

# Método de *Diferencia Honestamente Significativa* de Tukey (Tiuki)

22 de marzo de 2022

## 1. Recorrido

Sean una muestra  $X_1, \dots, X_n$  obtenida de la población  $X$ . Se define su recorrido (o rango) como

$$R = \max_i X_i - \min_i X_i = X_{(n)} - X_{(1)}$$

## 2. Distribución del recorrido estudentizado

Si  $X_i \hookrightarrow \mathcal{N}(0, 1)$  y  $Y \hookrightarrow \chi_r^2$  independiente de  $X_i$ , entonces

$$\frac{X_{(n)} - X_{(1)}}{\sqrt{Y/r}} \hookrightarrow Q_{n,r}$$

En R se usa `ptukey` para la función de distribución y `qtukey` para la función cuantil de  $Q_{n,r}$ .

## 3. Varias muestras

Sean las  $q$  poblaciones homoscedásticas independientes  $X_1, \dots, X_q$  con  $X_i \hookrightarrow \mathcal{N}(\mu_i, \sigma)$ . De cada población  $X_i$  se extrae una muestra  $\vec{X}_i = (X_{i1}, \dots, X_{in_i})$ .

Sea el estimador de la varianza intrapoblacional

$$\hat{S}^2 = \frac{\sum_i (n_i - 1) \hat{S}_i^2}{n - q} = \frac{\sum_i (n_i - 1) \frac{\sum_j (X_{ij} - \bar{X}_i)^2}{n_i - 1}}{n - q}$$

con  $n = \sum_i n_i$ , del que se sabe

$$\frac{(n - q) \hat{S}^2}{\sigma^2} \hookrightarrow \chi_{n-q}^2$$

Sea (1) el índice de la muestra con media muestral más baja y ( $q$ ) el de la más alta. Como cada  $\bar{X}_i$  es independiente de cada  $\hat{S}_i^2$  y aquéllos lo son entre sí, entonces  $\hat{S}^2$  es independiente de todos aquéllos y

$$\frac{\frac{\bar{X}_{(q)} - \mu_{(q)}}{\sigma/\sqrt{n_{(q)}}} - \frac{\bar{X}_{(1)} - \mu_{(1)}}{\sigma/\sqrt{n_{(1)}}}}{\sqrt{\frac{\hat{S}^2}{\sigma^2}}} = \frac{\frac{\bar{X}_{(q)} - \mu_{(q)}}{1/\sqrt{n_{(q)}}} - \frac{\bar{X}_{(1)} - \mu_{(1)}}{1/\sqrt{n_{(1)}}}}{\hat{S}} \hookrightarrow Q_{n,r}$$

Bajo  $H_0$  del ANOVA,  $\forall i \mu_i = \mu$  y

$$\frac{\frac{\bar{X}_{(q)} - \mu_{(q)}}{1/\sqrt{n_{(q)}}} - \frac{\bar{X}_{(1)} - \mu_{(1)}}{1/\sqrt{n_{(1)}}}}{\hat{S}} = \frac{\frac{\bar{X}_{(q)} - \mu}{1/\sqrt{n_{(q)}}} - \frac{\bar{X}_{(1)} - \mu}{1/\sqrt{n_{(1)}}}}{\hat{S}} \hookrightarrow Q_{n,r}$$

#### 4. Varias muestras equilibradas

Si el modelo es equilibrado,  $n_i = m = n/q \forall i$  y

$$\frac{\frac{\bar{X}_{(q)} - \mu}{1/\sqrt{n_{(q)}}} - \frac{\bar{X}_{(1)} - \mu}{1/\sqrt{n_{(1)}}}}{\hat{S}} = \frac{\frac{\bar{X}_{(q)} - \mu}{1/\sqrt{m}} - \frac{\bar{X}_{(1)} - \mu}{1/\sqrt{m}}}{\hat{S}} = \frac{\bar{X}_{(q)} - \bar{X}_{(1)}}{\hat{S}/\sqrt{m}} \hookrightarrow Q_{n,r} \quad (1)$$

La diferencia entre las medias más alejadas acota superiormente el valor absoluto de la diferencia entre cualquier par de medias:

$$\Pr \left[ \frac{|\bar{X}_i - \bar{X}_j|}{\hat{S}/\sqrt{m}} > k \right] \leq \Pr \left[ \frac{\bar{X}_{(q)} - \bar{X}_{(1)}}{\hat{S}/\sqrt{m}} > k \right] \quad (2)$$

por lo que haciendo  $k$  el cuantil de orden  $1 - \alpha$  de  $Q_{q,n-q}$  se definen regiones críticas a nivel  $\alpha$  para los múltiples contrastes  $H_0^{ij}: \mu_i = \mu_j$ ,  $H_1^{ij}: \mu_i \neq \mu_j$ .

#### 5. Varias muestras desequilibradas

La última fracción de (1) puede escribirse así:

$$\frac{\bar{X}_{(q)} - \bar{X}_{(1)}}{\hat{S}/\sqrt{m}} = \frac{\bar{X}_{(q)} - \bar{X}_{(1)}}{\sqrt{\frac{\hat{S}^2}{2} \left( \frac{1}{m} + \frac{1}{m} \right)}}$$

lo que recuerda a la expresión derivada de la diferencia de medias:

$$\frac{\bar{X}_i - \bar{X}_j}{\sqrt{S^2 \left( \frac{1}{n_i} + \frac{1}{n_j} \right)}} \xrightarrow{H_0} t_{n-q}$$

En caso de muestras desequilibradas, con  $n_i \neq n_j \exists i, j$ , se emplea una expresión aproximada de (2):

$$\Pr \left[ \frac{|\bar{X}_i - \bar{X}_j|}{\sqrt{\frac{\hat{S}^2}{2} \left( \frac{1}{n_i} + \frac{1}{n_j} \right)}} > k \right] \lesssim 1 - \alpha$$

## 6. Ejemplo en R

```
> cw <- chickwts
> levels(cw$feed) <- abbreviate(levels(cw$feed),2)
> a <- aov (weight ~ feed, cw)
> t <- TukeyHSD (a)
```

Tabla de Tukey:

```
> t
```

```
Tukey multiple comparisons of means
 95% family-wise confidence level
```

```
Fit: aov(formula = weight ~ feed, data = cw)
```

```
$feed
      diff      lwr      upr      p adj
hr-cs -163.383333 -232.346876 -94.41979 0.0000000
ln-cs -104.833333 -170.587491 -39.07918 0.0002100
mt-cs  -46.674242 -113.906207  20.55772 0.3324584
sy-cs  -77.154762 -140.517054 -13.79247 0.0083653
sn-cs   5.333333  -60.420825  71.08749 0.9998902
ln-hr   58.550000 -10.413543 127.51354 0.1413329
mt-hr  116.709091  46.335105 187.08308 0.0001062
sy-hr   86.228571  19.541684 152.91546 0.0042167
sn-hr  168.716667  99.753124 237.68021 0.0000000
mt-ln   58.159091  -9.072873 125.39106 0.1276965
sy-ln   27.678571 -35.683721  91.04086 0.7932853
sn-ln  110.166667  44.412509 175.92082 0.0000884
sy-mt  -30.480519 -95.375109  34.41407 0.7391356
sn-mt   52.007576 -15.224388 119.23954 0.2206962
sn-sy   82.488095  19.125803 145.85039 0.0038845
```

Sus elementos se calculan así:

```
> i <- "ln" # segunda fila de la tabla, por ejemplo
> j <- "cs"
> n <- table (cw$feed) ; q <- length(n) ; N <- sum(n)
> CME <- summary(a)[[1]][["Residuals","Mean Sq"]]
> Xi <- mean (cw$weight [cw$feed == i])
> Xj <- mean (cw$weight [cw$feed == j])
> alfa <- 0.05
> dif <- Xi-Xj
> den <- as.numeric (sqrt (CME/2 * (1/n[i] + 1/n[j])))
> ## as.numeric para quitar etiqueta
> int <- setNames (dif + c(-1,1) * den * qtukey (1-alfa, q, N-q),
```

```

+           c ("lwr", "upr"))
> c (diff = round (dif, 6),
+   int,
+   "p adj" = round (1 - ptukey (abs(dif)/den, q, N-q), 7))

      diff      lwr      upr      p adj
-104.83333 -170.58749  -39.07918   0.00021

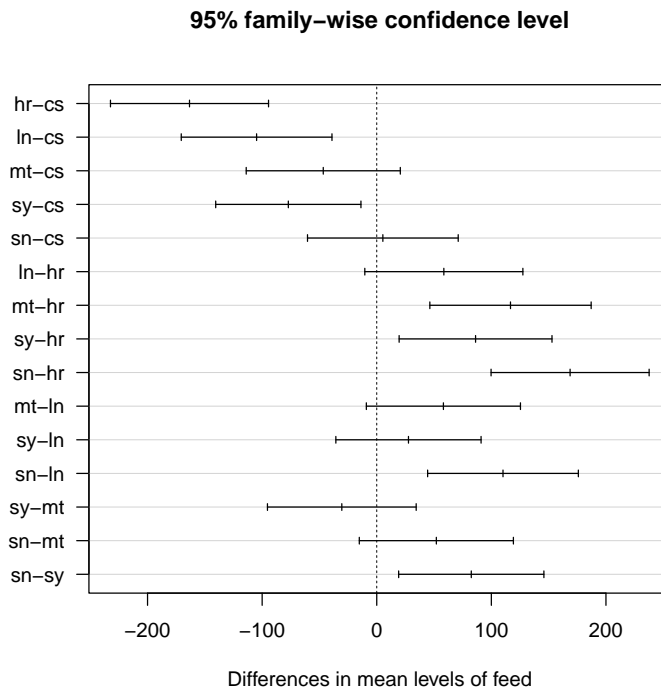
```

Representación gráfica de los intervalos (¿está el cero incluido en el intervalo?):

```

> par (las = 1) # etiquetas horizontales en el eje Y
> plot (t)

```



## 7. Simulación

Comparación de la tasa de rechazos usando ANOVA y Tukey:

```

> mu <- c (0, 0, 0.2) ; q <- length(mu)
> m <- 100 ; g <- factor (rep (1:q, each=m))
> table (data.frame (t (replicate (10000, {
+   x <- unlist (lapply (mu,

```

```

+                                     function (mui) rnorm (m, mui)))
+   a <- aov (x ~ g)
+   H1anova <- summary(a) [[1]] ["g", "Pr(>F)"] < alfa
+   H1tukey <- any (TukeyHSD(a)$g[, "p adj"] < alfa)
+   c (H1anova=H1anova, H1tukey=H1tukey)
+ }))))

```

```

      H1tukey
H1anova FALSE TRUE
FALSE  6974 129
TRUE   180 2717

```