

## Kapitel 10: "Comparing Two Means and Two Variances"

$X_1, \dots, X_m$  är stickprov på  $N(\mu_x, \sigma_x^2)$   
 $Y_1, \dots, Y_n$  är stickprov på  $N(\mu_y, \sigma_y^2)$  ↗ oberoende

Sats  $\bar{X} - \bar{Y}$  är  $N(\mu_x - \mu_y, \sigma_x^2/m + \sigma_y^2/n)$

(Bevisas mha MGF som i kapitel 7, se även kapitel 9.)

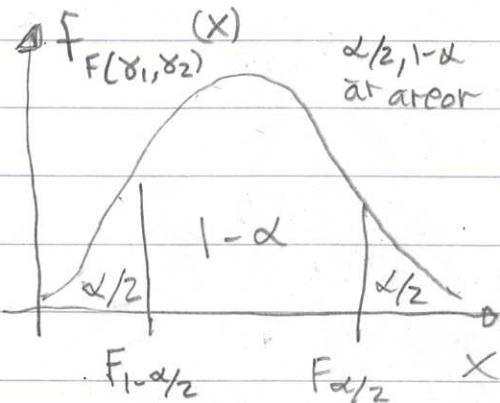
Definition Om  $X_{\gamma_1}^2$  och  $Y_{\gamma_2}^2$  är oberoende  $\chi^2(\gamma_1)$ -respektive  $\chi^2(\gamma_2)$ -fördelade är

$\frac{X_{\gamma_1}^2/\gamma_1}{Y_{\gamma_2}^2/\gamma_2}$   $F(\gamma_1, \gamma_2)$ -fördelad.

$F(\gamma_1, \gamma_2)$  variabler är positiva.

I bokens tabell IX finner man talet  $F_\alpha$  sådant att

$$P(F(\gamma_1, \gamma_2) > F_\alpha) = \alpha \text{ för } \alpha \text{ nära 0 och 1}$$



Sats  $\frac{S_x^2/\sigma_x^2}{S_y^2/\sigma_y^2}$  är  $F(m-1, n-1)$ -fördelad.

(Här är förstas  $S_x^2$  och  $S_y^2$  X- och Y-stickpovens stickprovsvarianter.)

Bevis Följer direkt av föregående definition med  $\gamma_1 = m-1$  och  $\gamma_2 = n-1$  eftersom  $(m-1)S_x^2/\sigma_x^2$  och  $(n-1)S_y^2/\sigma_y^2$  är  $\chi^2_{(m-1)}$  resp  $\chi^2_{(n-1)}$ -fördelade. □

Sats  $\frac{\sigma_x^2}{\sigma_y^2} \in \left[ \frac{S_x^2}{S_y^2} \frac{1}{F_{\alpha/2}}, \frac{S_x^2}{S_y^2} \frac{1}{F_{1-\alpha/2}} \right]$  med konfidensgrad  $1-\alpha$ .  
 $m-1, n-1$  frihetsgrader

Bevis Händelsen ovan är samma som  $\frac{S_x^2/\sigma_x^2}{S_y^2/\sigma_y^2} \in [F_{1-\alpha/2}, F_{\alpha/2}]$ . □

\*

Vi önskar testa  $H_0: \sigma_x^2 = \sigma_y^2$  mot en av

$$\begin{cases} H_1: \sigma_x^2 \neq \sigma_y^2 \\ H_1: \sigma_x^2 > \sigma_y^2 \\ H_1: \sigma_x^2 < \sigma_y^2 \end{cases}$$

m-1, n-1 frihetsgrader

Sats Förfasta  $H_0$  om  $\frac{S_x^2}{S_y^2} \notin \begin{cases} [F_{1-\alpha/2}, F_{\alpha/2}] \\ (0, F_\alpha] \\ [F_{1-\alpha}, \infty) \end{cases}$  med  $\alpha$ -fel.

(Om  $H_0$  ej förfastas enligt satsen så godtages  $H_0$ .)

Notera att de tre varianterna av testen är designade upptäcka motsvarande tre  $H_1$ -avvikelse från  $H_0$ .

Sats

Om  $\sigma_x^2 = \sigma_y^2$  är  $\sqrt{S_p^2(1/m+n)}$   $T(m+n-2)$ -fördelad

$$\text{där } S_p^2 = \frac{(m-1)S_x^2 + (n-1)S_y^2}{m+n-2}.$$

Bevis. Överkurs.  $\bar{X} - \bar{Y}$  är  $N(\mu_x - \mu_y, \sigma^2/m + \sigma^2/n)$  och  $((m-1)S_x^2 + (n-1)S_y^2)/\sigma^2$  är  $\chi^2(m+n-2)$  vilket ger

$$\frac{(\bar{X} - \bar{Y}) - (\mu_x - \mu_y)}{\sqrt{\sigma^2(1/m+n)}} / \sqrt{\frac{S_p^2/\sigma^2}{m+n-2}} = \frac{N(0,1)}{\sqrt{\chi^2_{(m+n-2)/(m+n-2)}}} = T(m+n-2). \square$$

Sats Om  $\sigma_x^2 = \sigma_y^2$  är

m+n-2 frihetsgrader

$$\mu_x - \mu_y \in \bar{X} - \bar{Y} \pm t_{\alpha/2} \sqrt{S_p^2(1/m+n)} \text{ med konfidensgrad } 1-\alpha.$$

Bevis. Följer direkt av föregående sats.  $\square$

Sats Om  $\sigma_x^2 = \sigma_y^2$  kan  $H_0: \mu_x = \mu_y$  testas med  $\alpha$ -felx genom

m+n-2 frihetsgrader

förfasta  $H_0$  mot  $\begin{cases} H_1: \mu_x \neq \mu_y & \text{om } \frac{\bar{X} - \bar{Y}}{\sqrt{S_p^2(1/m+n)}} \notin \{-t_{\alpha/2}, t_{\alpha/2}\} \\ H_1: \mu_x > \mu_y & \text{om } \frac{\bar{X} - \bar{Y}}{\sqrt{S_p^2(1/m+n)}} \notin (-\infty, t_\alpha] \\ H_1: \mu_x < \mu_y & \text{eller } \end{cases}$

Bevis. Följer också direkt av ovan.  $\square$

Man kan lätt modifiera testen ovan till test av

$H_0: \mu_x - \mu_y = \mu_0$  mot något av alternativen

$H_1: \mu_x - \mu_y \neq \mu_0$ ,  $H_1: \mu_x - \mu_y > \mu_0$  eller  $H_1: \mu_x - \mu_y < \mu_0$ .

- \* Om man ej vet att  $\sigma_x^2 = \sigma_y^2$  (=ofta realistiskt!) kan man ersätta

$$\sqrt{S_p^2(1/m + 1/n)} \text{ med } \sqrt{S_x^2/m + S_y^2/n}$$

i alla tre föregående resultat varefter de blir approximativ (istället exakt) sanna med

$$T(\chi)\text{-Fördelning med } \chi = \frac{(S_x^2/m + S_y^2/n)^2}{\frac{(S_x^2/m)^2}{m-1} + \frac{(S_y^2/n)^2}{n-1}}$$

- \* Om  $m=n$  behöver man ej veta att  $\sigma_x^2 = \sigma_y^2$  utan man kan studera sk. parade data

$$Z_i = X_i - Y_i, i=1, \dots, n \text{ som är } N(\underbrace{\mu_x - \mu_y}_{=0}, \underbrace{\sigma_x^2 + \sigma_y^2}_{=\sigma^2})$$

och sedan många dem testa  $H_0: \mu = \mu_x - \mu_y = \mu_0$  mot  $H_1: \mu_x - \mu_y \neq \mu_0$  enligt metodiken från kapitel 8.

## Wilcoxon rank summe test

10.5

Liksom i slutet av kapitel 8 tittar vi nu på icke-parametriska tester där vi ej behöver antaga att våra stickprov är N-fördelade.

- \*  $X_1, \dots, X_m$  har teoretisk median  $M_x$   
 $Y_1, \dots, Y_n$  har teoretisk median  $M_y$   $m \leq n$
  - \*
- Vi vill testa  $H_0: M_x = M_y$  mot en av
- $H_1: M_x \neq M_y$   
 $H_1: M_x > M_y$   
 $H_1: M_x < M_y$
- 1) Ordna alla data  $X_1, \dots, X_m, Y_1, \dots, Y_n$  i växande ordning  $Z_1 < Z_2 < \dots < Z_{m+n}$   $\rightarrow$  samma siffror ordnade
  - 2) Till dela  $Z_1, \dots, Z_{m+n}$  rankerna  $R=1, \dots, R_{m+n}=m+n$  och låt  $W_x$  och  $W_y$  vara summan av alla  $X$ -respektive  $Y$ -ranker så att
- $$W_x + W_y = \frac{(m+n)(m+n+1)}{2}$$
- 3) Om  $H_0$  är sann bör  $W_x \approx \frac{m(m+n+1)}{2}$ . Vad som är tillräcklig avvikelse från detta för kvarna förkasta  $H_0$  framgår av tabell I i boken.
  - 4) Tabell I visar onormalt små värde och sedan är

$$\text{onormalt} = \frac{m(m+n+1)}{2} + \left( \frac{m(m+n+1)}{2} - \text{litet värde} \right) = m(m+n+1) - \text{litet värde}$$

Exempel (Example 10.6.1 i M&A) Jag har två datamaterial med  $m=12$  och  $n=15$

Brand B		Brand A	
69.3	52.6	28.6	30.6
56.0	34.4	25.1	31.8
22.1	60.2	26.4	41.6
47.6	43.8	34.9	21.1
53.2		29.8	36.0
48.1		28.4	37.9
23.2		38.5	13.9
13.8		30.2	

Jag utför steg 1-2 i Wilcoxon ranksumme test:

Observation	13.8	13.9	21.1	22.1	23.2	25.1	26.4	28.4	28.6
Brand	B	A	A	B	B	A	A	A	A
Rank	1	2	3	4	5	6	7	8	9

Observation	29.8	30.2	30.6	31.8	34.4	34.9	36.0	37.9	38.5
Brand	A	A	A	A	B	A	A	A	A
Rank	10	11	12	13	14	15	16	17	18

Observation	41.6	43.8	47.6	48.1	52.6	53.2	56.0	60.2	69.3
Brand	A	B	B	B	B	B	B	B	B
Rank	19	20	21	22	23	24	25	26	27

$$W_B = 1 + 4 + 5 + 14 + 20 + 21 + 22 + 23 + 24 + 25 + 26 + 27 = 212$$

och önskar testa  $H_0: M_B = M_A$  mot  $H_1: M_B > M_A$  ensidigt.  
Enligt tabell X är ett  $\alpha=0.05$  onormalt stort värde

$$m(m+n+1) - \text{onormalt litet värde} = 12(12+15+1) - 134,202 = 201,8$$

så  $H_0$  kan förkastas och  $H_1$  godtas.

Om  $m = n$  kan man studera parade data

$$Z_i = X_i - Y_i \text{ för } i=1, \dots, m$$

så kan  $H_0: M_x = M_y$  testas mot  $H_1: M_x \stackrel{+}{<} M_y$   
genom använda Wilcoxon signed rank test  
från slutet kapitel 8 för Z-data och testa  
 $H_0: M_z = 0$  mot  $H_1: M_z \stackrel{+}{<} 0$  mha. denna/dessa.

## Kapitel 11: "Simple Linear Regression and Correlation"

- \* Linjär regression. Jag har ett datamaterial

$(X_i, Y_i)_{i=1}^n$  där  $\bar{Y}_i = \beta_0 + \beta_1 X_i + E_i$  med

$\begin{cases} X_1, \dots, X_n \text{ kända deterministiska data} \\ \beta_0 \text{ och } \beta_1 \text{ parametrar med okända värden} \\ E_1, \dots, E_n \text{ oberoende } N(0, \sigma^2) \text{ variabler med okänt värde} \\ \bar{Y}_1, \dots, \bar{Y}_n \text{ kända observerade stokastiska data} \end{cases}$

- \* Uppgiften är att skatta  $\beta_0$  och  $\beta_1$ .
- \* Lösning: Minsta kvadratmetoden - minimera

$$SSE = \sum_{i=1}^n (\bar{Y}_i - (\beta_0 + \beta_1 X_i))^2 \text{ map } \beta_0 \text{ och } \beta_1$$

$$\frac{\partial SSE}{\partial \beta_0} = -2 \sum_{i=1}^n (\bar{Y}_i - (\beta_0 + \beta_1 X_i)) = 0$$

$$\frac{\partial SSE}{\partial \beta_1} = -2 \sum_{i=1}^n X_i (\bar{Y}_i - (\beta_0 + \beta_1 X_i)) = 0 \quad \text{ger}$$

Sats:

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X} \quad \text{och} \quad \hat{\beta}_1 = \frac{\sum_{i=1}^n X_i \bar{Y}_i - n \bar{X} \bar{Y}}{\sum_{i=1}^n X_i^2 - n \bar{X}^2}$$

(Beweis Bara att flytta om i ekvationerna ovan.  $\square$ )

Sats

$$\hat{\beta}_0 \text{ är } N\left(\beta_0, \frac{\sum_{i=1}^n x_i^2}{n \sum_{i=1}^n (x_i - \bar{x})^2} \sigma^2\right), \quad \hat{\beta}_1 \text{ är } N\left(\beta_1, \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}\right)$$

(Boken sid 388 står en del fakta som kan göra enklare beräkning av summor men de är lite utdaterade - det är länge sedan datorer blev billigare använda än utbildade människor.)

(Bevis. Örekurs. Detta går göra själv med noga användning av det jag lärt er men är "risigt". □)

Sats  $\hat{\sigma}^2 = \frac{SSE}{n-2}$  med  $(\hat{\beta}_0, \hat{\beta}_1) = (\hat{\beta}_0, \hat{\beta}_1)$  är vr-skattning av  $\sigma^2$ .

Sats  $SSE/\sigma^2$  är  $\chi^2_{(n-2)}$  så att följande är  $T(n-2)$ :

$$(\hat{\beta}_0 - \beta_0) / \sqrt{\sum_{i=1}^n x_i^2 \frac{SSE}{n-2} / \sum_{i=1}^n (x_i - \bar{x})^2} \stackrel{D}{=} (\hat{\beta}_1 - \beta_1) / \sqrt{\frac{SSE}{n-2} / \sum_{i=1}^n (x_i - \bar{x})^2}$$

(Mha satsen ovan går göra konfidensintervall för  $\beta_0$  och  $\beta_1$  i analogi med kapitel 8 samt testa hypoteser angående värdet för  $\beta_0$  och  $\beta_1$ .)

Sats  $n-2$  frihetsgrader,

$$Y|X = E(Y) = \beta_0 + \beta_1 X \in \hat{\beta}_0 + \hat{\beta}_1 X \pm t_{\alpha/2} \sqrt{\frac{1}{n} + \frac{(X - \bar{X})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \frac{SSE}{n-2}}$$

$$Y \in \hat{\beta}_0 + \hat{\beta}_1 X \pm t_{\alpha/2} \sqrt{1 + \frac{1}{n} + \frac{(X - \bar{X})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \frac{SSE}{n-2}} \text{ med konfidensgrad } 1 - \alpha.$$

(Konfidensintervallet för  $Y$  gäller för en "ny" observation av  $Y$  variabeln, inte  $Y_1, \dots, Y_n$ .)

Korrelation. Nu ändras grundförutsättningarna:

$(X_i, Y_i)_{i=1}^n$  är stickprov på  $(X, Y)$

(Här kan  $X$  och  $Y$  (liksom då  $X_i$  och  $Y_i$ ) bero av varandra medan  $(X_1, Y_1), \dots, (X_n, Y_n)$  är oberoende.)

Sats  $\text{Cov}(X, Y)$

$$P = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}} \text{ skattas med } \hat{P} = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2}}$$

Sats

$H_0: p=0$  förkastas mot

$$\begin{cases} H_1: p \neq 0 \\ H_1: p > 0 \text{ om } \frac{\hat{P}\sqrt{n-2}}{\sqrt{1-\hat{P}^2}} \notin \left[-t_{\alpha/2}, t_{\alpha/2}\right] \\ H_1: p < 0 \end{cases}$$

med  $\alpha$ -fel  $\alpha$

$n-2$  frihetsgrader

Sats  $\frac{1}{2} \ln\left(\frac{1+\hat{P}}{1-\hat{P}}\right) \approx N\left(\frac{1}{2} \ln\left(\frac{1+p}{1-p}\right), \frac{1}{n-3}\right)$  så att

$$\frac{1}{2} \ln\left(\frac{1+p}{1-p}\right) \in \frac{1}{2} \ln\left(\frac{1+\hat{P}}{1-\hat{P}}\right) \pm \lambda_{\alpha/2} \sqrt{\frac{1}{n-3}} \text{ approximativt med konfidensgrad } 1-\alpha.$$

Genom "lösa ut"  $p$  fås konfidensintervall för  $p$ .

Man kan testa tex  $H_0: p=p_0$  mot  $H_1: p \neq p_0$  med  $\alpha$ -fel  $\alpha$  genom förkasta  $H_0$  om  $\frac{1}{2} \ln\left(\frac{1+p_0}{1-p_0}\right) \notin$  intervallet.