

1) Forskningsmidler

a) H_0 : Tildeling av prosjektmidler og kjønn er uavhengige

H_1 : Tildeling av prosjektmidler og kjønn er avhengige

Kontingenstabell:

	Menn	Kvinner	
Nei	1345	1011	2356
Ja	290	177	467
	1635	1188	2823

Rate under uavhengighet: $p = \frac{467}{2823}$
 (for tildeling av prosjektmidler)

⇒ Forventede frekvenser under uavhengighet

	Menn	Kvinner
Nei	$E(X_{11})$	$E(X_{12})$
Ja	$E(X_{21})$	$E(X_{22})$

$$E(X_{11}) = (1-p) \cdot 1635 = 1364.5$$

$$E(X_{12}) = (1-p) \cdot 1188 = 991.5$$

$$E(X_{21}) = p \cdot 1635 = 270.5$$

$$E(X_{22}) = p \cdot 1188 = 196.5$$

⇒ Vi skal bruke en χ^2 -test. Hvis kobling av prosjekt medler og kjønn er uavhengige, skal de observerte frekvenser passes bra til de forventede frekvenser, siden de har blitt kalkulert under den antakelsen av uavhengighet.

Testobservatoren blir:

$$d^2 = \sum_{i=1}^2 \sum_{j=1}^2 \frac{(X_{ij} - np_{ij})^2}{np_{ij}}$$

$$d_2 = \frac{(1345 - 1364.5)^2}{1364.5} + \frac{(1011 - 991.5)^2}{991.5} + \frac{(290 - 270.5)^2}{270.5} + \frac{(177 - 196.5)^2}{196.5}$$

$$d_2 = 4.00$$

Vi forkaster H_0 når d_2 er "stor." Vi må sammenligne med den 0.05 kvantil av en χ^2 -fordeling med 1 frihetsgrad $\chi^2_{0.05, 1} = 3.841$.

$\Rightarrow d_2 > 3.841$ Vi forkaster H_0 . \Rightarrow Data støtter at utdeling av prosjektmidler og kjønn er uavhengige.

b)

p_K : Sann andel av kvinner som får prosjektmidler

p_M : Sann andel av menn som får prosjektmidler

$\Rightarrow H_0: p_F = p_M$ $H_1: p_F \neq p_M$

Testobservator blir

$$z = \frac{\frac{177}{1188} - \frac{290}{1635}}{\sqrt{\frac{467}{2823} \left(1 - \frac{467}{2823}\right) \left(\frac{1}{1635} + \frac{1}{1188}\right)}}$$

$$= -2.00$$

(Teorem 9.4.1)

$$z = \frac{\frac{x}{n} - \frac{y}{m}}{p_e(1-p_e)\left(\frac{1}{n} + \frac{1}{m}\right)}$$

$$p_e = \frac{x+y}{n+m}$$

Med signifikansnivå $\alpha = 0.05$ blir kritiske verdier den 0.025 og 0.975 kvantil fra en standard normal fordeling som er -1.96 og 1.96 .

$\Rightarrow z < -1.96 \Rightarrow$ Vi forkaster H_0

Forskjellen mellom p_F og p_M er statistisk signifikant

c) Begge analyser forkaster H_0 . Det er ingen teoretisk grunn til å foretrekke en analyse over den annen. Både er ekvivalente, som mener at hvis χ^2 -testen forkaster H_0 at two egenskaper er uavhengige, må z -testen forkaste H_0 at de to binomiske andeler er like. Test observatoren til χ^2 test er den kvadraten av testobservatoren til z -test, og den kritisk verdi til χ^2 test er kvadraten av den kritisk verden til z -test.

$$\text{Her } d_2 = 4.00 = (2.00)^2 = (\text{obs. } z)$$

$$\text{og } \chi^2_{0.05, 1} = 3.84 = (\pm 1.96)^2 = z_{0.025}^2 = z_{0.975}^2$$

z -score i deloppgave c b) er den vert med bruk av en normal approximation.

Beräkna om ekvivalens:

	Man	Kvinna	
Ja	k_{11}	k_{12}	
Nei	k_{21}	k_{22}	
	n_1	n_2	n

Kontingenstabell

$$\Rightarrow d_2 = \frac{(k_{11} - \overbrace{p \cdot n_1}^{E(X_{11})})^2}{p \cdot n_1} + \frac{(k_{12} - \overbrace{p \cdot n_2}^{E(X_{12})})^2}{p \cdot n_2} + \frac{(k_{21} - \overbrace{q \cdot n_1}^{E(X_{21})})^2}{q \cdot n_1} + \frac{(k_{22} - \overbrace{q \cdot n_2}^{E(X_{22})})^2}{q \cdot n_2}$$

↑
Test för
oavhengighet

med $p = \frac{k_{11} + k_{12}}{n_1 + n_2}$ $q = 1 - p = \frac{k_{21} + k_{22}}{n_1 + n_2}$

$$d_2 = \frac{n_1^2 \left(\frac{k_{11}}{n_1} - p \right)^2}{n_1 p} + \frac{n_2^2 \left(\frac{k_{12}}{n_2} - p \right)^2}{n_2 p} + \frac{n_1^2 \left(\frac{k_{21}}{n_1} - q \right)^2}{n_1 q} + \frac{n_2^2 \left(\frac{k_{22}}{n_2} - q \right)^2}{n_2 q}$$

$$= n_1 \left[\frac{\left(\frac{k_{11}}{n_1} - p \right)^2}{p} + \frac{\left(\frac{k_{21}}{n_1} - q \right)^2}{q} \right] + n_2 \left[\frac{\left(\frac{k_{12}}{n_2} - p \right)^2}{p} + \frac{\left(\frac{k_{22}}{n_2} - q \right)^2}{q} \right]$$

$$q = (1-p)$$

$$= \frac{n_1 \left(\frac{k_{11}}{n_1} - p \right)^2 + n_2 \left(\frac{k_{12}}{n_2} - p \right)^2}{p \cdot q}$$

Beweis: $p = \frac{k_{11} + k_{12}}{n_1 + n_2}$

$$= \frac{n_1 \left(\frac{k_{11}}{n_1} - \frac{k_{11} + k_{12}}{n_1 + n_2} \right)^2 + n_2 \left(\frac{k_{12}}{n_2} - \frac{k_{11} + k_{12}}{n_1 + n_2} \right)^2}{p \cdot q}$$

$$= \frac{n_1 \cdot \left(\frac{k_{11}}{n_1} \cdot n_2 - n_2 \cdot \frac{k_{12}}{n_2} \right)^2 + n_2 \cdot \left(\frac{k_{12}}{n_2} \cdot n_1 + n_1 \cdot \frac{k_{11}}{n_1} \right)^2}{p \cdot q \cdot (n_1 + n_2)^2}$$

$$= \frac{n_1 \cdot n_2^2 \left(\frac{k_{11}}{n_1} - \frac{k_{12}}{n_2} \right)^2 + n_2 \cdot n_1^2 \left(\frac{k_{12}}{n_2} - \frac{k_{11}}{n_1} \right)^2}{p \cdot q \cdot (n_1 + n_2)^2}$$

$$= \frac{(n_1 \cdot n_2^2 + n_2 \cdot n_1^2) \left(\frac{k_{11}}{n_1} - \frac{k_{12}}{n_2} \right)^2}{p \cdot q \cdot (n_1 + n_2)^2}$$

$$= \frac{\left(\frac{k_{11}}{n_1} - \frac{k_{12}}{n_2} \right)^2}{p \cdot q \cdot \frac{(n_1 + n_2)^2}{n_1 n_2 (n_1 + n_2)}} = \frac{\left(\frac{k_{11}}{n_1} - \frac{k_{12}}{n_2} \right)^2}{p \cdot (1-p) \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}$$

$$= z^2$$

\Rightarrow Testobservatoren av χ^2 testen for uavhengighet er den kvadrat av testobservatoren vi bruker i z-testen for å sammenligne to binomiske andeler.

Dessuten hvis $X \sim N(0,1)$

$$\Rightarrow X^2 \sim \chi^2_1$$

Verdren til

\Rightarrow Teststatistikk til χ^2 testen er den kvadrat av verdren til teststatistikk i z-test.

2) Integralberegning:

$$f_{T_n}(t) = \frac{\Gamma(\frac{n+1}{2})}{\sqrt{n\pi} \Gamma(\frac{n}{2}) (1 + \frac{t^2}{n})^{(n+1)/2}}$$

med $\Gamma(1) = 1$ og $\Gamma(\frac{1}{2}) = \sqrt{\pi}$

Mål: Evaluer $\int_0^{\infty} \frac{1}{1+t^2} dt$

Vi vet at $\int_{-\infty}^{\infty} f_{T_n}(t) dt \stackrel{!}{=} 1$ for alle n .

Brek $n=1$

$$\Rightarrow f_{T_1}(t) = \frac{1}{\sqrt{\pi} \sqrt{\pi} (1+t^2)} = \frac{1}{\pi (1+t^2)}$$

$$\Rightarrow \int_{-\infty}^{\infty} \frac{1}{\pi(1+t^2)} dt = 1$$

$$\Rightarrow \int_{-\infty}^{\infty} \frac{1}{1+t^2} dt = \pi$$

Vi vet at en Student t fordeling er symmetrisk:

$$\Rightarrow 2 \cdot \int_0^{\infty} \frac{1}{1+t^2} dt = \pi$$

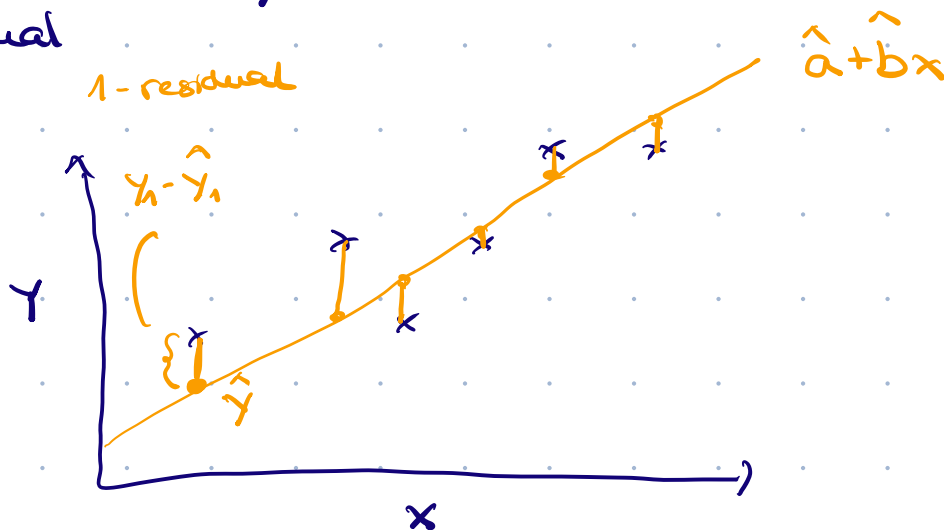
$$\Rightarrow \int_0^{\infty} \frac{1}{1+t^2} dt = \frac{\pi}{2} \underline{\underline{=}}$$

Enkelt lineær model

a) Antag at \hat{a} og \hat{b} er de mindste kvadrat estimatorene til en sæt af observationerne $(x_1, y_1), \dots, (x_n, y_n)$.

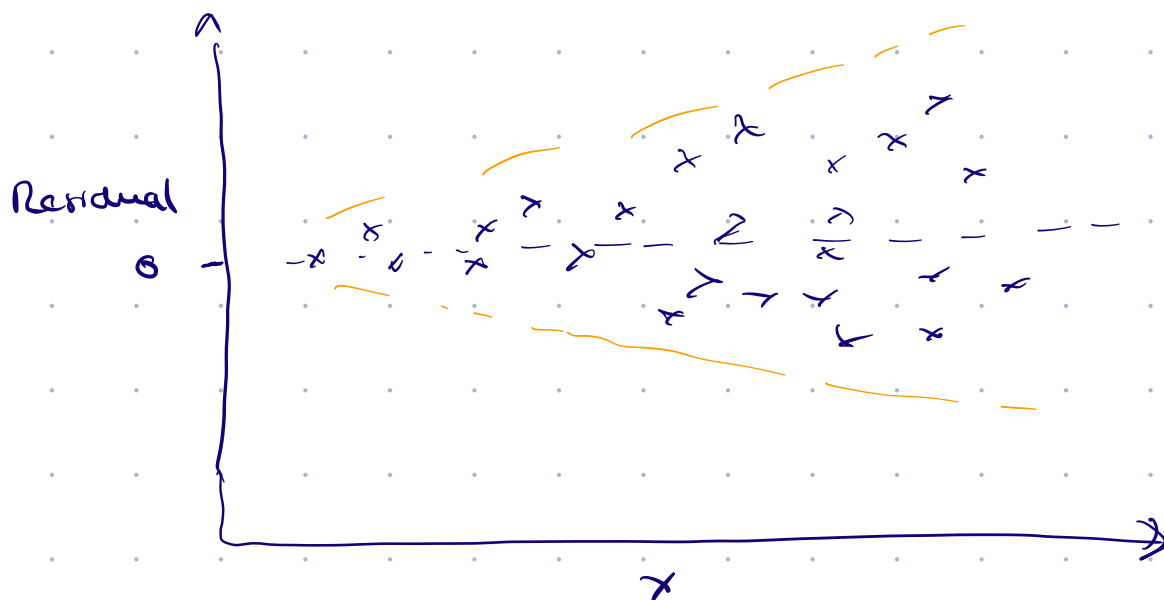
$$\hat{y} = \hat{a} + \hat{b}x \quad \text{kaldes ofte "fitted value"}$$

Differensen $y_i - \hat{y}_i = y_i - (\hat{a} + \hat{b}x_i)$ kaldes residual

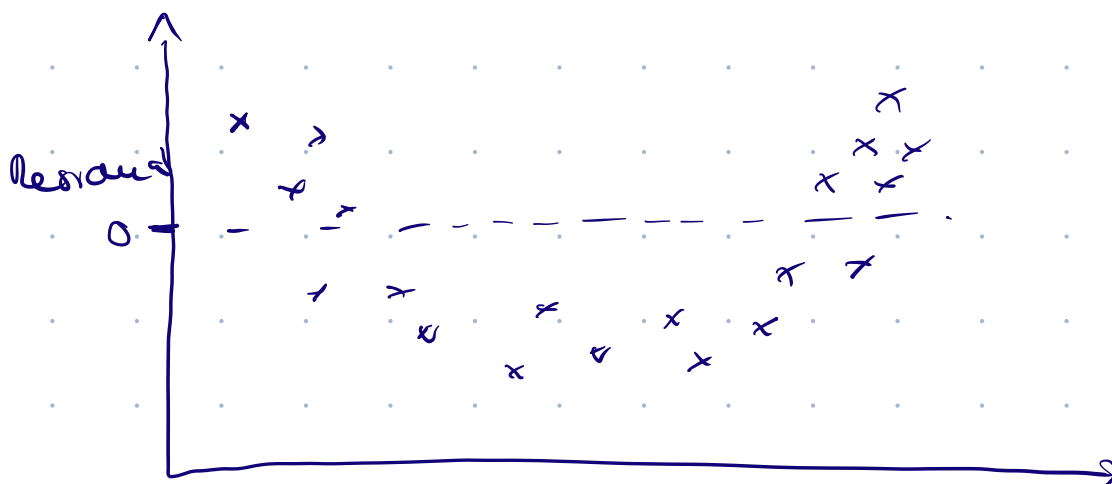


Det betyder at differensen mellem den observerede værdi y_i og den værdi \hat{y}_i som ligger på regressionslinjen når $x = x_i$ kaldes the i -th residual.

Her to eksempler af residual plots som viser at en egenskab er en lineær model er brudt.



Her ser man mere varians for større verdier av x
 men varians skal bli den samme for alle x
 (homoskedastisk)



Den residualplott har et klart mønster. Det virker
 sann at det er ikke en linear sammenheng mellom
 x og y .

b) Vi må finne likelihood funksjonen og sett $\hat{\beta}$ til β

$$L(\beta) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi} \sigma^2} \exp\left(-\frac{(Y_i - \beta x_i)^2}{2\sigma^2}\right)$$

$$\ln L(\beta) = -\frac{n}{2} \ln(2\pi) - \frac{n}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (Y_i - \beta x_i)^2$$

Deriverer

$$\frac{\partial}{\partial \beta} \ln L(\beta, \sigma^2) = \frac{1}{\sigma^2} \sum_{i=1}^n (Y_i - \beta x_i) \cdot (-x_i) \stackrel{!}{=} 0$$

$$-\sum_{i=1}^n (Y_i x_i) + \beta \sum_{i=1}^n x_i^2 \stackrel{!}{=} 0 \Rightarrow \hat{\beta} = \frac{\sum_{i=1}^n Y_i x_i}{\sum_{i=1}^n x_i^2}$$

$$\frac{\partial}{\partial \sigma^2} \ln L(\beta, \sigma^2) = -\frac{n}{2\sigma^2} + \frac{\sum_{i=1}^n (Y_i - \beta x_i)^2}{2\sigma^4} \stackrel{!}{=} 0$$

$$\Rightarrow \frac{1}{\sigma^2} \sum_{i=1}^n (Y_i - \beta x_i)^2 = n \Rightarrow \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{\beta} x_i)^2$$

$\hat{\beta}$ er en linear kombinasjon av uavhengige normal-

fordelte stokastiske variabler Y_i , $i=1, \dots, n$ og

er derfor normalfordelt:

$$E(\hat{\beta}) = \frac{\sum_{i=1}^n E(Y_i) x_i}{\sum_{i=1}^n x_i^2} = \frac{\sum_{i=1}^n \beta x_i^2}{\sum_{i=1}^n x_i^2} = \beta \Rightarrow \text{forventningsrett}$$

$$\text{Var}(\hat{\beta}) = \text{Var}\left(\frac{\sum Y_i x_i}{\sum x_i^2}\right) = \frac{1}{(\sum x_i^2)^2} \sum x_i^2 \text{Var}(Y_i) = \frac{\sigma^2}{\sum x_i^2}$$

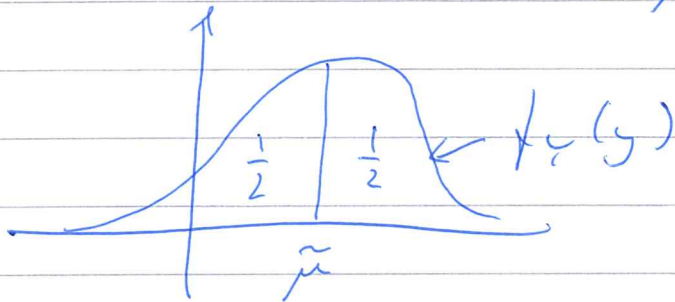
$$\Rightarrow \hat{\beta} \sim N\left(\beta, \frac{\sigma^2}{\sum x_i^2}\right)$$

Løsning

Pris på fjellhytter

a) Medianen er slik at

$$\int_{-\infty}^{\tilde{\mu}} f_Y(y) dy = \int_{\tilde{\mu}}^{\infty} f_Y(y) dy = \frac{1}{2}$$



Her er:

$$H_0: \tilde{\mu} = 2500 \text{ mot } H_1: \tilde{\mu} > 2500$$

La X = antall hytter med salgpris > 2500 .

Under H_0 er $X \sim \text{bin}(12, 0.5)$

H_0 forkastes hvis X er stor, og p -verdien blir da siden det observeres $X = 9$,

$$p\text{-verdi} = P(X \geq 9) = 1 - P(X \leq 8) \stackrel{\text{tabell}}{=} 1 - 0.927 = \underline{0.073}$$

H_0 forkastes dermed ikke siden signifikansen er 5%.

b) Y-2500:

	-510	-110	-50	80	90	110	175	220	590	940
Rang:	8	4.5	1	2	3	4.5	6	7	9	10
	1900	3000								
Rang:	11	12								

$$W_+ = \text{sum pos. ranger} = 2+3+\dots+12 = 64.5$$

$$W_- = \text{" neg " } = 8+4.5+1 = 13.5$$

Kan forkaste H_0 hvis W_- er liten
~~p-verdi~~ Tabell gir

$$P(W_- \leq 17) = 0.046$$

$$P(W_- \leq 18) = 0.055$$

des 5% test forkaster H_0 hvis $W_- \leq 17$.

Vi har $W_- = 13.5$ så H_0 forkastes.

$$p\text{-verdi} = P(W_- \leq 13.5) \approx \frac{0.021 + 0.026}{2} = 0.0235$$

(NB: Tabellen er korrekt

bare for obs. uten "ties", men dette er OK som en tilnærning).

Det er også mulig å bruke tilnærning til normalfordeling:

$$E(W_-) (= E(W_+)) = \frac{n(n+1)}{4} = \frac{12 \cdot 13}{4} = 39$$

$$\text{Var}(W_-) (= \text{Var}(W_+)) = \frac{n(n+1)(2n+1)}{24} = \frac{12 \cdot 13 \cdot 25}{24} = 162.5$$

-3-

$$p\text{-verdi} = P(W_- \leq 13.5) = \Phi\left(\frac{13.5 - 39}{\sqrt{162.5}}\right)$$

$$= \Phi(-2.00) = 0.0228$$

(stemmer godt med den vi fant)

c) La μ_1, μ_2, μ_3 være forventet salgpris for henholdsvis Oppdal, Hafjell og Tysil.

Skal teste $H_0: \mu_1 = \mu_2 = \mu_3$ mot $H_1: \text{ikke alle like}$

forts.

c) Ranger alle observasjoner:

Oppdal		Hafjell		Trysil	
1990	1	2660	8	2070	2
2390	3	2810	12	2730	10
2450	4	2880	13	3080	16
2580	5	2900	14	3150	17
2590	6	2950	15	3230	19
2610	7	3290	20	3370	21
2675	9	4830	26	3620	23
2770	11	5320	27	4125	24
3190	18				
3440	22	$R_{.2} = 135$		$R_{.3} = 132$	
4400	25				
5500	28				

$$R_{.1} = 139$$

$$\text{Sjekk: } 139 + 135 + 132 = 406 = \frac{28 \cdot 29}{2}$$

$$B = \frac{12}{28 \cdot 29} \left(\frac{139^2}{12} + \frac{135^2}{8} + \frac{132^2}{8} \right) - 3 \cdot (28 + 1)$$

$$= 2.65$$

Under H_0 er $B \sim \chi_{3-1}^2 = \chi_2^2$. Med 5% sign. nivå forkastes hvis $B \geq 5.991$. Alltså forkastes ikke H_0 .

d) Her kunne Wilcoxon's to-utvalgstest brukes. Antagelsene ville være de samme som for Kruskal-Wallis testen. Testobservatoren ville være summen av rangene for de to ene av områdene, og det ville forkastes ved høye eller lave verdier for denne.

Hvis situasjonen er som i c), men at observasjonene er normalfordelte med samme varians, ville vi kunne bruke en-faktor design som i kap 12 i boka. Testobservatoren er da F-fordelt under H_0 (se boka).