

Statistik 1 för biologer, logopeders och psykologer

Föreläsningar, del 6

Innehåll

- 1 Analys av korstabeller
- 2 Variansanalys och försöksplanering

Innehåll

- 1 Analys av korstabeller
- 2 Variansanalys och försöksplanering

Korstabeller

- Vi har tidigare under kursen redan bekantat oss med korstabeller.
- I en korstabell redovisar man fördelningen på två eller flere, vanligen kvalitativa variabler.
- Också om man har att göra med i princip kvantitativa variabler kan det ibland vara skäl att övergå till att studera enbart fördelningen på olika klasser.

Korstabeller

Exempel på en korstabell:

Rökvanor	Socioekonomisk status			Summa
	Hög	Medel	Låg	
Röker	51	22	43	116
Rökt tidigare	92	21	28	141
Aldrig rökt	68	9	22	99
Summa	211	52	93	356

Test av oberoende med χ^2 -test

- Ett liknande χ^2 -test som användes för test av fördelning kan även användas för att testa *oberoende* av två korstabulerade variabler.
- Värdet på testvariabeln beräknas enligt formeln

$$\chi^2 = \sum_{i=1}^k \sum_{j=1}^l \frac{(O_{ij} - E_{ij})^2}{E_{ij}},$$

där O_{ij} = observerad cellfrekvens och E_{ij} = förväntad cellfrekvens under antagande att H_0 är sann. Summeringen sker över de k raderna och l kolumnerna i korstabellen.

- Då H_0 är sann, följer testvariabeln en χ^2 -fördelning med frihetsgraderna $df = (k - 1)(l - 1)$.
- H_0 förkastas enligt samma principer som i det tidigare introducerade χ^2 -testet.

Test av oberoende med χ^2 -test

- Hypoteserna i ett χ^2 -oberoendetest är

H_0 : Inget samband, dvs. variablerna oberoende av varandra.

H_1 : Det finns ett samband.

- Vi låter p_{ij} beteckna sannolikheten för att en observation tillhör en viss cell i korstabellen. Om variablerna är *oberoende* blir sannolikheten för att en observation tillhör en viss cell

$$\begin{aligned} p_{ij} &= P(\text{observationen tillhör rad } i \text{ och kolumn } j) \\ &= P(\text{observationen tillhör rad } i) \cdot P(\text{observationen tillhör kolumn } j) . \end{aligned}$$

- De *förväntade* cellfrekvenserna räknas alltså

$$E_{ij} = \frac{(\text{summan av rad } i) \cdot (\text{summan av kolumn } j)}{n} .$$

Test av oberoende med χ^2 -test – exempel

Exempel.

- För att undersöka effekten av ett nytt vaccin på en sjukdom ges 70 frivilliga försökspersoner vaccinet. I undersökningen ingår även en lika stor kontrollgrupp.
- De sammanlagt 140 personerna följs upp under en viss tid och man erhåller följande resultat:

	Har insjuknat	Har ej insjuknat	Tot.
Har vaccinerats	20	50	70
Har ej vaccinerats	40	30	70
Tot.	60	80	140

- Följande hypoteser formuleras:

H_0 : Vaccinet har ingen effekt.

H_1 : Vaccinet förebygger sjukdomen.

Test av oberoende med χ^2 -test – exempel (forts.)

Exempel.

- Om vaccin och sjukdom är oberoende skulle vi förvänta oss följande:

	Har insjuknat	Har ej insjuknat	Tot.
Har vaccinerats	$\frac{70 \cdot 60}{140} = 30$	$\frac{70 \cdot 80}{140} = 40$	70
Har ej vaccinerats	$\frac{70 \cdot 60}{140} = 30$	$\frac{70 \cdot 80}{140} = 40$	70
Tot.	60	80	140

- Från de två korstabellerna räknar vi sedan värdet på testvariabeln

$$\chi^2 = \frac{(20 - 30)^2}{30} + \frac{(40 - 30)^2}{30} + \frac{(50 - 40)^2}{40} + \frac{(30 - 40)^2}{40} = 11.7 .$$

- Med signifikansnivån $\alpha = 0.01$ och frihetsgraderna

$df = (2 - 1)(2 - 1) = 1$ förkastar vi H_0 eftersom

$$\chi^2 = 11.7 > \chi_{\alpha}^2 = 6.635.$$

Vidare om korstabeller: betingat oberoende och Simpson's paradox

- Två variabler som till synes verkar beroende kan vara oberoende om man tar hänsyn till en tredje variabel. Detta kallas **betingat oberoende**.
- I följande introduktion till analys av korstabeller behandlas även betingat oberoende:

[http://web.abo.fi/fak/mnf/mate/kurser/statistik1/
AnalysAvKorstabeller.pdf](http://web.abo.fi/fak/mnf/mate/kurser/statistik1/AnalysAvKorstabeller.pdf).

- Texten tar även upp det sk. **Simpson's paradoxet**. För definition och flera exempel, se

http://en.wikipedia.org/wiki/Simpson's_paradox.

- Allmänt om betingat oberoende:

[http://web.abo.fi/fak/mnf/mate/kurser/statistik1/
MarginelltBetingat.pdf](http://web.abo.fi/fak/mnf/mate/kurser/statistik1/MarginelltBetingat.pdf)

Oddsquot

- Ett mått som ibland används för att beskriva graden av ett samband mellan två korstabulerade *dikotoma* (=av typen ja/nej) variabler är **oddsquoten** (förkortas ofta *OR* från *odds ratio*).
- Vi definierar först **odds** för händelsen A :

$$\frac{P(A \text{ inträffar})}{P(A \text{ inträffar ej})} = \frac{P(A)}{P(A^c)} = \frac{P(A)}{1 - P(A)} .$$

- Om det är lika sannolikt att händelsen A inträffar som att den inte inträffar får oddsquoten värdet 1.
- Odds kan även beräknas för betingade sannolikheter. T.ex. räknas oddsquoten för att insjukna i en viss sjukdom givet att man har vaccinerats:

$$\frac{P(\text{Sjuk} | \text{Har vaccinerats})}{P(\text{Ej sjuk} | \text{Har vaccinerats})} = \frac{P(\text{Sjuk} | \text{Har vaccinerats})}{1 - P(\text{Sjuk} | \text{Har vaccinerats})} .$$

Oddskvot

- Oddskvoten definieras som kvoten mellan två odds Vi kan då t.ex. räkna

$$OR = \frac{\frac{P(\text{Sjuk}|\text{Har vaccinerats})}{P(\text{Ej sjuk}|\text{Har vaccinerats})}}{\frac{P(\text{Sjuk}|\text{Har ej vaccinerats})}{P(\text{Ej sjuk}|\text{Har ej vaccinerats})}}$$

vilket talar om för oss hur stort oddset för att insjukna är då man blivit vaccinerad i förhållande till motsvarande odds då man ej blivit vaccinerad.

- Resultatet tolkas på följande sätt
 - $OR < 1$: vaccinering *minskar* oddset för att insjukna
 - $OR = 1$: vaccinering *påverkar inte* oddset för att insjukna
 - $OR > 1$: vaccinering *ökar* oddset för att insjukna.

Oddsquot och relativ risk

- Ett mått som är ganska likt oddsquoten är den sk. **relativa risken** (förkortas ofta *RR*), vilken i fallet ovan räknas som

$$RR = \frac{P(\text{Sjuk}|\text{Har vaccinerats})}{P(\text{Sjuk}|\text{Har ej vaccinerats})}$$

- För små värden är $OR \approx RR$.
- Fastän *RR* är något enklare och mera intuitiv än *OR* är den senare ofta mera användbar i statistiska analyser.
- Mera om oddsquoten och dess egenskaper:
<http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=1127651>

Innehåll

- 1 Analys av korstabeller
- 2 Variansanalys och försöksplanering

Variansanalys

- Vi har tidigare i form av ett t -test bekantat oss med jämförelse av väntevärden från två normalfördelade populationer med lika standardavvikelse.
- **Variansanalys** (förkortas ofta ANOVA från *analysis of variance*) är en generalisering av det ovannämnda testet för två eller flera väntevärden.
- Namnet kan te sig något vilseledande då det ju är väntevärden man testar. Det har dock sitt ursprung i att den sk. F -testvariabeln som metoden bygger på kan tolkas som en kvot av varianser bildade på två olika sätt.

Variansanalys

- I variansanalys antar vi alltså att vi har k oberoende (möjligtvis olika stora) stickprov från lika många *normalfördelade* populationer med lika standardavvikelse, dvs. $\sigma_1 = \sigma_2 = \dots = \sigma_k = \sigma$.
- Hypoteserna kan formuleras på följande sätt:

$$H_0 : \mu_1 = \mu_2 = \dots = \mu_k$$

H_1 : Åtminstone ett av väntevärdena skiljer sig från de andra.

- F -testvariabeln är en kvot av spridningen *mellan* och spridningen *inom* de k grupperna. Om H_0 är sann följer variabeln en F -fördelning med frihetsgraderna $df = (k - 1, n - k)$, där n är det totala antalet observationer.
- H_0 förkastas på signifikansnivån α om testvariabelns värde överskrider det kritiska värdet f_α i en F -fördelning med frihetsgraderna $(k - 1, n - k)$.

Variansanalys – exempel

Exempel.

Vi har fyra stickprov från normalt fördelade populationer. Av tabellen nedan framgår storleken, medelvärdet och variansen för respektive stickprov:

n_i	7	5	6	6
\bar{x}_i	2.4	3.3	3.4	2.6
s_i^2	0.25	0.34	0.10	0.05

Enligt H_0 har alla populationer samma väntevärde. Om det nu visar sig att spridningen mellan stickproven är stor jämfört med spridningen inom stickproven har vi orsak att tro att H_0 är falsk. Spridningen mellan stickproven är 1.53 medan den genomsnittliga spridningen inom stickproven är 0.18 (vi hoppar här över uträkningarna). F -testvariabeln får då värdet $1.53/0.18 = 8.66$. Det kritiska värdet för en F -fördelning med frihetsgraderna $(4 - 1, 24 - 4)$ på signifikansnivån $\alpha = 0.05$ är 3.10. Eftersom $8.66 > 3.10$ förkastar vi H_0 .

Experimentella försök

- Variansanalysen är en metod som primärt utvecklats för att analysera resultat av **experimentella försök**.
- I ett typiskt experimentellt försök jämför man effekten av olika behandlingar, vilket här ska uppfattas som en allmän benämning på något som man utsätter försöksenheter för.
- Till skillnad från ett **icke-experimentellt försök** kan man i ett experimentellt försök påverka vilka enheter som får vilken behandling.
- Ett väl utfört experimentellt försök medger säkrare slutsatser än ett icke-experimentellt försök.

Exempel på experimentella försök

Exempel.

- För att jämföra två medicinska preparat A och B ger en läkare det ena till en patientgrupp och det andra till en annan patientgrupp och jämför resultaten.
- Ett företag överväger att ersätta nuvarande tillverkningsmetod A med en ny metod B. För att undersöka om den nya metoden är bättre än den gamla tillverkar man ett antal enheter enligt vardera metoden och jämför resultatet.
- För att jämföra tre olika undervisningsmetoder delas eleverna i en årskurs i början av terminen slumpmässigt in i tre grupper. I slutet av terminen jämför man den genomsnittliga utvecklingen bland eleverna i de tre grupperna.

Flervägs variansanalys

- Den typ av variansanalys vi ovan har diskuterat kallas för **envägs variansanalys** eftersom indelningen i grupper sker enligt *en* typ av behandling (t.ex. medicinskt preparat, tillverkningsmetod, undervisningsmetod...).
- Analyserar man kombinationer av *flera* typers behandlingar använder man sk. **flervägs variansanalys**.
- Mera om variansanalys, se http://en.wikipedia.org/wiki/Analysis_of_variance

Försöksplanering

- Med **försöksplanering** strävar man efter att kontrollera effekten av faktorer som kan påverka tillförlitligheten av statistiska analyser i samband med ett experimentellt försök.
- Några för försöksplanering centrala begrepp och tekniker är:
 - **Randomisering**, dvs. slumpmässig allokering av försöksenheter i olika behandlingsgrupper för att minska på inverkan av okända systematiska fel på slutsatserna.
 - **Replikering**, vilket betyder att man gör upprepade mätningar av samma enhet för att få en uppfattning av mätfelet.
 - Indelning av liknande enheter i **block** för uppnå bättre precision ifall det finns stor variation bland enheterna.

Försöksplaner

- Ett par vanliga **försöksplaner** är:
 - Fullständigt randomiserat experiment (CRD, *Completely Randomized Design*), se t.ex.
http://courses.ncssm.edu/math/Stat_Inst/PDFS/RanDesgn.pdf
 - Randomiserat blockexperiment (RBD, *Randomized Block Design*), se t.ex.
http://en.wikipedia.org/wiki/Randomized_block_design