

## 14.3 Tukey's T-procedure

12

John Wilder Tukey (1915 - 2000)

method initially published in 1949

### 14.3.1 Tukey's pairwise comparisons theorem

13

#### Lemma 14.1 Studentized range

Let  $T_1, \dots, T_m$  be a random sample from  $N(\mu, \sigma^2)$ ,  $\sigma^2 > 0$ . Let  $R = \max_{j=1, \dots, m} T_j - \min_{j=1, \dots, m} T_j$

be the range of the sample. Let  $s^2$  be the estimator of  $\sigma^2$  such that  $s^2$  and  $T = (T_1, \dots, T_m)^T$  are independent and  $\frac{vs^2}{\sigma^2} \sim \chi^2_v$  for some  $v > 0$ .

Let  $Q = \frac{R}{s}$ . The distribution of  $Q$  then depends on neither  $\mu$ , nor  $\sigma$ .

Proof: 
$$\frac{R}{s} = \frac{\frac{1}{\sigma} \max(T_j - \mu) - \frac{1}{\sigma} \min(T_j - \mu)}{\frac{s}{\sigma}} = \frac{\max\left(\frac{T_j - \mu}{\sigma}\right) - \min\left(\frac{T_j - \mu}{\sigma}\right)}{\frac{s}{\sigma}}$$

•  $\forall j=1, \dots, m$   $\frac{T_j - \mu}{\sigma} \sim N(0, 1)$ : distribution of numerator of  $Q$  depends neither on  $\mu$ , nor on  $\sigma$

•  $\frac{s}{\sigma} \equiv$  transformation of random variable having  $\chi^2_v \rightarrow$  distribution depends on neither  $\mu$ , nor on  $\sigma$

• numerator & denominator of  $Q$  are independent  $\Rightarrow$  joint distribution of ~~the~~ numer. & denom. of  $Q$

15 depends on neither  $\mu$ , nor on  $\sigma$

=> distribution of  $\frac{\text{numer.}}{\text{denom.}}$  depends on neither  $\mu$ , nor on  $\sigma$

14

Distribution of  $\frac{R}{S}$  still depends on  $m$  and  $v$ .  
sample size of  $T$   $\nearrow$   
degr. of freedom of  $\chi^2$   $\nearrow$

Def 14.5 Studentized range

The random variable  $Q = \frac{R}{S}$  from Lemma 4.1 will be called studentized range of size  $m$  with  $v$  degrees of freedom and its distribution will be denoted as  $q_{m,v}$ .

14

NOTATION:

- $0 < p < 1$  :  $q_{m,v}(p)$  =  $p \cdot 100\%$  quantile of distribution  $q_{m,v}$
- $CDF_{q_{m,v}}(\cdot)$   $\equiv$  cdf of a random variable  $Q \sim q_{m,v}$

Illustrations

$q_{m,v}$

15-16

Theorem 14.2 Tukey's pairwise comparisons theorem, balanced version

Let  $T_1, \dots, T_m$  be independent random variables and let  $T_j \sim \mathcal{N}(\mu_j, \sigma^2)$ ,  $j=1, \dots, m$ , where  $\sigma^2 > 0$  is a known constant. and let  $S^2$  be the estimator of  $\sigma^2$  such that  $S^2$  and  $T = (T_1, \dots, T_m)^T$  are independent and  $\frac{\nu S^2}{\sigma^2} \sim \chi^2_\nu$  for some  $\nu > 0$ .

Then

$$P(\forall j \neq l : |T_j - T_l - (\mu_j - \mu_l)| < q_{m, \nu}(1-\alpha) \cdot \sqrt{\nu^2 S^2}) = 1 - \alpha$$

Proof:

- $\frac{T_j - \mu_j}{\sigma}$  iid  $\sim \mathcal{N}(0, 1)$
- $R := \max_j \frac{T_j - \mu_j}{\sigma} - \min_j \frac{T_j - \mu_j}{\sigma}$
- $\Rightarrow \frac{R}{S} \sim q_{m, \nu}$  (studentized range)

• For any  $0 < \alpha < 1$

$$1-\alpha = P \left( \frac{\max_j (T_j - \mu_j) - \min_j (T_j - \mu_j)}{s} < q_{m,v}^{(1-\alpha)} \right)$$

$$= P \left( \frac{\max_j (T_j - \mu_j) - \min_j (T_j - \mu_j)}{n \cdot s} < q_{m,v}^{(1-\alpha)} \right)$$

$$= P \left( \max_j (T_j - \mu_j) - \min_j (T_j - \mu_j) < n \cdot s q_{m,v}^{(1-\alpha)} \right)$$

$$= P \left( \forall j \neq l \quad |(T_j - \mu_j) - (T_l - \mu_l)| < n \cdot s q_{m,v}^{(1-\alpha)} \right)$$

$$= P \left( \forall j \neq l \quad |T_j - T_l - (\mu_j - \mu_l)| < \sqrt{n^2 s^2} \cdot q_{m,v}^{(1-\alpha)} \right)$$

□

Theorem 14.3 Tukey's ..., general version

18

$$T_j \sim N(\mu_j, n^2 s^2) \quad \rightarrow \quad T_j \sim N(\mu_j, n_j^2 s^2)$$

$$\sqrt{n^2 s^2} \quad \rightarrow \quad \sqrt{\frac{n_j^2 + n_l^2}{2}} \cdot s$$

$$1-\alpha = P(\dots) \quad \rightarrow \quad 1-\alpha \leq P(\dots)$$

• Tukey (1953), Kramer (1956) formulated it.

• Hayter (1984) proved.

Proof is not an easy adaptation of the balanced version proof.

18

## 14.3.2 Tukey's honest significance differences (HSD)

→ application of Tukey's pairwise comparisons theorem to practical problems under the following assumptions: 19

### ASSUMPTIONS

\*  $T = (T_1, \dots, T_m)^T \sim N_m(\mu, \sigma^2 V)$

•  $\mu = (\mu_1, \dots, \mu_m)^T \in \mathbb{R}^m$ ,  $\sigma^2 > 0$ : unknown param.

•  $V = \text{diag}(v_1^2, \dots, v_m^2)$ : known diag. matrix

\*  $S^2$ : estimator of  $\sigma^2$

•  $S^2$  and  $T$  independent

•  $\frac{v S^2}{\sigma^2} \sim \chi^2_\nu$  for some  $\nu > 0$

Example:  $G$ -sample problem (homoscedastic)

$Y_{ij} \stackrel{i.i.d.}{\sim} N(\mu_i, \sigma^2)$ ,  $i=1, \dots, G$ ,  $j=1, \dots, n_i$

( $m = G$ )  
(samples independ.)

\*  $T_i := \frac{1}{n_i} \sum_{j=1}^{n_i} Y_{ij} = \bar{Y}_{i\cdot}$

$(T_1, \dots, T_G)^T = (\bar{Y}_{1\cdot}, \dots, \bar{Y}_{G\cdot})^T \sim N\left(\begin{pmatrix} \mu_1 \\ \vdots \\ \mu_G \end{pmatrix}, \sigma^2 V\right)$

$V = \begin{pmatrix} \frac{1}{n_1} & & \mathbb{0} \\ & \ddots & \\ \mathbb{0} & & \frac{1}{n_G} \end{pmatrix}$

\*  $S^2 := \frac{1}{n-G} \sum_{i=1}^G \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{i\cdot})^2$

$$* \frac{(n-g) \cdot S^2}{\sigma^2} \sim \chi_{n-g}^2, \quad S^2 \perp\!\!\!\perp \begin{pmatrix} T_1, \dots, T_g \\ \bar{Y}_1, \dots, \bar{Y}_g \end{pmatrix}^T$$

End example

Interest in solving the multiple comparison problem:

$$H_{j,l}: \underbrace{\mu_j - \mu_l}_{\theta_{j,l}} = \theta_{j,l}^0, \quad \begin{matrix} j=1, \dots, m-1 \\ l=j+1, \dots, m \end{matrix}$$

$\rightarrow m^* = \binom{m}{2}$  elementary hypotheses

$$\theta^0 = (\theta_{1,2}^0, \theta_{1,3}^0, \dots, \theta_{m-1,m}^0)^T \in \mathbb{R}^{m^*}$$

given

most common  $\theta^0 = \mathbf{0}_{m^*}$

i.e. (global) hypothesis  $H_0$  is

$$H_0: \mu_1 = \dots = \mu_m$$

$$H_{j,l}: \mu_j = \mu_l, \quad \begin{matrix} j=1, \dots, m-1 \\ l=j+1, \dots, m \end{matrix}$$

Let us use Tukey's Theorem:

$$1-\alpha \stackrel{=}{=} \underset{\text{balanced case}}{P(\forall j \neq l \quad |T_j - T_l - (\mu_j - \mu_l)| < q_{m,r}(1-\alpha) \cdot \sqrt{\frac{\sigma_j^2 + \sigma_l^2}{2}})}$$

$$1-\alpha \stackrel{=}{=} P(\forall j \neq l \quad \left| \frac{T_j - T_l - (\mu_j - \mu_l)}{\sqrt{\frac{\sigma_j^2 + \sigma_l^2}{2}}} \right| < q_{m,r}(1-\alpha))$$

Let for given  $\theta_{j,l}^0$   $T_{j,l}(\theta_{j,l}^0) := \frac{T_j - T_l - \theta_{j,l}^0}{\sqrt{\frac{\sigma_j^2 + \sigma_l^2}{2}}} (=: T_{j,l}^0)$

Then  $\forall \theta^0 = (\theta_{1,2}^0, \dots, \theta_{m-1,m}^0)^T, \forall 0 < \alpha < 1$

$$\begin{aligned} 1-\alpha &\stackrel{=}{=} P(\forall j \neq l \quad |T_{j,l}^0| < q_{m,r}(1-\alpha); \theta = \theta^0) = \\ &= P(\forall j \neq l \quad \left| \frac{T_j - T_l - \theta_{j,l}^0}{\sqrt{\frac{\sigma_j^2 + \sigma_l^2}{2}}} \right| < q_{m,r}(1-\alpha); \theta = \theta^0) = \\ &= P(\forall j \neq l \quad (\theta_{j,l}^{TL}(\alpha), \theta_{j,l}^{TU}(\alpha)) \ni \theta_{j,l}^0; \theta = \theta^0), \end{aligned}$$

where  $\theta_{j,l}^{TL}(\alpha) = T_j - T_l - q_{m,r}(1-\alpha) \sqrt{\frac{\sigma_j^2 + \sigma_l^2}{2}}$

$\theta_{j,l}^{TU}(\alpha) = T_j - T_l + q_{m,r}(1-\alpha) \sqrt{\frac{\sigma_j^2 + \sigma_l^2}{2}}$

It seems that we have just derived a set of simultaneous confidence

intervals for  $\theta_{j,l} = \mu_j - \mu_l, j=1, \dots, m-1, l=j+1, \dots, m$

# Theorem 14.4 Tukey's honest significance differences 20

Short: • simultaneous confidence intervals  
for  $\theta_{j,l} = \mu_j - \mu_l$ ,  $j=1, \dots, m-1$   
 $l=j+1, \dots, m$

• balanced case

$$1-\alpha = P(\forall j \neq l \mid \theta_{j,l}^{TL}(\alpha), \theta_{j,l}^{TU}(\alpha)) \ni \theta_{j,l}^0; \theta = \theta^0)$$

• unbalanced case

$$1-\alpha \leq P(\forall j \neq l \mid \theta_{j,l}^{TL}(\alpha), \theta_{j,l}^{TU}(\alpha)) \ni \theta_{j,l}^0; \theta = \theta^0)$$

• P-values adjusted for multiple comparison:

$$p_{j,l}^T = 1 - \text{CDF}_{q, m, v}(|t_{j,l}^0|), j < l$$

$$t_{j,l}^0 = \text{value of } T_{j,l}(\theta_{j,l}^0) = \frac{T_j - T_l - \theta_{j,l}^0}{\sqrt{\frac{v_1 + v_2}{2} s^2}}$$

attained with given data

$$H_{j,l}: \theta_{j,l} = \theta_{j,l}^0$$

$\mu_j - \mu_l$

Proof: The only issue to clarify are  
the p-values adjusted for multiple  
comparison.



We have

$$H_{j,e} \mid (\theta_{j,e}^{TL}(\alpha), \theta_{j,e}^{TU}(\alpha)) \neq \theta_{j,e}^0$$

$\Leftrightarrow H_{j,e}$  rejected by MCP

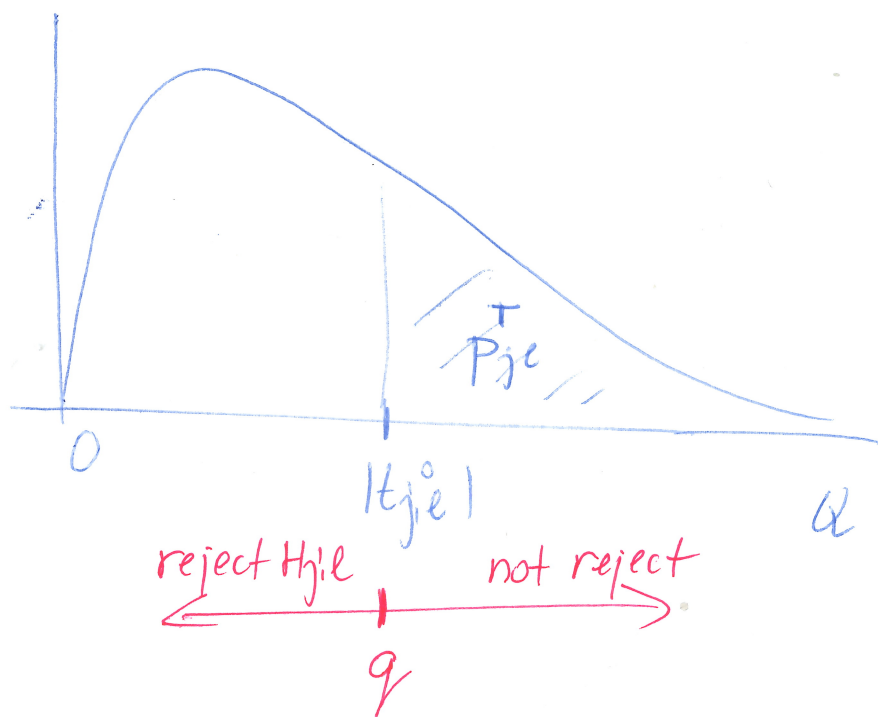
$$\Leftrightarrow |T_{j,e}(\theta_{j,e}^0)| \geq q_{m,r}(1-\alpha)$$

Hence  $p_{j,e}^T = \inf \{ \alpha : |T_{j,e}(\theta_{j,e}^0)| \geq q_{m,r}(1-\alpha) \}$

$$= \inf \{ \alpha : |T_{j,e}(\theta_{j,e}^0)| \geq q_{m,r}(1-\alpha) \}$$

Let  $t_{j,e}^0 =$  value of  $T_{j,e}(\theta_{j,e}^0)$  attained with given data

$$p_{j,e}^T = \inf \{ \alpha : |t_{j,e}^0| \geq q_{m,r}(1-\alpha) \}$$



$$P_{j,e}^T = 1 - \text{CDF}_{q_{m,r}}(|t_{j,e}^0|)$$

### 14.3.3 Tukey's HSD in a linear model 21

Consider a linear model  $Y|X \sim N_n(X\beta, \sigma^2 I_n)$   
 $\text{rank}(X_{n \times k}) = k < n$

•  $L = \begin{pmatrix} l_1^T \\ \vdots \\ l_m^T \end{pmatrix}$  given matrix

•  $\eta_i = L\beta = (l_1^T \beta, \dots, l_m^T \beta)^T = \begin{matrix} \text{parameters} \\ \text{of interest} \end{matrix}$   
 $\eta_1, \dots, \eta_m$

•  $L$  must be such that  
 $(\sigma^2)^{-1} V := L(X^T X)^{-1} L^T$  is diagonal

→ properties of LSE under normality

•  $T = \hat{\eta} = L \hat{\beta} = (l_1^T \hat{\beta}, \dots, l_m^T \hat{\beta})^T \sim N_m(\eta, \sigma^2 V)$   
(given  $X$ )

•  $\frac{(n-k) \text{MSE}}{\sigma^2} \sim \chi_{n-k}^2$  (both given  $X$  as well as unconditionally)

•  $T$  and MSE independent (given  $X$ )

→ all subsequent inference will be conditional given  $X$

Tukey's pairwise comparisons theorems allow for inference on

21

$$\begin{aligned}\theta_{j,l} &= \eta_j - \eta_l, \quad j < l \\ &= (\mathbf{c}_j - \mathbf{c}_l)^T \beta\end{aligned}$$

• simultaneous CI's for  $\theta_{j,l}$  are

$$\theta_{j,l}^{TL}(\alpha) = \hat{\eta}_j - \hat{\eta}_l - q_{m, m-k}(1-\alpha) \sqrt{\frac{\sigma_j^2 + \sigma_l^2}{2} \text{MSE}}$$

$$\theta_{j,l}^{TU}(\alpha) = \hat{\eta}_j - \hat{\eta}_l + q_{m, m-k}(1-\alpha) \sqrt{\frac{\sigma_j^2 + \sigma_l^2}{2} \text{MSE}}$$

• p-values adjusted for multiple comparison to test elementary hypotheses

$H_{j,l}^0: \theta_{j,l} = \theta_{j,l}^0, \quad j < l$  based on statistics

$$T_{j,l}^0 = T_{j,l}^0(\theta_{j,l}^0) = \frac{\hat{\eta}_j - \hat{\eta}_l - \theta_{j,l}^0}{\sqrt{\frac{\sigma_j^2 + \sigma_l^2}{2} \text{MSE}}}$$

$$p_{j,l}^T = 1 - \text{CDF}_{q, m, m-k}(|t_{j,l}^0|)$$

## One-way classification

(Homoscedastic) G-sample problem

$$Y_{gj} \stackrel{iid}{\sim} N(m_g, \sigma^2), \quad g=1, \dots, G, \quad j=1, \dots, n_g$$

$Y_{gj}$  all independent

Lemma 13.1 (or previous knowledge)

$$\bar{Y} := \begin{pmatrix} \bar{Y}_1 \\ \vdots \\ \bar{Y}_G \end{pmatrix} \sim N_G \left( \begin{pmatrix} m_1 \\ \vdots \\ m_G \end{pmatrix}, \sigma^2 \underbrace{\begin{pmatrix} \frac{1}{n_1} & & 0 \\ & \ddots & \\ 0 & & \frac{1}{n_G} \end{pmatrix}}_V \right)$$

$$\frac{(n-G) \text{Mse}}{\sigma^2} \sim \chi^2_{n-G}$$

Mse and  $\bar{Y}$  independent (given covariates, given groups allocation)

Tukey's simultaneous conf. intervals

for  $m_g - m_h, g < h$ :

$$\bar{Y}_g - \bar{Y}_h \pm q_{G, n-G} (1-\alpha) \sqrt{\frac{1}{2} \left( \frac{1}{n_g} + \frac{1}{n_h} \right) \text{Mse}}$$

(lead coverage of  $1-\alpha$  in a balanced case)

→ R function TukeyHSD

# Two-way classification, BALANCED data

23

$$Y_{g,h,j} \stackrel{iid}{\sim} N(\mu_{gh}, \sigma^2), \quad g=1, \dots, G$$
$$h=1, \dots, H, \quad j=1, \dots, n_{gh}$$

$Y_{g,h,j}$  all independent

$= J$   
(balanced)

Perhaps additivity can be assumed for effect of the grouping factors.

## Lemma 13.2

- under both additivity and interaction model (=no structure)

$$\mathbb{T} = \begin{pmatrix} \bar{Y}_{1\cdot} \\ \vdots \\ \bar{Y}_{G\cdot} \end{pmatrix} \sim N_G \left( \begin{pmatrix} \bar{\mu}_{1\cdot} \\ \vdots \\ \bar{\mu}_{G\cdot} \end{pmatrix}, \sigma^2 \begin{pmatrix} \frac{1}{J \cdot H} & 0 \\ 0 & \frac{1}{J \cdot H} \end{pmatrix} \right)$$
$$\bar{Y}_{g\cdot} = \frac{1}{J \cdot H} \sum_{h=1}^H \sum_{j=1}^J Y_{g,h,j}$$
$$\bar{\mu}_{g\cdot} = \frac{1}{H} \sum_{h=1}^H \mu_{gh}$$

$$\frac{v_e^* \text{MSE}^*}{\sigma^2} \sim \chi^2_{v_e^*}, \quad \text{MSE}^* \text{ and } \mathbb{T} \text{ independent (given the grouping factors)}$$

depending on which structure is assumed for the group means:

$$v_e^* = n - (G+H-1) \quad : \text{additive model}$$

$$v_e^* = n - G \cdot H \quad : \text{interaction model}$$

Tukey's simultaneous confidence intervals  
for  $\bar{m}_{g_1} - \bar{m}_{g_2}$ ,  $g_1 < g_2$

$$\bar{Y}_{g_1} - \bar{Y}_{g_2} \pm q_{\alpha, ve}^* (1-\alpha) \sqrt{\frac{1}{J_{iH}} \cdot MSE^*}$$

---

Remember that in the additive model

$$\bar{m}_{g_1} - \bar{m}_{g_2} = m_{g_1, h} - m_{g_2, h}$$

for any  $h=1, \dots, H$

→ R function TukeyHSD

---

UNBALANCED data

- TukeyHSD is still producing something
- What? → see notes