ANOVA

- Draw a straight line of between 20cm and 25 cm on a sheet of plain white card (only you know its exact length)
  - Collect 6 to 10 volunteers from each of Class VII, Class X and Class XII. Ask each volunteer to estimate independently the length of the line.
- Do differences in year means appear to outweigh differences within years?

What is/ are the Factor(s) and Levels here?

# Example 4 : Two-Factor ANOVA

- Make a list of 10 food/household items purchased regularly by your family.
- Obtain the current prices of the items in three different shops; preferably a small 'corner' shop, a small supermarket and a large supermarket or hyper market.
- Compare total shop prices.

What is/ are the Factor(s) and Levels here?

# Variants of ANOVA

Based on the number of Independent Variables and Dependent Variables considered for the study, there are different variants of ANOVA

1. **One-way ANOVA:** Only one independent variable (factor) with greater than 2 levels.
2. **Two-way ANOVA:** Two independent variables (i.e., factors).
3. **Three-way ANOVA:** Three independent variables (i.e., factors).
4. **Multivariate ANOVA:** It is used to test the significance of the effect of more independent variables.



# One-way ANOVA

# One-way ANOVA

- The purpose of the procedure is to compare sample means of  $k$  populations.
- In general, One-way ANOVA technique can be used to study the effect of  $k (> 2)$  levels of a single factor.
- To determine if different levels of the factor affect measured observations differently, the following hypotheses are tested.

$$H_0: \mu_i = \mu \quad \text{all } i = 1, 2, \dots, k$$

$$H_1: \mu_i \neq \mu \quad \text{some } i = 1, 2, \dots, k$$

That is, at least one equality is not satisfied

where  $\mu_i$  is the population mean for a level  $i$ .

# Assumptions

- When applying one-way analysis of variance, there are three key assumptions that should be satisfied as follows.
  1. The observations are obtained independently and randomly from the populations defined by the factor levels.
  2. The population at each factor level is (approximately) normally distributed.
  3. These normal populations have a common variance,  $\sigma^2$ .
- Thus, for factor level  $i$ , the population is assumed to have a distribution which is  $N(\mu_i, \sigma^2)$ .

# One-way ANOVA

| Level | Observations |          |       |          | Total    | Average        |
|-------|--------------|----------|-------|----------|----------|----------------|
| 1     | $y_{11}$     | $y_{12}$ | ..... | $y_{1n}$ | $y_{1.}$ | $\bar{y}_{1.}$ |
| 2     | $y_{21}$     | $y_{22}$ | ..... | $y_{2n}$ | $y_{2.}$ | $\bar{y}_{2.}$ |
| .     | .            | .        | ..... | .        | .        | .....          |
| .     | .            | .        | ..... | .        | .        | .....          |
| .     | .            | .        | ..... | .        | .        | .....          |
| $k$   | $y_{k1}$     | $y_{k2}$ |       | $y_{kn}$ | $y_{k.}$ | $\bar{y}_{k.}$ |
|       |              |          |       |          | $y_{..}$ | $\bar{y}_{..}$ |

An entry in the table (e.g.,  $y_{ij}$ ) represents the  $j^{th}$  observation taken under the factor at level  $i$ .

- There will be, in general,  $n$  observations under the  $i^{th}$  level.
- $y_{i.}$  represents the total of the observations under the  $i^{th}$  level.
- $\bar{y}_{i.}$  represent the average of the observation under the  $i^{th}$  level.
- $y_{..}$  represent the grand total of all the observation under the factor.
- $\bar{y}_{..}$  represent the average grand total of all the observation under the factor.

# One-way ANOVA

Expressed symbolically,

$$y_{i.} = \sum_{j=1}^{n_i} y_{ij} \quad i = 1, 2, \dots, k$$

$$\bar{y}_{i..} = \frac{y_{i.}}{n_i}$$

$$y_{..} = \sum_{i=1}^k \sum_{j=1}^{n_i} y_{ij} \quad \bar{y}_{..} = y_{..}/N$$

Here,  $N$  is the total observations, that is,  $N = n_1 + n_2 + \dots + n_k$

# Overall Variability in Data

The correlated sum of squares for each factor level

$$SS_i = \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{i.})^2 \text{ for } i = 1, 2, \dots, k$$

The corrected sum of squares for each factor level

$$SS_i = \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{..})^2$$

Alternatively, it can be prove using the computational form that

$$SS_i = \sum_{j=1}^{n_i} y_{ij}^2 - \frac{(y_{i.})^2}{n_i}$$

# Overall Variability in Data

We then calculate a pooled sum of squares

$$SS_p = \sum_{i=1}^k SS_i$$

Finally, the pooled sample of variance is

$$s_p = \frac{SS_p}{\text{pooled degree of freedom}} = \frac{SS_p}{\sum n_i - k}$$

Note that if the individual variances are available, the same can be computed as

$$s_p = \frac{\sum_{i=1}^k (n_i - 1) s_i^2}{\sum n_i - k}$$

where  $s_i^2$  are the variances for each sample. This is also called **variance within samples** and also popularly be denoted as  $\hat{\sigma}_W^2$

# Example 5: Variance within Samples

- The table below shows the lifetimes under controlled conditions, in hours in excess of 1000 hours, of samples of 60W electric light bulbs of three different brands.

| Brand |    |    |
|-------|----|----|
| 1     | 2  | 3  |
| 16    | 18 | 26 |
| 15    | 22 | 31 |
| 13    | 20 | 24 |
| 21    | 16 | 30 |
| 15    | 24 | 24 |

# Solution : Variance within Samples

- Here, there is one factor (brand) at three levels (1, 2 and 3). Also the sample sizes are all equal (to 5).
- The sample mean and variance (divisor ( $n - 1$ )) for each level are as follows.

|                | Brand |      |      |
|----------------|-------|------|------|
|                | 1     | 2    | 3    |
| Sample Size    | 5     | 5    | 5    |
| Sum            | 80    | 100  | 135  |
| Sum of squares | 1316  | 2040 | 3689 |
| Mean           | 16    | 20   | 27   |
| Variance       | 9     | 10   | 11   |

# Solution : Variance within Samples

- A pooled estimate of variance then can be calculated as follows.

$$\hat{\sigma}_W^2 = \frac{(5 - 1) \times 9 + (5 - 1) \times 10 + (5 - 1) \times 11}{5 + 5 + 5 - 3} = 10$$

- This quantity is called the **variance within samples**.
- It is an estimate of  $\sigma^2$  based on  $\nu = 5 + 5 + 5 - 3 = 12$  degrees of freedom.

# Heuristic Justification of ANOVA

- From the sampling distribution of the mean, we know that a sample mean computed from a random sample of size  $n$  from a population with mean  $\mu$  and variance  $\sigma^2$  is a random variable with mean  $\mu$  and variance  $\sigma^2/n$  [Central Limit Theorem].
- Let us see, what we can conclude in case of  $k$  ( $k > 1$ ) populations, which may have different  $\mu_i$  but have the same variance  $\sigma^2$ .

# Heuristic Justification of ANOVA

- If the null hypothesis is true, that is, each of the  $\mu_i$  has the same value, say,  $\mu$ , then the distribution of each of the  $k$  sample means,  $\bar{y}_i$ , will have mean  $\mu$  and variance  $\sigma^2/n$ .
- It then follows that, if we calculate a variance using the sample means as observations,

$$\hat{\sigma}_B^2 = \sum(\bar{y}_i - \bar{y}_{..})^2 / (k - 1)$$

- Then the quantity is an estimate of  $\sigma^2/n$ .
- Hence,  $n\hat{\sigma}_B^2$  is an estimate of  $\sigma^2$ .
  - This estimate has  $k-1$  degree of freedom and is independent of the pooled estimate of  $\sigma^2$ .

# Heuristic Justification of ANOVA

- Out of several sampling distributions, the F-distribution describes the ratio of two independent estimates of a common variance.
- The parameters of the distribution are the degrees of freedom of the numerator and denominator variances, respectively.
- If the null hypothesis of equal mean is true, then we can compute the two estimates of  $\sigma^2$  namely

$$\hat{\sigma}_B^2 = \sum(\bar{y}_{i.} - \bar{y}_{..})^2 / (k - 1) \text{ and } s_p^2, \text{ the pooled variance.}$$

- Therefore, the ratio  $\frac{n\hat{\sigma}_B^2}{s_p^2}$  has the F-distribution with degrees of freedom  $(k-1)$  and  $n - k$ .

# Heuristic Justification of ANOVA

- Thus, the procedure for testing the hypothesis.

$$H_0: \mu_i = \mu \quad \text{all } i = 1, 2, \dots, k$$

$H_1$ : at least one equality is not satisfied

- We are to reject  $H_0$ , if the calculated value of  $F = \frac{n\hat{\sigma}_B^2}{s_p^2}$  exceeds  $\alpha$  (confidence level) of the F-distributions with  $(k-1)$  and  $n - k$  degrees of freedom.

# Example 6: F-Test

| Set 1           |                  |                  | Set 2           |                  |                  |
|-----------------|------------------|------------------|-----------------|------------------|------------------|
| Sample 1        | Sample 2         | Sample 3         | Sample 1        | Sample 2         | Sample 3         |
| 5.7             | 9.4              | 14.2             | 3.0             | 5.0              | 11.0             |
| 5.9             | 9.8              | 14.4             | 4.0             | 7.0              | 13.0             |
| 6.0             | 10.0             | 15.0             | 6.0             | 10.0             | 16.0             |
| 6.1             | 10.2             | 15.6             | 8.0             | 13.0             | 17.0             |
| 6.3             | 10.6             | 15.8             | 9.0             | 15.0             | 18.0             |
| $\bar{y} = 6.0$ | $\bar{y} = 10.0$ | $\bar{y} = 15.0$ | $\bar{y} = 6.0$ | $\bar{y} = 10.0$ | $\bar{y} = 15.0$ |

- For both sets, the value of  $n\hat{\sigma}_B^2$  is 101.67. However, for Set 1,  $s_p^2 = 0.250$  while for Set 2,  $s_p^2 = 10.67$ . Thus for Set 1,  $F = 406.67$  and for Set 2,  $F = 9.53$ .
- This confirms that the relative magnitude of the two variances is the important factor for detecting difference among means.

# Example 7: Variance between Samples

- The table below shows the lifetimes under controlled conditions, in hours in excess of 1000 hours, of samples of 60W electric light bulbs of three different brands.

| Brand |    |    |
|-------|----|----|
| 1     | 2  | 3  |
| 16    | 18 | 26 |
| 15    | 22 | 31 |
| 13    | 20 | 24 |
| 21    | 16 | 30 |
| 15    | 24 | 24 |

- Assuming all lifetimes to be normally distributed with common variance, test, at the 1% significance level, the hypothesis that there is no difference between the three brands with respect to mean lifetime.

# Solution : Variance between Samples

- The variability between samples may be estimated from the three sample means as follows.

|                | Brand |      |    |
|----------------|-------|------|----|
|                | 1     | 2    | 3  |
| Sample Mean    | 16    | 20   | 27 |
| Sum            |       | 63   |    |
| Sum of squares |       | 1385 |    |
| Mean           |       | 21   |    |
| Variance       |       | 31   |    |

- This variance (divisor  $(n - 1)$ ), denoted by  $\hat{\sigma}_B^2$  is called the **variance between sample means**. Since it calculated using sample means, it is an estimate of

$$\frac{\sigma^2}{5} \text{ (that is } \frac{\sigma^2}{n} \text{ in general)}$$

based upon  $(3 - 1) = 2$  degrees of freedom, but only if the null hypothesis is true.

- If  $H_0$  is false, then the subsequent 'large' differences between the sample means **will result in  $5\hat{\sigma}_B^2$  being an inflated estimate of  $\sigma^2$** .

# Solution : F-Test

- The two estimates of  $\sigma^2$ ,  $\widehat{\sigma}_B^2$  and  $\widehat{\sigma}_W^2$ , may be tested for equality using the  $F$ -test with

$$F = \frac{5\widehat{\sigma}_B^2}{\widehat{\sigma}_W^2}$$

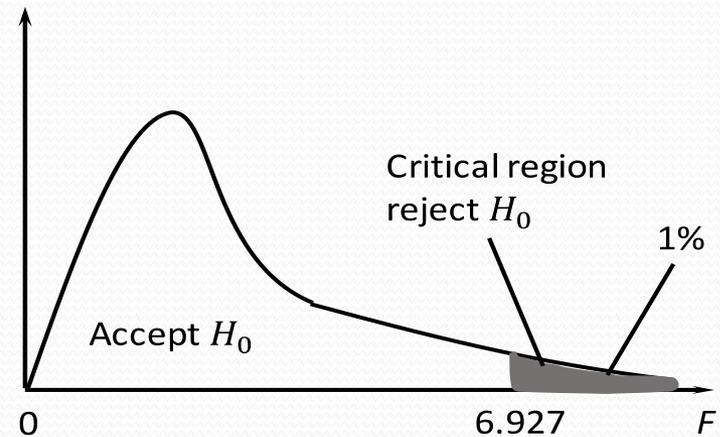
as lifetimes may be assumed to be normally distributed.

- Recall that the  $F$ -test requires the two variances to be independently distributed (from independent samples). Although this is by no means obvious here (both were calculated from the same data),  $\widehat{\sigma}_W^2$  and  $\widehat{\sigma}_B^2$  are in fact independently distributed.
- The test is always one-sided, upper-tail, since if  $H_0$  is false,  $\widehat{\sigma}_W^2$  is inflated whereas  $5\widehat{\sigma}_B^2$  is unaffected.
- **Thus in analysis of variance, the convention of placing the larger sample variance in the numerator of the F-statistic is NOT applied.**

# Solution

- The solution is thus summarized and completed as follows.

- $H_0: \mu_i = \mu$  all  $i = 1, 2, 3$
- $H_1: \mu_i \neq \mu$  some  $i = 1, 2, 3$
- Significance level,  $\alpha = 0.01$
- Degrees of freedom,  $\nu_1 = 2, \nu_2 = 12$
- Critical region is  $F > 6.927$



- Test statistic is  $F = \frac{5\hat{\sigma}_B^2}{\hat{\sigma}_W^2} = \frac{155}{10} = 15.5$

- This value does lie in the critical region. There is evidence, at the 1% significance level, that **the true mean lifetimes of the three brands of bulb do differ.**

# Notation and computational formulae

- In essence, given a population a single factor of  $k$  levels, we have to calculate two estimations for  $\sigma^2$ .

- Sampling variance between groups with  $(k-1)$  degree of freedom

$$n\hat{\sigma}_B^2 = n \sum (\bar{y}_{i.} - \bar{y}_{..})^2 / (k - 1).$$

- Sampling variance within groups with  $(n-k)$  degree of freedom

$$\hat{\sigma}_W^2 = \frac{\sum_{i=1}^k SS_i}{\sum n_i - k}$$

# Notation and computational formulae

- The calculations undertaken in the previous example are somewhat cumbersome, and are prone to inaccuracy with non-integer sample means. They also require considerable changes when the sample sizes are unequal. Equivalent computational formulae are available which cater for both equal and unequal sample sizes.
- First, some notation.

|                                             |                                           |
|---------------------------------------------|-------------------------------------------|
| Number of samples (or levels)               | $= k$                                     |
| Number of observations in $i$ th sample     | $= n_i, \quad i = 1, 2, \dots, k$         |
| Total number of observations                | $= n = \sum_i n_i$                        |
| $j$ – th observation in $i$ -th sample      | $= y_{ij}, \quad j = 1, 2, \dots, n_i$    |
| Sum of $n_i$ observations in $i$ –th sample | $= T_i = \sum_j y_{ij}$                   |
| Sum of all $n$ observations                 | $= T = \sum_i T_i = \sum_i \sum_j y_{ij}$ |

# Notation and computational formulae

- The computational formulae now follow.

|                                 |                                                   |
|---------------------------------|---------------------------------------------------|
| Total sum of squares,           | $SS_T = \sum_i \sum_j x_{ij}^2 - \frac{T^2}{n}$   |
| Between samples sum of squares, | $SS_B = \sum_i \frac{T_i^2}{n_i} - \frac{T^2}{n}$ |
| Within samples sum of squares,  | $SS_W = SS_T - SS_B$                              |

- A mean square (or unbiased variance estimate) is given by

(sum of squares)  $\div$  (degrees of freedom)

e.g. 
$$\hat{\sigma}^2 = \frac{(x - \bar{x})^2}{n-1}$$

Hence

|                              |                           |
|------------------------------|---------------------------|
| Total mean square,           | $MS_T = \frac{SS_T}{n-1}$ |
| Between samples mean square, | $MS_B = \frac{SS_B}{k-1}$ |
| Within samples mean square,  | $MS_W = \frac{SS_W}{n-k}$ |

- Note that for the degrees of freedom:  $(k-1) + (n-k) = (n-1)$**

# Example 8: F-Test using Formula

- For the previous example on 60W electric light bulbs, use these computational formulae to show the following.

(a)  $SS_T = 430$

(b)  $SS_B = 310$

(c)  $MS_B = 155 (5\hat{\sigma}_B^2)$

(d)  $MS_W = 10 (\hat{\sigma}_W^2)$

- Note that  $F = \frac{MS_B}{MS_W} = \frac{155}{10} = 15.5$  as previously.

# ANOVA Table

- It is convenient to summarize the results of an analysis of variance in a table. For a one factor analysis this takes the following form.

| Source of variation | Sum of squares | Degrees of freedom | Mean square | F ratio             |
|---------------------|----------------|--------------------|-------------|---------------------|
| Between samples     | $SS_B$         | $k - 1$            | $MS_B$      | $\frac{MS_B}{MS_W}$ |
| Within samples      | $SS_W$         | $n - k$            | $MS_W$      |                     |
| <b>Total</b>        | $SS_T$         | $n - 1$            |             |                     |

# Example 9: F-Test for unbalanced

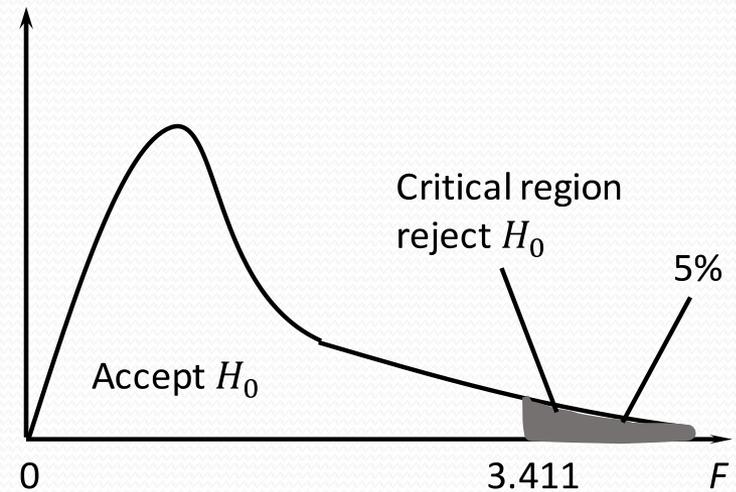
- In a comparison of the cleaning action of four detergents, 20 pieces of white cloth were first soiled with India ink. The cloths were then washed under controlled conditions with 5 pieces washed by each of the detergents. Unfortunately three pieces of cloth were 'lost' in the course of the experiment. Whiteness readings, made on the 17 remaining pieces of cloth, are shown below.

| Detergent |    |    |    |
|-----------|----|----|----|
| A         | B  | C  | D  |
| 77        | 74 | 73 | 76 |
| 81        | 66 | 78 | 85 |
| 61        | 58 | 57 | 77 |
| 76        |    | 69 | 64 |
| 69        |    | 63 |    |

- Assuming all whiteness readings to be normally distributed with common variance, test the hypothesis of no difference between the four brands as regards mean whiteness readings after washing.

# Solution

- $H_0: \mu_i = \mu$  all  $i = 1, 2, 3$
- $H_1: \mu_i \neq \mu$  some  $i = 1, 2, 3$
- Significance level,  $\alpha = 0.05$  (say)
- Degrees of freedom,  $v_1 = k - 1 = 3$ ,  
and  $v_2 = n - k = 17 - 4 = 13$
- Critical region is  $F > 3.411$



# Solution

|       | A   | B   | C   | D   | Total      |
|-------|-----|-----|-----|-----|------------|
| $n_i$ | 5   | 3   | 5   | 4   | $17 = n$   |
| $T_i$ | 364 | 198 | 340 | 302 | $1204 = T$ |

$$\sum_i \sum_j y_{ij}^2 = 86362$$

$$SS_T = 86362 - \frac{1204^2}{17} = 1090.47$$

$$SS_B = \left( \frac{364^2}{5} + \frac{198^2}{3} + \frac{340^2}{5} + \frac{302^2}{4} \right) - \frac{1204^2}{17} = 216.67$$

$$SS_W = 1090.47 - 216.67 = 873.80$$

# Solution

- The ANOVA table is now as follows.

| Source of variation | Sum of squares | Degrees of freedom | Mean square | F ratio |
|---------------------|----------------|--------------------|-------------|---------|
| Between detergents  | 216.67         | 3                  | 72.22       | 1.07    |
| Within detergents   | 873.80         | 13                 | 67.22       |         |
| <b>Total</b>        | 1090.47        | 16                 |             |         |

- The F ratio of 1.07 does not lie in the critical region.
- Thus there is no evidence, at the 5% significance level, to suggest a difference between the four brands as regards mean whiteness after washing.



# Two-way ANOVA

# Two way (factor) ANOVA

- This is an extension of the one factor situation to take account of a second factor.
- The levels of this second factor are often determined by groupings of subjects or units used in the investigation. As such it is often called a blocking factor because it places subjects or units into homogeneous groups called blocks. The design itself is then called a randomised block design.

# Example 10: Two-factor Analysis

- A computer manufacturer wishes to compare the speed of four of the firm's compilers. The manufacturer can use one of two experimental designs.
  - a) Use 20 similar programs, randomly allocating 5 programs to each compiler.
  - b) Use 4 copies of any 5 programs, allocating 1 copy of each program to each compiler.
- Which of (a) and (b) would you recommend, and why?

# Solution

- In (a), although the 20 programs are similar, any differences between them may affect the compilation times and hence perhaps any conclusions. Thus in the 'worst scenario', the 5 programs allocated to what is really the fastest compiler could be the 5 requiring the longest compilation times, resulting in the compiler appearing to be the slowest! If used, the results would require a one factor analysis of variance; the factor being compiler at 4 levels.
- In (b), since all 5 programs are run on each compiler, differences between programs should not affect the results. Indeed it may be advantageous to use 5 programs that differ markedly so that comparisons of compilation times are more general. For this design, there are two factors; compiler (4 levels) and program (5 levels). The factor of principal interest is compiler whereas the other factor, program, may be considered as a blocking factor as it creates 5 blocks each containing 4 copies of the same program.
- Thus (b) is the better designed investigation.

# Solution

- The actual compilation times, in milliseconds, for this two factor (randomised block) design are shown in the following table.

|           | Compiler |       |       |       |
|-----------|----------|-------|-------|-------|
|           | 1        | 2     | 3     | 4     |
| Program A | 29.21    | 28.25 | 28.20 | 28.62 |
| Program B | 26.18    | 26.02 | 26.22 | 25.56 |
| Program C | 30.91    | 30.18 | 30.52 | 30.09 |
| Program D | 25.14    | 25.26 | 25.20 | 25.02 |
| Program E | 26.16    | 25.16 | 25.26 | 25.46 |

# Assumptions and Interaction

- The three assumptions for a two factor analysis of variance when there is only one observed measurement at each combination of levels of the two factors are as follows.
  1. The population at each factor level combination is (approximately) normally distributed.
  2. These normal populations have a common variance,  $\sigma^2$ .
  3. The effect of one factor is the same at all levels of the other factor.
- Hence from assumptions 1 and 2, when one factor is at level  $i$  and the other at level  $j$ , the population has a distribution which is
$$N(\mu_{ij}, \sigma^2)$$
- Assumption 3 is equivalent to stating that there is no interaction between the two factors.

# Assumptions and Interaction

- Now interaction exists when the effect of one factor depends upon the level of the other factor. For example consider the effects of the two factors: sugar (levels none and 2 teaspoons), and stirring (levels none and 1 minute), *on the sweetness of a cup of tea.*
- Stirring has no effect on sweetness if sugar is not added but certainly does have an effect if sugar is added. Similarly, adding sugar has little effect on sweetness unless the tea is stirred.
- Hence factors sugar and stirring are said to interact.
- Interaction can only be assessed if more than one measurement is taken at each combination of the factor levels. Since such situations are beyond the scope of this text, it will always be assumed that interaction between the two factors does not exist.

# Assumptions and Interaction

- Thus, for example, since it would be most unusual to find one compiler particularly suited to one program, the assumption of no interaction between compilers and programs appears reasonable.

# Notation and Computational Formulae

- As illustrated earlier, the data for a two-way ANOVA can be displayed in a two-way table. It is thus convenient, in general, to label the factors as a **row factor** and a **column factor**.
- Notation, similar to that for the one factor case, is then as follows.

|                                                             |                                    |
|-------------------------------------------------------------|------------------------------------|
| Number of levels of row factor                              | = $r$                              |
| Number of levels of column factor                           | = $c$                              |
| Total number of observations                                | = $rc$                             |
| Observation in (i j-th cell of table                        | = $x_{ij}$                         |
| (ith level of row factor and<br>jth level of column factor) | $i=1,2,\dots,r$<br>$j=1,2,\dots,c$ |

# Notation and computational formulae

Sum of  $c$  observations in  $i$ -th row

$$= T_{Ri} = \sum_j x_{ij}$$

Sum of  $r$  observations in  $j$ -th column

$$= T_{Cj} = \sum_i x_{ij}$$

Sum of all  $rc$  observations

$$= T = \sum_i \sum_j x_{ij} = \sum_i T_{Ri} = \sum_j T_{Cj}$$

- These lead to the following computational formulae which again are similar to those for one-way ANOVA except that there is an additional sum of squares, etc. for the second factor.

# Notation and computational formulae

Total sum of squares,

$$SS_T = \sum_i \sum_j x_{ij}^2 - \frac{T^2}{rc}$$

Between rows sum of squares,

$$SS_R = \sum_i \frac{T_{Ri}^2}{c} - \frac{T^2}{rc}$$

Between columns sum of squares,

$$SS_C = \sum_j \frac{T_{Cj}^2}{r} - \frac{T^2}{rc}$$

Error (residual) sum of squares,

$$SS_E = SS_T - SS_R - SS_C$$

What are the degrees of freedom for SST , SSR and SSC when there are 20 observations in a table of 5 rows and 4 columns?  
What is the degrees of freedom of SSE ?

# ANOVA Table and Hypothesis Test

For a two factor analysis of variance this takes the following form.

| Source of variation | Sum of squares | Degrees of freedom | Mean square | F ratio             |
|---------------------|----------------|--------------------|-------------|---------------------|
| Between rows        | $SS_R$         | $r - 1$            | $MS_R$      | $\frac{MS_R}{MS_C}$ |
| Between columns     | $SS_C$         | $c - 1$            | $MS_C$      | $\frac{MS_C}{MS_E}$ |
| Error (residual)    | $SS_E$         | $(r-1)(c-1)$       | $MS_E$      |                     |
| Total               | $SS_T$         | $rc - 1$           |             |                     |

- Notes :

1. The three sums of squares,  $SS_R$ ,  $SS_C$  and  $SS_E$  are independently distributed.
2. For the degrees of freedom:  $(r-1) + (c-1) + (r-1)(c-1) = rc - 1$

# ANOVA Table and Hypothesis Test

- Using the F ratios, tests for significant row effects and for significant column effects can be undertaken.

| Ho: no effect due to row factor                           | Ho: no effect due to column factor                        |
|-----------------------------------------------------------|-----------------------------------------------------------|
| H <sub>1</sub> : an effect due to row factor              | H <sub>1</sub> : an effect due to column factor           |
| Critical region,<br>$F > F^{\alpha}_{[(r-1),(r-1)(c-1)]}$ | Critical region,<br>$F > F^{\alpha}_{[(c-1),(r-1)(c-1)]}$ |
| Test statistic,<br>$F_r = \frac{MS_R}{MS_E}$              | Test statistic,<br>$F_r = \frac{MS_C}{MS_E}$              |

# Example 11: Two-way ANOVA

- Returning to the compilation times, in milliseconds, for each of five programs, run on four compilers.
- Test, at the 1% significance level, the hypothesis that there is no difference between the performance of the four compilers.
- Has the use of programs as a blocking factor proved worthwhile? Explain.
- The data, given earlier, are reproduced below.

|           | Compiler |       |       |       |
|-----------|----------|-------|-------|-------|
|           | 1        | 2     | 3     | 4     |
| Program A | 29.21    | 28.25 | 28.20 | 28.62 |
| Program B | 26.18    | 26.02 | 26.22 | 25.56 |
| Program C | 30.91    | 30.18 | 30.52 | 30.09 |
| Program D | 25.14    | 25.26 | 25.20 | 25.02 |
| Program E | 26.16    | 25.16 | 25.26 | 25.46 |

# Solution : Dataset

- To ease computations, these data have been transformed (coded) by

$$x = 100 \times (\text{time} - 25)$$

to give the following table of values and totals.

|                                             | Compiler                  |            |             |            | Row(totals) ( $T_{R_i}$ ) |
|---------------------------------------------|---------------------------|------------|-------------|------------|---------------------------|
|                                             | 1                         | 2          | 3           | 4          |                           |
| Program A                                   | 421                       | 325        | 320         | 362        | 1428                      |
| Program B                                   | 118                       | 102        | 122         | 56         | 398                       |
| Program C                                   | 591                       | 518        | 552         | 509        | 2170                      |
| Program D                                   | 14                        | 26         | 20          | 2          | 62                        |
| Program E                                   | 116                       | 14         | 26          | 46         | 202                       |
| <b>Column totals (<math>T_{C_j}</math>)</b> | <b>1260</b>               | <b>985</b> | <b>1040</b> | <b>975</b> | <b>4260 = T</b>           |
|                                             | $\sum x_{ij}^2 = 1757768$ |            |             |            |                           |

# Solution : Parameters

- The sums of squares are now calculated as follows.  
(Rows = Programs, Columns = Compilers)

- $SS_T = 1757768 = \frac{4260^2}{20} = 850388$

- $SS_R = \frac{1}{4} (1428^2 + 398^2 + 2170^2 + 62^2 + 202^2) - \frac{4260^2}{20} = 830404$

- $SS_C = \frac{1}{5} (1260^2 + 985^2 + 1040^2 + 975^2) - \frac{4260^2}{20} = 10630$

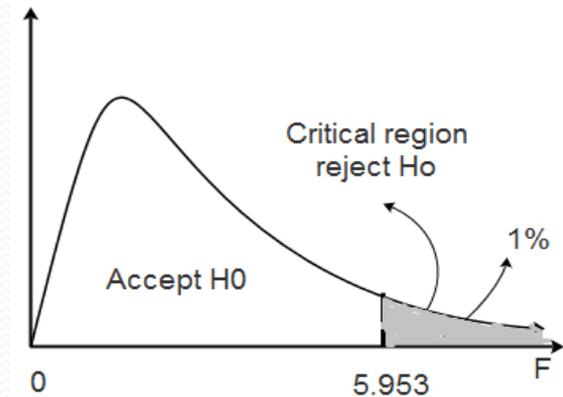
- $SS_E = 850388 - 830404 - 10630 = 9354$

# Solution: ANOVA Table

| Source of variation | Sum of squares | Degrees of freedom | Mean square | F ratio |
|---------------------|----------------|--------------------|-------------|---------|
| Between programs    | 830404         | 4                  | 207601.0    | 266.33  |
| Between compilers   | 10630          | 3                  | 3543.3      | 4.55    |
| Error (residual)    | 9354           | 12                 | 779.5       |         |
| Total               | 850388         | 19                 |             |         |

# Solution : Hypothesis Test

- $H_0$ : no effect on compilation times due to compilers
- $H_1$ : an effect on compilation times due to compilers
- Significance level,  $\alpha = 0.001$
- Degrees of freedom,  $v_1 = c - 1 = 3$   
and  $v_2 = (r - 1)(c - 1) = 4 \times 3 = 12$
- Critical region is  $F > 5.953$
- Test statistic  $FC = 4.55$



- This value does not lie in the critical region. Thus there is no evidence, at the 1% significance level, to suggest a difference in compilation times between the four compilers.

# Reference

- Design and Analysis of Experiments (8<sup>th</sup> Edition), Douglas C. Montgomery, John Wiley & Sons, 2013.



# Any question?

You may also send your question(s) at [tanujitisi@gmail.com](mailto:tanujitisi@gmail.com)