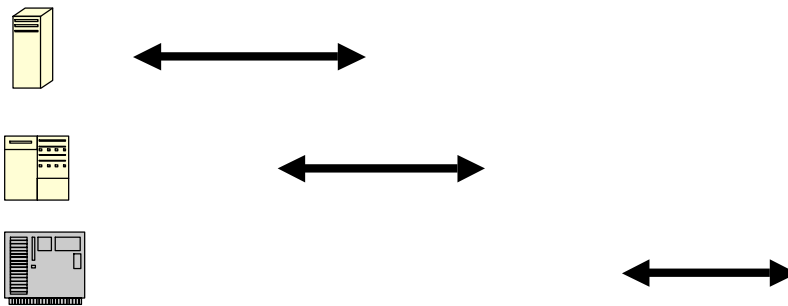# ANOVA- Analyisis of Variance

## CS 700

---

# Comparing alternatives

❏ Comparing two alternatives
  ➢ use confidence intervals
❏ Comparing more than two alternatives
  ➢ ANOVA
    • Analysis of Variance

## Comparing More Than Two Alternatives

❑ Naïve approach
  ➢ Compare confidence intervals

## One-Factor Analysis of Variance (ANOVA)

❑ Very general technique
  ➢ Look at total *variation* in a set of measurements
  ➢ Divide into meaningful components
❑ Also called
  ➢ One-way classification
  ➢ One-factor experimental design
❑ Introduce basic concept with one-factor ANOVA
❑ Generalize later with *design of experiments*

## One-Factor Analysis of Variance (ANOVA)

❑ Separates total variation observed in a set of measurements into:
   1. Variation within one system
      - Due to random measurement errors
   2. Variation between systems
      - Due to real differences + random error

❑ Is variation(2) statistically > variation(1)?

5

## ANOVA

❑ Make $n$ measurements of $k$ alternatives
❑ $y_{ij}$ = $i$th measurment on $j$th alternative
❑ Assumes errors are:
   - Independent
   - Gaussian (normal)

6

3

## Measurements for All Alternatives

| Measurements | Alternatives | | | | | |
|---|---|---|---|---|---|---|
| | 1 | 2 | ... | $j$ | ... | $k$ |
| 1 | $y_{11}$ | $y_{12}$ | ... | $y_{1j}$ | ... | $y_{k1}$ |
| 2 | $y_{21}$ | $y_{22}$ | ... | $y_{2j}$ | ... | $y_{2k}$ |
| ... | ... | ... | ... | ... | ... | ... |
| $i$ | $y_{i1}$ | $y_{i2}$ | ... | $y_{ij}$ | ... | $y_{ik}$ |
| ... | ... | ... | ... | ... | ... | ... |
| $n$ | $y_{n1}$ | $y_{n2}$ | ... | $y_{nj}$ | ... | $y_{nk}$ |
| Col mean | $y_{.1}$ | $y_{.2}$ | ... | $y_{.j}$ | ... | $y_{.k}$ |
| Effect | $\alpha_1$ | $\alpha_2$ | ... | $\alpha_j$ | ... | $\alpha_k$ |

## Column Means

❑ Column means are average values of all measurements within a single alternative
  ➢ Average performance of one alternative

$$\bar{y}_{.j} = \frac{\sum_{i=1}^{n} y_{ij}}{n}$$

# Column Means

| Measurem ents | Alternatives | | | | | |
|---|---|---|---|---|---|---|
| | 1 | 2 | ... | $j$ | ... | $k$ |
| 1 | $y_{11}$ | $y_{12}$ | ... | $y_{1j}$ | ... | $y_{k1}$ |
| 2 | $y_{21}$ | $y_{22}$ | ... | $y_{2j}$ | ... | $y_{2k}$ |
| ... | ... | ... | ... | ... | ... | ... |
| $i$ | $y_{i1}$ | $y_{i2}$ | ... | $y_{ij}$ | ... | $y_{ik}$ |
| ... | ... | ... | ... | ... | ... | ... |
| $n$ | $y_{n1}$ | $y_{n2}$ | ... | $y_{nj}$ | ... | $y_{nk}$ |
| Col mean | $y_{.1}$ | $y_{.2}$ | ... | $y_{.j}$ | ... | $y_{.k}$ |
| Effect | $\alpha_1$ | $\alpha_2$ | ... | $\alpha_j$ | ... | $\alpha_k$ |

9

# Deviation From Column Mean

$$y_{ij} = \bar{y}_{.j} + e_{ij}$$

$$e_{ij} = \text{deviation of } y_{ij} \text{ from column mean}$$

$$= \text{error in measurements}$$

10

## Error = Deviation From Column Mean

| Measurements | Alternatives | | | | | |
|---|---|---|---|---|---|---|
| | 1 | 2 | … | $j$ | … | $k$ |
| 1 | $y_{11}$ | $y_{12}$ | … | $y_{1j}$ | … | $y_{k1}$ |
| 2 | $y_{21}$ | $y_{22}$ | … | $y_{2j}$ | … | $y_{2k}$ |
| … | … | … | … | … | … | … |
| $i$ | $y_{i1}$ | $y_{i2}$ | … | $y_{ij}$ | … | $y_{ik}$ |
| … | … | … | … | … | … | … |
| $n$ | $y_{n1}$ | $y_{n2}$ | … | $y_{nj}$ | … | $y_{nk}$ |
| Col mean | $y_{.1}$ | $y_{.2}$ | … | $y_{.j}$ | … | $y_{.k}$ |
| Effect | $\alpha_1$ | $\alpha_2$ | … | $\alpha_j$ | … | $\alpha_k$ |

11

## Overall Mean

❑ Average of all measurements made of all alternatives

$$\overline{y}_{..} = \frac{\sum_{j=1}^{k} \sum_{i=1}^{n} y_{ij}}{kn}$$

12

## Overall Mean

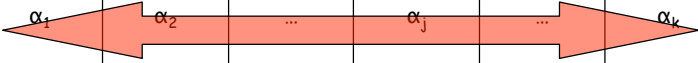| Measurements | Alternatives | | | | | |
|---|---|---|---|---|---|---|
| | 1 | 2 | ... | j | ... | k |
| 1 | $y_{11}$ | $y_{12}$ | ... | $y_{1j}$ | ... | $y_{k1}$ |
| 2 | $y_{21}$ | $y_{22}$ | ... | $y_{2j}$ | ... | $y_{2k}$ |
| ... | ... | ... | ... | ... | ... | ... |
| i | $y_{i1}$ | $y_{i2}$ | ... | $y_{ij}$ | ... | $y_{ik}$ |
| ... | ... | ... | ... | ... | ... | ... |
| n | $y_{n1}$ | $y_{n2}$ | ... | $y_{nj}$ | ... | $y_{nk}$ |
| Col mean | $y_{.1}$ | $y_{.2}$ | ... | $y_{.j}$ | ... | $y_{.k}$ |
| Effect | $\alpha_1$ | $\alpha_2$ | ... | $\alpha_j$ | ... | $\alpha_k$ |

13

## Deviation From Overall Mean

$$\bar{y}_{.j} = \bar{y}_{..} + \alpha_j$$

$\alpha_j$ = deviation of column mean from overall mean

= effect of alternative $j$

14

## Effect = Deviation From Overall Mean

| Measurements | Alternatives | | | | | |
|---|---|---|---|---|---|---|
| | 1 | 2 | ... | j | ... | k |
| 1 | $y_{11}$ | $y_{12}$ | ... | $y_{1j}$ | ... | $y_{k1}$ |
| 2 | $y_{21}$ | $y_{22}$ | ... | $y_{2j}$ | ... | $y_{2k}$ |
| ... | ... | ... | ... | ... | ... | ... |
| i | $y_{i1}$ | $y_{i2}$ | ... | $y_{ij}$ | ... | $y_{ik}$ |
| ... | ... | ... | ... | ... | ... | ... |
| n | $y_{n1}$ | $y_{n2}$ | ... | $y_{nj}$ | ... | $y_{nk}$ |
| Col mean | $y_{.1}$ | $y_{.2}$ | ... | $y_{.j}$ | ... | $y_{.k}$ |
| Effect | $\alpha_1$ | $\alpha_2$ | ... | $\alpha_j$ | ... | $\alpha_k$ |

## Effects and Errors

- ❏ *Effect* is distance from overall mean
  - ➢ Horizontally across alternatives
- ❏ *Error* is distance from column mean
  - ➢ Vertically within one alternative
  - ➢ Error across alternatives, too
- ❏ Individual measurements are then:

$$y_{ij} = \overline{y}_{..} + \alpha_j + e_{ij}$$

## Sum of Squares of Differences:  SSE

$$y_{ij} = \bar{y}_{.j} + e_{ij}$$

$$e_{ij} = y_{ij} - \bar{y}_{.j}$$

$$SSE = \sum_{j=1}^{k} \sum_{i=1}^{n} \left(e_{ij}\right)^2 = \sum_{j=1}^{k} \sum_{i=1}^{n} \left(y_{ij} - \bar{y}_{.j}\right)^2$$

17

## Sum of Squares of Differences:  SSA

$$\bar{y}_{.j} = \bar{y}_{..} + \alpha_j$$

$$\alpha_j = \bar{y}_{.j} - \bar{y}_{..}$$

$$SSA = n \sum_{j=1}^{k} \left(\alpha_j\right)^2 = n \sum_{j=1}^{k} \left(\bar{y}_{.j} - \bar{y}_{..}\right)^2$$

18

## Sum of Squares of Differences:  SST

$$y_{ij} = \bar{y}_{..} + \alpha_j + e_{ij}$$

$$t_{ij} = \alpha_j + e_{ij} = y_{ij} - \bar{y}_{..}$$

$$SST = \sum_{j=1}^{k} \sum_{i=1}^{n} \left( t_{ij} \right)^2 = \sum_{j=1}^{k} \sum_{i=1}^{n} \left( y_{ij} - \bar{y}_{..} \right)^2$$

19

## Sum of Squares of Differences

$$SSA = n \sum_{j=1}^{k} \left( \bar{y}_{.j} - \bar{y}_{..} \right)^2$$

$$SSE = \sum_{j=1}^{k} \sum_{i=1}^{n} \left( y_{ij} - \bar{y}_{.j} \right)^2$$

$$SST = \sum_{j=1}^{k} \sum_{i=1}^{n} \left( y_{ij} - \bar{y}_{..} \right)^2$$

20

## Sum of Squares of Differences

- **SST** = differences between each measurement and overall mean
- **SSA** = variation due to effects of alternatives
- **SSE** = variation due to errors in measurments

$$SST = SSA + SSE$$

## ANOVA – Fundamental Idea

- Separates variation in measured values into:
  1. Variation due to effects of alternatives
     SSA – variation across columns
  2. Variation due to errors
     SSE – variation within a single column
- If differences among alternatives are due to real differences, SSA should be statistically > SSE

## Comparing SSE and SSA

❑ Simple approach
  ➤ *SSA* / *SST* = fraction of total variation explained by differences among alternatives
  ➤ *SSE* / *SST* = fraction of total variation due to experimental error
❑ But is it statistically significant?

23

## Statistically Comparing SSE and SSA

$$\text{Variance} = \text{mean square value}$$

$$= \frac{\text{total variation}}{\text{degrees of freedom}}$$

$$s_x^2 = \frac{SSx}{df}$$

24

# Degrees of Freedom

- $df(SSA) = k - 1$, since $k$ alternatives
- $df(SSE) = k(n - 1)$, since $k$ alternatives, each with $(n - 1)$ $df$
- $df(SST) = df(SSA) + df(SSE) = kn - 1$

25

# Degrees of Freedom for Effects

| Measurements | Alternatives | | | | | |
|---|---|---|---|---|---|---|
| | 1 | 2 | … | $j$ | … | $k$ |
| 1 | $y_{11}$ | $y_{12}$ | … | $y_{1j}$ | … | $y_{k1}$ |
| 2 | $y_{21}$ | $y_{22}$ | … | $y_{2j}$ | … | $y_{2k}$ |
| … | … | … | … | … | … | … |
| $i$ | $y_{i1}$ | $y_{i2}$ | … | $y_{ij}$ | … | $y_{ik}$ |
| … | … | … | … | … | … | … |
| $n$ | $y_{n1}$ | $y_{n2}$ | … | $y_{nj}$ | … | $y_{nk}$ |
| Col mean | $y_{.1}$ | $y_{.2}$ | … | $y_{.j}$ | … | $y_{.k}$ |
| Effect | $\alpha_1$ | $\alpha_2$ | … | $\alpha_j$ | … | $\alpha_k$ |

26

13

## Degrees of Freedom for Errors

| Measurements | Alternatives | | | | | |
|---|---|---|---|---|---|---|
| | 1 | 2 | ... | $j$ | ... | $k$ |
| 1 | $y_{11}$ | $y_{12}$ | ... | $y_{1j}$ | ... | $y_{k1}$ |
| 2 | $y_{21}$ | $y_{22}$ | ... | $y_{2j}$ | ... | $y_{2k}$ |
| ... | ... | ... | ... | ... | ... | ... |
| $i$ | $y_{i1}$ | $y_{i2}$ | ... | $y_{ij}$ | ... | $y_{ik}$ |
| ... | ... | ... | ... | ... | ... | ... |
| $n$ | $y_{n1}$ | $y_{n2}$ | ... | $y_{nj}$ | ... | $y_{nk}$ |
| Col mean | $y_{.1}$ | $y_{.2}$ | ... | $y_{.j}$ | ... | $y_{.k}$ |
| Effect | $\alpha_1$ | $\alpha_2$ | ... | $\alpha_j$ | ... | $\alpha_k$ |

27

## Degrees of Freedom for Errors

| Measurements | Alternatives | | | | | |
|---|---|---|---|---|---|---|
| | 1 | 2 | ... | $j$ | ... | $k$ |
| 1 | $y_{11}$ | $y_{12}$ | ... | $y_{1j}$ | ... | $y_{k1}$ |
| 2 | $y_{21}$ | $y_{22}$ | ... | $y_{2j}$ | ... | $y_{2k}$ |
| ... | ... | ... | ... | ... | ... | ... |
| $i$ | $y_{i1}$ | $y_{i2}$ | ... | $y_{ij}$ | ... | $y_{ik}$ |
| ... | ... | ... | ... | ... | ... | ... |
| $n$ | $y_{n1}$ | $y_{n2}$ | ... | $y_{nj}$ | ... | $y_{nk}$ |
| Col mean | $y_{.1}$ | $y_{.2}$ | ... | $y_{.j}$ | ... | $y_{.k}$ |
| Effect | $\alpha_1$ | $\alpha_2$ | ... | $\alpha_j$ | ... | $\alpha_k$ |

28

## Variances from Sum of Squares (Mean Square Value)

$$s_a^2 = \frac{SSA}{k-1}$$

$$s_e^2 = \frac{SSE}{k(n-1)}$$

29

## Comparing Variances

❑ Use F-test to compare ratio of variances

$$F = \frac{s_a^2}{s_e^2}$$

$$F_{[1-\alpha;\,df(num),\,df(denom)]} = \text{tabulated critical values}$$

30

## F-test

- If $F_{computed} > F_{table}$
  $\rightarrow$ We have $(1 - \alpha) * 100\%$ confidence that variation due to actual differences in alternatives, SSA, is statistically greater than variation due to errors, SSE.

## ANOVA Summary

| Variation | Alternatives | Error | Total |
|---|---|---|---|
| Sum of squares | $SSA$ | $SSE$ | $SST$ |
| Deg freedom | $k-1$ | $k(n-1)$ | $kn-1$ |
| Mean square | $s_a^2 = SSA/(k-1)$ | $s_e^2 = SSE/[k(n-1)]$ | |
| Computed $F$ | $s_a^2/s_e^2$ | | |
| Tabulated $F$ | $F_{[1-\alpha;(k-1),k(n-1)]}$ | | |

## ANOVA Example

| Measurements | Alternatives | | | Overall mean |
| --- | --- | --- | --- | --- |
| | 1 | 2 | 3 | |
| 1 | 0.0972 | 0.1382 | 0.7966 | |
| 2 | 0.0971 | 0.1432 | 0.5300 | |
| 3 | 0.0969 | 0.1382 | 0.5152 | |
| 4 | 0.1954 | 0.1730 | 0.6675 | |
| 5 | 0.0974 | 0.1383 | 0.5298 | |
| Column mean | 0.1168 | 0.1462 | 0.6078 | 0.2903 |
| Effects | -0.1735 | -0.1441 | 0.3175 | |

## ANOVA Example

| Variation | Alternatives | Error | Total |
| --- | --- | --- | --- |
| Sum of squares | $SSA = 0.7585$ | $SSE = 0.0685$ | $SST = 0.8270$ |
| Deg freedom | $k - 1 = 2$ | $k(n-1) = 12$ | $kn - 1 = 14$ |
| Mean square | $s_a^2 = 0.3793$ | $s_e^2 = 0.0057$ | |
| Computed $F$ | $0.3793/0.0057 = 66.4$ | | |
| Tabulated $F$ | $F_{[0.95;2,12]} = 3.89$ | | |

# Conclusions from example

- SSA/SST = 0.7585/0.8270 = 0.917
  - → 91.7% of total variation in measurements is due to differences among alternatives
- SSE/SST = 0.0685/0.8270 = 0.083
  - → 8.3% of total variation in measurements is due to noise in measurements
- Computed *F* statistic > tabulated *F* statistic
  - → 95% confidence that differences among alternatives are statistically significant.

# Contrasts

- ANOVA tells us that there is a statistically significant difference among alternatives
- But it does *not* tell us *where* difference is
- Use method of contrasts to compare subsets of alternatives
  - ➤ A vs B
  - ➤ {A, B} vs {C}
  - ➤ Etc.

## Contrasts

❑ Contrast = linear combination of *effects* of alternatives

$$c = \sum_{j=1}^{k} w_j \alpha_j$$

$$\sum_{j=1}^{k} w_j = 0$$

37

## Contrasts

❑ E.g. Compare effect of system 1 to effect of system 2

$$w_1 = 1$$

$$w_2 = -1$$

$$w_3 = 0$$

$$c = (1)\alpha_1 + (-1)\alpha_2 + (0)\alpha_3$$

$$= \alpha_1 - \alpha_2$$

38

## Construct confidence interval for contrasts

- ❑ Need
  - ➢ Estimate of variance
  - ➢ Appropriate value from *t* table
- ❑ Compute confidence interval as before
- ❑ If interval includes 0
  - ➢ Then no statistically significant difference exists between the alternatives included in the contrast

## Variance of random variables

- ❑ Recall that, for independent random variables $X_1$ and $X_2$

$$\mathrm{Var}[X_1 + X_2] = \mathrm{Var}[X_1] + \mathrm{Var}[X_2]$$

$$\mathrm{Var}[aX_1] = a^2 \, \mathrm{Var}[X_1]$$

## Variance of a contrast $c$

$$\text{Var}[c] = \text{Var}[\sum_{j=1}^{k}(w_j\alpha_j)]$$
$$= \sum_{j=1}^{k}\text{Var}[w_j\alpha_j]$$
$$= \sum_{j=1}^{k}w_j^2\text{Var}[\alpha_j]$$

$$s_c^2 = \frac{\sum_{j=1}^{k}(w_j^2 s_e^2)}{kn}$$

$$s_e^2 = \frac{SSE}{k(n-1)}$$

$$df(s_c^2) = k(n-1)$$

❑ Assumes variation due to errors is equally distributed among *kn* total measurements

## Confidence interval for contrasts

$$(c_1, c_2) = c \mp t_{1-\alpha/2;k(n-1)}s_c$$

$$s_c = \sqrt{\frac{\sum_{j=1}^{k}(w_j^2 s_e^2)}{kn}}$$

$$s_e^2 = \frac{SSE}{k(n-1)}$$

## Example

❑ 90% confidence interval for contrast of [Sys1- Sys2]

$$\alpha_1 = -0.1735$$

$$\alpha_2 = -0.1441$$

$$\alpha_3 = 0.3175$$

$$c_{[1-2]} = -0.1735 - (-0.1441) = -0.0294$$

$$s_c = s_e \sqrt{\frac{1^2 + (-1)^2 + 0^2}{3(5)}} = 0.0275$$

$$90\% : (c_1, c_2) = (-0.0784, 0.0196)$$

## Summary

❑ Use one-factor ANOVA to separate total variation into:
- Variation within one system
  - Due to random errors
- Variation between systems
  - Due to real differences (+ random error)

❑ Is the variation due to real differences *statistically* greater than the variation due to errors?

❑ Use contrasts to compare effects of subsets of alternatives