

Föreläsning 9.

* Miniprojekt 2 - deadline 24/5 kl. 16.⁰⁰

+ Uppgift a-d efter laboration 3.

+ Uppgift e efter laboration 4.

* Quizto 2 - deadline 28/5 kl 23⁵⁹.

* Tenta 3/6 08⁰⁰-13⁰⁰; Sparta: D

Outline

* Repetition

* Icke-parametriska test

* Inferens om proportioner

* Analys av kategoridata

Repetition

Utgångspunkt: $\bar{X} \in N(\mu_{\bar{X}}, \sigma_{\bar{X}}^2)$

$H_0: \mu_x - \mu_y = 0$ $\bar{Y} \in N(\mu_{\bar{Y}}, \sigma_{\bar{Y}}^2)$

$H_1: \mu_x - \mu_y \neq 0$ (eller > 0 eller < 0)

1) Kan man se \bar{X}_k och \bar{Y}_k som parvisa mätningar?

Ja! • Låt $Z_k = \bar{X}_k - \bar{Y}_k$ (eller $\bar{Y}_k - \bar{X}_k$)

$Z_k \in N(\Delta, \sigma_z^2)$

- Sätt upp hypotes.

- Testa. Om σ_z känd, använd normalkvantil $Z_{1-\alpha}$ ($Z_{1-\alpha/2}$ om dubbelsidigt)

Om σ_z okänd skatta med S_z , använd t-kvantil $t_{1-\alpha, n-1}$ ($t_{1-\alpha/2, n-1}$ om dubbelsidigt).

Repetition

Nej! Alltså inte parvisa!

- σ_x och σ_y kända:

Testa med normalkvantil och med varians

$$\frac{\sigma_x^2}{n_x} + \frac{\sigma_y^2}{n_y}$$

- σ_x och σ_y okända:

- Kan man anta att $\sigma_x = \sigma_y$?

Ja! Testa med t-kvantil med

$n_x + n_y - 2$ frihetsgrader och med

varians: $S^2 \left(\frac{1}{n_x} + \frac{1}{n_y} \right)$

$$\text{där } S^2 = \frac{(n_x - 1)S_x^2 + (n_y - 1)S_y^2}{n_x + n_y - 2}$$

Repetition

Nej! σ_x och σ_y antas vara olika.
(överkurs)

- Testa med F-kvantil med

$$f = \frac{\left[\frac{s_x^2}{n_x} + \frac{s_y^2}{n_y} \right]^2}{\frac{\left[\frac{s_x^2}{n_x} \right]^2}{n_x - 1} + \frac{\left[\frac{s_y^2}{n_y} \right]^2}{n_y - 1}} \quad \text{frihetsgrader}$$

och varians $\frac{s_x^2}{n_x} + \frac{s_y^2}{n_y}$.

Repetition

Hur ser konfidensintervallen ut för de olika scenarierna? Antag att variansen är okänd (annars använd normalkventil), konfidensnivå α , samt att vi testar:

$$H_0: \mu_x - \mu_y = 0$$

$$H_1: \mu_x - \mu_y \neq 0$$

Parvis stickprov: $Z_{ik} = X_{ik} - Y_{ik} \sim N(\Delta, \sigma_z^2)$

$$I_\alpha = \left(\bar{Z} - t_{1-\alpha/2, n-1} \cdot \frac{s_z}{\sqrt{n}}, \bar{Z} + t_{1-\alpha/2, n-1} \cdot \frac{s_z}{\sqrt{n}} \right)$$

TVå oberoende stickprov:

Antag $\sigma_x = \sigma_y$

$$I_{\mu_x - \mu_y} = \left(\bar{x} - \bar{y} \pm t_{1-\alpha/2, n_x + n_y - 2} s \sqrt{\frac{1}{n_x} + \frac{1}{n_y}} \right)$$

$$\text{med } s^2 = \frac{(n_x - 1)S_x^2 + (n_y - 1)S_y^2}{n_x + n_y - 2}$$

Repetition

Två oberoende stickprov: (Överkurs)

Antag $\sigma_x \neq \sigma_y$

$$I_{\mu_x - \mu_y} = (\bar{x} - \bar{y} \pm t_{1-\alpha/2, f} \sqrt{\frac{s_x^2}{n_x} + \frac{s_y^2}{n_y}})$$

med

$$f = \frac{\left[\frac{s_x^2}{n_x} + \frac{s_y^2}{n_y} \right]^2}{\frac{\left[\frac{s_x^2}{n_x} \right]^2}{n_x - 1} + \frac{\left[\frac{s_y^2}{n_y} \right]^2}{n_y - 1}} \quad \text{frihetsgrader}$$

Repetition

* Ett bra sätt att testa om $\sigma_X = \sigma_Y$ är att ställa upp följande hypotes och testa:

$$H_0: \sigma_X^2 = \sigma_Y^2$$

$$H_1: \sigma_X^2 \neq \sigma_Y^2$$

Man kan visa att $\frac{S_X^2}{S_Y^2}$ tillhör F-fördelningen, dvs.

$$\frac{S_X^2}{S_Y^2} \in F_{(n_X-1), (n_Y-1)}$$

Vi förkastar H_0 om $\frac{S_X^2}{S_Y^2} > \bar{F}_{1-\alpha/2, (n_X-1), (n_Y-1)}$

Icke-parametriska test.

- * Om vi vill testa en hypotes men inte kan göra ett normalfördelningsantagande, hur gör vi då?
- * Istället för att titta på de verkliga värdena så tittar vi på tecknet (+/-)
- * Hypotesen blir då:
 - H_0 : medianen är 0
 - H_1 : medianen är (större, mindre) skild från 0.
- * Löser det med hjälp av p-värdesmetoden och att vi låter $X \in \text{Bin}(n, 0.5)$.

 BIostatistisk Grundkurs, MASB11
 OH-BILDER 8

ICKEPARAMETRISKA TEST:

EXEMPEL: Man ville undersöka om energiförbrukning i vila är annorlunda hos personer drabbade av cystisk fibros jämfört med friska personer. Tio par matchades ihop, en i varje par hade sjukdomen medan den andre var frisk. För övrigt var personerna i varje par lika beträffande kön, ålder, vikt och längd. Resultat i energiförbrukning (kcal/dag):

Par	1	2	3	4	5	6	7	8	9	10
CF	1153	1132	1165	1460	1634	1493	1358	1453	1185	1824
Friska	996	1080	1182	1452	1162	1619	1140	1123	1113	1463
Skillnad										
CF-friska	157	52	-17	8	472	-126	218	330	72	361

- (a) Hur gör vi om det är rimligt att tänka sig att differenserna kommer från en normalfördelning?
- (b) Hur gör vi utan antaganden om normalfördelning?

EXEMPEL: I en skola ville man göra en liten pilotstudie för att se om en annorlunda idrottsträning på kort tid skulle kunna påverka skolbarnens fysiska prestationer. Man valde ut 16 barn, som var likvärdiga beträffande den fysiska kapaciteten. Barnen delades slumpmässigt in i två grupper. Under en månad följde hälften av barnen (grupp A) den normala undervisningen i ämnet Idrott och hälsa, medan de övriga barnen (grupp B) dessutom fick delta i den speciella träningen. När en månad hade gått, fick barnen vid ett gemensamt tillfälle springa en kort terrängbana och deras tider noterades. Två barn i grupp A var sjuka under testdagen. Resultat (sekunder):

Grupp A	64	62	73	54	66	71		
Grupp B	53	74	70	59	42	38	48	60

SITUATION	NORMALFÖRDELNING	ICKE-PARAMETRISKT
Ett stickprov	t-test	teckentest
Två stickprov, matchade data	t-test på differenserna	teckentest på differenser
Två stickprov, oberoende data	t-test för två ober stickprov	Wilcoxon-Mann-Whitney's test
Flera stickprov, oberoende data	ensidig variansanalys	Kruskal-Wallis test

FÖR- OCH NACKDELAR MED ICKE-PARAMETRISKA TEST

- (+) Behöver inte göra antaganden om fördelning hos data
- (+) Fungerar för små stickprov
- (+) "Robust" mot outliers — d.v.s. påverkas inte så mycket av kraftigt avvikande värden i ett datamaterial

- (-) Är inte lika "känsliga" som (har mindre styrka än) de test som baseras på normalfördelning — d.v.s. det behövs ett större stickprov för att förkasta en felaktig H_0
- (-) Nollhypotesen är oftast inte lika specificerad som i "traditionella" test
- (-) Utnyttjar inte all information om fördelningen som ges i data — baseras oftast på ranger, inte på de aktuella värdena i mätningarna

Interens om proportioner

* En proportion anger hur stor andel av datan som innehar en viss egenskap.

T.ex. Andelen som röstade i senaste valet eller andelen elever i kursen MASB11 som klarade mozquito 1 på första försöket.

$$\hat{p} = \frac{x}{n}$$

(85.8% röstade i valet 2014 och 10% klarade mozquito 1 vid första försöket.)

Interens om proportioner

* I de flesta fallen är $X \in \text{Bin}(n, p)$
och p skattas som $p = \frac{\bar{X}}{n}$.

Interens om proportioner

* I de flesta fallen är $X \in \text{Bin}(n, p)$
och p skattas som $p = \frac{X}{n}$.

* Om vi utnyttjar detta, då kan vi
skapa konfidensintervall för p
samt göra hypotestest för p .

Interens om proportioner

* I de flesta fallen är $X \in \text{Bin}(n, p)$
och p skattas som $p = \frac{X}{n}$.

* Om vi utnyttjar detta, då kan vi
skapa konfidensintervall för p
samt göra hypotestest för p .

Kom ihåg att om $X \in \text{Bin}(n, p)$
då är $E[X] = np$ och $\text{Var}(X) = np(1-p)$

Interens om proportioner

* I de flesta fallen är $X \in \text{Bin}(n, p)$
och p skattas som $p = \frac{X}{n}$.

* Om vi utnyttjar detta, då kan vi
skapa konfidensintervall för p
samt göra hypotestest för p .

Kom ihåg att om $X \in \text{Bin}(n, p)$
då är $E[X] = np$ och $\text{Var}(X) = np(1-p)$

$$\hat{p} = \frac{X}{n} :$$

$$E[\hat{p}] = E\left[\frac{X}{n}\right] = \frac{E[X]}{n} = \frac{np}{n} = p$$

$$\text{Var}[\hat{p}] = \text{Var}\left[\frac{X}{n}\right] = \frac{np(1-p)}{n^2} = \frac{p(1-p)}{n}$$

Interens om proportioner

• Så om $X \in \text{Bin}(n, p)$ då är

$$E[\hat{p}] = p, \quad \text{Var}[\hat{p}] = \frac{p(1-p)}{n}$$

• En rimlig skattning av p är $\frac{x}{n} = \hat{p}_{\text{obs}}$.

Alltså:

$$E[\hat{p}] = \hat{p}_{\text{obs}}, \quad \text{Var}[\hat{p}] = \frac{\hat{p}_{\text{obs}}(1-\hat{p}_{\text{obs}})}{n}$$

Interens om proportioner

• Så om $X \in \text{Bin}(n, p)$ då är

$$E[\hat{p}] = p, \quad \text{Var}[\hat{p}] = \frac{p(1-p)}{n}$$

• En rimlig skattning av p är $\frac{x}{n} = \hat{p}_{\text{obs}}$.

Alltså:

$$E[\hat{p}] = \hat{p}_{\text{obs}}, \quad \text{Var}[\hat{p}] = \frac{\hat{p}_{\text{obs}}(1-\hat{p}_{\text{obs}})}{n}$$

• Om $np(1-p) \geq 10$ då kan vi

normalapproximera $X \stackrel{\text{approx}}{\sim} N(np, np(1-p))$

vilket ger $p \stackrel{\text{approx}}{\sim} N\left(p, \frac{p(1-p)}{n}\right)$.

Interens om proportioner

• Så om $X \in \text{Bin}(n, p)$ då är

$$E[\hat{p}] = p, \quad \text{Var}[\hat{p}] = \frac{p(1-p)}{n}$$

• En rimlig skattning av p är $\frac{x}{n} = \hat{p}_{\text{obs}}$.

Alltså:

$$E[\hat{p}] = \hat{p}_{\text{obs}}, \quad \text{Var}[\hat{p}] = \frac{\hat{p}_{\text{obs}}(1-\hat{p}_{\text{obs}})}{n}$$

• Om $np(1-p) \geq 10$ då kan vi

normalapproximera $X \stackrel{\text{d}}{\sim} N(np, np(1-p))$

vilket ger $\hat{p} \stackrel{\text{d}}{\sim} N\left(p, \frac{p(1-p)}{n}\right)$.

• Ett konfidensintervall för p blir då

$$I_p = \left(\hat{p}_{\text{obs}} \pm z_{1-\alpha/2} \sqrt{\frac{\hat{p}_{\text{obs}}(1-\hat{p}_{\text{obs}})}{n}} \right)$$

OBS! Ej t -fördelningen; vi har normalapprox!

Interens om proportioner

* Vi kan också göra hypotestest.

$$H_0: p = p_0$$

$$H_1: p > p_0$$

* Vi kan testa hypotesen genom:

$$Z = \frac{\hat{p}_{obs} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} \stackrel{d}{\sim} N(0, 1)$$

Testa mot normalkvantilen

$Z_{1-\alpha}$ (ensidig hypotes)

$Z_{1-\alpha/2}$ (dubbelsidig hypotes).

Interens om proportioner

* Om vi istället har två proportioner
vi vill jämföra, säg

$$X \in \text{Bin}(n_x, p_x) \text{ och } Y \in \text{Bin}(n_y, p_y)$$

vi vill testa följande hypotes:

$$H_0: p_x = p_y$$

$$H_1: p_x \neq p_y$$

Anta att $n_x p_x (1-p_x) > 10$ och $n_y p_y (1-p_y) > 10$.

Interens om proportioner

* Om vi istället har två proportioner
vi vill jämföra, säg

$$X \in \text{Bin}(n_x, p_x) \text{ och } Y \in \text{Bin}(n_y, p_y)$$

vi vill testa följande hypotes:

$$H_0: p_x = p_y$$

$$H_1: p_x \neq p_y$$

Anta att $n_x p_x (1-p_x) > 10$ och $n_y p_y (1-p_y) > 10$.

Under H_0 så gäller att $p_0 = p_x = p_y$.

Så vi kan skatta p_0 gemensamt:

$$\hat{p}_0 = \frac{x + y}{n_x + n_y} \quad \text{och testa}$$

$$Z = \frac{\hat{p}_x - \hat{p}_0}{\sqrt{\hat{p}_0 (1 - \hat{p}_0) \left(\frac{1}{n_x} + \frac{1}{n_y} \right)}} \quad \text{mot normalförd.}$$

Finns i formelsamlingen!

INFERENS OM PROPORTIONER:

EXEMPEL: I Stockholms län gjorde man 1990 en undersökning av förekomsten av pollenallergi bland vissa känsliga grupper. Man valde slumpmässigt ut 500 personer i åldern 20–64 år och av dessa hade 23 % pollenallergi.

- Vad kan vi säga om andelen pollenallergiker i populationen?
- År 1994 gjordes motsvarande undersökning och 500 nya personer valdes ut. Då hade 29% pollenallergi. Kan man rimligen säga att det skett en förändring av benägenheten för denna typ av allergi under perioden?

ANALYS AV KATEGORIDATA

EXEMPEL: Varje individ i en viss population hör i genetiskt hänseende till en av fyra kategorier K_1, K_2, K_3, K_4 . Teoretiskt skall de fyra kategoriernas storlekar förhålla sig som 9 : 3 : 3 : 1. Vid en undersökning av 160 slumpmässigt utvalda ur populationen fick man följande resultat:

kategori	K_1	K_2	K_3	K_4
frekvens	78	42	27	13

Talar de observerade data emot teorin?

EXEMPEL: Finns det ett samband mellan blodgrupp och risken för magsår? Blodgruppen bestämdes för 1655 magsårspatienter och för en kontrollgrupp om 10000 personer från samma stad. Resultat:

	0	A	B	AB	Totalt
Magsårspatienter	911	579	124	41	1655
Kontrollgrupp	4578	4219	890	313	10000
Totalt	5489	4798	1014	354	11655

Interens om proportioner

Lösning: $X =$ "Antal pollenallergiker" $\in \text{Bin}(500, p_1)$

$p_1 =$ andelen pollenallergiker

$$\hat{p}_1 = 23\%$$

Kan vi normalapproximera?

$$n \cdot p_1 (1 - p_1) = 500 \cdot 0.23 \cdot 0.77 = 88.55 > 10$$

Ja!

$$I_{p_1} = (\hat{p}_1 \pm z_{1-\alpha/2} \sqrt{\text{Var}(\hat{p}_1)})$$

(Låt $\alpha = 0.05$)

$$\frac{\hat{p}_1 (1 - \hat{p}_1)}{n}$$

$$= (0.23 \pm 1.96 \sqrt{\frac{0.23 \cdot (1 - 0.23)}{500}})$$

$$= (0.19, 0.27)$$

Interens om proportioner

Lösning: \mathbb{Y} = "Antal pollerallergiker 1994"

$$\mathbb{Y} \in \text{Bin}(500, p_2)$$

$$1990 \quad \hat{p}_1 = 0.23 \quad n_1 = 500$$

$$1994 \quad \hat{p}_2 = 0.29 \quad n_2 = 500$$

Vi vill testa

$$H_0: p_1 = p_2$$

$$H_1: p_1 \neq p_2$$

Under H_0 är $p_0 = p_1 = p_2$.

$$\hat{p}_0 = \frac{x + y}{n_1 + n_2} = \frac{\hat{p}_1 n_1 + \hat{p}_2 n_2}{n_1 + n_2} = 0.26.$$

$$Z = \left| \frac{\hat{p}_1 - \hat{p}_2}{\sqrt{\hat{p}_0 (1 - \hat{p}_0) \left(\frac{1}{n_1} + \frac{1}{n_2}\right)}} \right| = 2.1628 > 1.96 \quad = z_{0.975}$$

H_0 kan förkastas på nivå $\alpha = 0.05$.

Interens om proportioner

* Om stickprovet är litet så att normalapproximation ej funkar:

Använd direktmetoden!

Exempel: $X \in \text{Bin}(8, p)$, $x=5$

$$H_0 : p = 1/2$$

$$H_1 : p > 1/2$$

$$\begin{aligned} p\text{-värde} &= P(\text{"Få det vi fick eller värre"} | H_0 \text{ sant}) \\ &= P(X \geq 5 | p = 1/2) = 1 - P(X \leq 4 | p = 1/2) \\ &= 1 - 0.6875 = 0.3125. \end{aligned}$$

H_0 kan ej förkastas på någon rimlig nivå.

Analys av kategoridata

Ett försök kan utfalla på "k" olika sätt.

Gör "n" oberoende försök, räkna hur många försök som hamnar i varje kategori.

Leder oftast till χ^2 -test.

Vi ska gå igenom 3 situationer:

- 1) Test av modellanpassning
- 2) Homogenitetstest.
- 3) Oberoende-test.

Analys av kategoridata

Test av modellanpassning: Man har

k - kategorier

O_i - antalet observationer i kategori i

samt en hypotes kring modellen: tex.

$$H_0: p_1 = \frac{9}{16}, p_2 = \frac{5}{16}, p_3 = \frac{3}{16}, p_4 = \frac{1}{16}$$

H_1 : några av dessa sannolikheter är fel.

Man skapar sen vad man borde ha

sett om H_0 var sann:

$$E_i = n \cdot p_i \quad \text{for varje kategori.}$$

	Kategori					
	1	2	3	4	...	k
O_i	120	48	32	107	...	99
E_i	np_1	np_2	...			np_k

Test:

$$\chi^2 = \sum \frac{(O_i - E_i)^2}{E_i}$$

$$\text{mot } \chi^2_{(1-\alpha, k-1)}$$

Frhetsgrader =
antal kategorier - 1.

Analys av kategoridata

Homogenitetsbtest

* Vi har k kategorier och r stickprov där vi har bestämt n_r (dvs antalet mätningar i stickprov r) innan.

Stickprov	Kategori				Summa
	1	2	...	k	
1	O_{11}	O_{12}	...	O_{1k}	$n_{1.}$
2	O_{21}	O_{22}	...	O_{2k}	$n_{2.}$
...
r	O_{r1}	O_{r2}	...	O_{rk}	$n_{r.}$
Summa	$n_{.1}$	$n_{.2}$...	$n_{.k}$	n

Vi testar:

$$H_0: p_{11} = p_{21} = p_{31} = \dots = p_{rk}$$

H_1 : Några p skiljer sig.

Analys av kategoridata

$$H_0: p_{11} = p_{21} = p_{31} = \dots = p_{rk}$$

H_1 : Några p skiljer sig.

Vi skapar en ny tabell som det borde vara om H_0 var sann.

$$E_{ij} = \frac{n_{i.} \cdot n_{.j}}{n}$$

Strickprov	Kategori				Summa
	1	2	...	k	
1	E_{11}	E_{12}	...	E_{1k}	$n_{1.}$
2	E_{21}	E_{22}	...	E_{2k}	$n_{2.}$
...
r	E_{r1}	E_{r2}	...	E_{rk}	$n_{r.}$
Summa	$n_{.1}$	$n_{.2}$...	$n_{.k}$	n

Vi testar:
$$\chi^2 = \sum_i \sum_j \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

och jämför med $\chi^2_{1-\alpha, (r-1)(k-1)}$

Med $(r-1)(k-1)$ frihetsgrader.

Analys av kategoridata

Oberoendetest: (Typ samma som homogenitetstest)

Nu testar vi

$$H_0: P_{ij} = P_i \cdot P_j \text{ för alla } i \text{ och } j.$$

Strickprov	Kategori				Summa
	1	2	...	k	
1	O_{11}	O_{12}	...	O_{1k}	$n_{1.}$
2	O_{21}	O_{22}	...	O_{2k}	$n_{2.}$
...
r	O_{r1}	O_{r2}	...	O_{rk}	$n_{r.}$
Summa	$n_{.1}$	$n_{.2}$...	$n_{.k}$	n

Strickprov	Kategori				Summa
	1	2	...	k	
1	P_{11}	P_{12}	...	P_{1k}	$P_{1.}$
2	P_{21}	P_{22}	...	P_{2k}	$P_{2.}$
...
r	P_{r1}	P_{r2}	...	P_{rk}	$P_{r.}$
Summa	$P_{.1}$	$P_{.2}$...	$P_{.k}$	1

Analys av kategoridata

Om H_0 är sann så borde vi få:

$$E_{ij} = \frac{n_{i.} \cdot n_{.j}}{n}$$

Strckprov	Kategori				Summa
	1	2	...	k	
1	E_{11}	E_{12}	...	E_{1k}	$n_{1.}$
2	E_{21}	E_{22}	...	E_{2k}	$n_{2.}$
⋮	⋮	⋮	⋮	⋮	⋮
r	E_{r1}	E_{r2}	...	E_{rk}	$n_{r.}$
Summa	$n_{.1}$	$n_{.2}$...	$n_{.k}$	n

Vi testar (igen)

$$\chi^2 = \sum_j \sum_i \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

och jämför med $\chi^2_{1-\alpha, (r-1)(k-1)}$

INFERENS OM PROPORTIONER:

EXEMPEL: I Stockholms län gjorde man 1990 en undersökning av förekomsten av pollenallergi bland vissa känsliga grupper. Man valde slumpmässigt ut 500 personer i åldern 20–64 år och av dessa hade 23 % pollenallergi.

- Vad kan vi säga om andelen pollenallergiker i populationen?
- År 1994 gjordes motsvarande undersökning och 500 nya personer valdes ut. Då hade 29% pollenallergi. Kan man rimligen säga att det skett en förändring av benägenheten för denna typ av allergi under perioden?

ANALYS AV KATEGORIDATA

EXEMPEL: Varje individ i en viss population hör i genetiskt hänseende till en av fyra kategorier K_1, K_2, K_3, K_4 . Teoretiskt skall de fyra kategoriernas storlekar förhålla sig som 9 : 3 : 3 : 1. Vid en undersökning av 160 slumpmässigt utvalda ur populationen fick man följande resultat:

kategori	K_1	K_2	K_3	K_4
frekvens	78	42	27	13

Talar de observerade data emot teorin?

EXEMPEL: Finns det ett samband mellan blodgrupp och risken för magsår? Blodgruppen bestämdes för 1655 magsårspatienter och för en kontrollgrupp om 10000 personer från samma stad. Resultat:

	0	A	B	AB	Totalt
Magsårspatienter	911	579	124	41	1655
Kontrollgrupp	4578	4219	890	313	10000
Totalt	5489	4798	1014	354	11655

Test av modellanpassning.

$$H_0: P_1 = \frac{9}{16}, P_2 = \frac{3}{16}, P_3 = \frac{3}{16}, P_4 = \frac{1}{16}$$

H_1 : Inte som i H_0 .

$$n = 160, \alpha = 0.05, K = 4 \text{ (Antal kategorier)}$$

Test av modellanpassning.

$$H_0: P_1 = \frac{9}{16}, P_2 = \frac{3}{16}, P_3 = \frac{3}{16}, P_4 = \frac{1}{16}$$

H_1 : Inte som H_0 .

$$n = 160, \alpha = 0.05, K = 4 \text{ (Antal kategorier)}$$

O-tabell:

Kategori	K_1	K_2	K_3	K_4
Frekvens	78	42	27	13

E-tabell: (Enligt H_0)

Kategori	K_1	K_2	K_3	K_4
Frekvens	90	30	30	10

Test av modellanpassning.

$$H_0: P_1 = \frac{9}{16}, P_2 = \frac{3}{16}, P_3 = \frac{3}{16}, P_4 = \frac{1}{16}$$

H_1 : Inte som H_0 .

$n = 160$, $\alpha = 0.05$, $K = 4$ (Antal kategorier)

O-tabell:

E-tabell: (Enligt H_0)

Kategori	K_1	K_2	K_3	K_4	Kategori	K_1	K_2	K_3	K_4
Frekvens	78	42	27	13	Frekvens	90	30	30	10

$$\begin{aligned} \text{Test: } \chi^2 &= \sum_{i=1}^4 \frac{(O_i - E_i)^2}{E_i} = \frac{(78-90)^2}{90} + \frac{(42-30)^2}{30} + \\ &+ \frac{(27-30)^2}{30} + \frac{(13-10)^2}{10} = 7.6 \end{aligned}$$

Test av modellanpassning.

$$H_0: P_1 = \frac{9}{16}, P_2 = \frac{3}{16}, P_3 = \frac{3}{16}, P_4 = \frac{1}{16}$$

H_1 : Inte som H_0 .

$n = 160$, $\alpha = 0.05$, $k = 4$ (Antal kategorier)

O-tabell:

E-tabell: (Enligt H_0)

Kategori	K_1	K_2	K_3	K_4	Kategori	K_1	K_2	K_3	K_4
Frekvens	78	42	27	13	Frekvens	90	30	30	10

$$\text{Test: } \chi^2 = \sum_{i=1}^4 \frac{(O_i - E_i)^2}{E_i} = \frac{(78-90)^2}{90} + \frac{(42-30)^2}{30} + \frac{(27-30)^2}{30} + \frac{(13-10)^2}{10} = 7.6$$

$$\chi_{1-\alpha, k-1}^2 = \chi_{0.95, 3}^2 \stackrel{\text{tabell/dator}}{=} 7.815$$

Test av modellanpassning.

$$H_0: P_1 = \frac{9}{16}, P_2 = \frac{3}{16}, P_3 = \frac{3}{16}, P_4 = \frac{1}{16}$$

H_1 : Inte som H_0 .

$n = 160$, $\alpha = 0.05$, $k = 4$ (Antal kategorier)

O-tabell:

E-tabell: (Enligt H_0)

Kategori	K_1	K_2	K_3	K_4	Kategori	K_1	K_2	K_3	K_4
Frekvens	78	42	27	13	Frekvens	90	30	30	10

$$\text{Test: } \chi^2 = \sum_{i=1}^4 \frac{(O_i - E_i)^2}{E_i} = \frac{(78-90)^2}{90} + \frac{(42-30)^2}{30} + \frac{(27-30)^2}{30} + \frac{(13-10)^2}{10} = 7.6$$

$$\chi_{1-\alpha, k-1}^2 = \chi_{0.95, 3}^2 \stackrel{\text{tabell/dator}}{=} 7.815$$

Eftersom $\chi^2 = 7.6 < 7.815 = \chi_{0.95, 3}^2$

så kan vi inte förkasta H_0 på nivån

$$\alpha = 0.05.$$

INFERENS OM PROPORTIONER:

EXEMPEL: I Stockholms län gjorde man 1990 en undersökning av förekomsten av pollenallergi bland vissa känsliga grupper. Man valde slumpmässigt ut 500 personer i åldern 20–64 år och av dessa hade 23 % pollenallergi.

- Vad kan vi säga om andelen pollenallergiker i populationen?
- År 1994 gjordes motsvarande undersökning och 500 nya personer valdes ut. Då hade 29% pollenallergi. Kan man rimligen säga att det skett en förändring av benägenheten för denna typ av allergi under perioden?

ANALYS AV KATEGORIDATA

EXEMPEL: Varje individ i en viss population hör i genetiskt hänseende till en av fyra kategorier K_1, K_2, K_3, K_4 . Teoretiskt skall de fyra kategoriernas storlekar förhålla sig som 9 : 3 : 3 : 1. Vid en undersökning av 160 slumpmässigt utvalda ur populationen fick man följande resultat:

kategori	K_1	K_2	K_3	K_4
frekvens	78	42	27	13

Talar de observerade data emot teorin?

EXEMPEL: Finns det ett samband mellan blodgrupp och risken för magsår? Blodgruppen bestämdes för 1655 magsårspatienter och för en kontrollgrupp om 10000 personer från samma stad. Resultat:

	0	A	B	AB	Totalt
Magsårspatienter	911	579	124	41	1655
Kontrollgrupp	4578	4219	890	313	10000
Totalt	5489	4798	1014	354	11655

Homogenitetstest

$$H_0: p_0 = p_A = p_B = p_{AB}$$

H_1 : H_0 stämmer ej.

r - rader

c - kolumner

O-tabell

	O	A	B	AB	Totalt
Magsär	911	579	124	41	1655
Kontrollgrupp	4578	4219	890	313	10000
Totalt	5489	4798	1014	354	11655

Homogenitetstest

$$H_0: p_0 = p_A = p_B = p_{AB}$$

$H_1: H_0$ stämmer ej.

r - rader

c - kolumner

O-tabell

	O	A	B	AB	Totalt
Magsär	911	579	124	41	1655
Kontrollgrupp	4578	4219	890	313	10000
Totalt	5489	4798	1014	354	11655

E-tabell

	O	A	B	AB	Totalt
Magsär	779.4	681.3	144.0	50.3	1655
Kontrollgrupp	4709.6	4116.7	870.0	303.3	10000
Totalt	5489	4798	1014	354	11655

Homogenitetstest

$$H_0: p_0 = p_A = p_B = p_{AB}$$

H_1 : H_0 stämmer ej.

r - rader

c - kolumner

O-tabell

	O	A	B	AB	Totalt
Magsär	911	579	124	41	1655
Kontrollgrupp	4578	4219	890	313	10000
Totalt	5489	4798	1014	354	11655

E-tabell

	O	A	B	AB	Totalt
Magsär	779.4	681.3	144.6	50.3	1655
Kontrollgrupp	4709.6	4116.7	870.0	303.3	10000
Totalt	5489	4798	1014	354	11655

$$\chi^2 = \sum_{c=1}^4 \sum_{r=1}^2 \frac{(O_{cr} - E_{cr})^2}{E_{cr}} = \frac{(911 - 779.4)^2}{779.4} + \dots$$
$$= 49.0155.$$

Homogenitetstest

$$H_0: p_0 = p_A = p_B = p_{AB}$$

H_1 : H_0 stämmer ej.

r - rader

c - kolumner

O-tabell

	O	A	B	AB	Totalt
Magsär	911	579	124	41	1655
Kontrollgrupp	4578	4219	890	313	10000
Totalt	5489	4798	1014	354	11655

E-tabell

	O	A	B	AB	Totalt
Magsär	779.4	681.3	144.0	50.3	1655
Kontrollgrupp	4709.6	4116.7	870.0	303.3	10000
Totalt	5489	4798	1014	354	11655

$$\chi^2 = \sum_{c=1}^4 \sum_{r=1}^2 \frac{(O_{cr} - E_{cr})^2}{E_{cr}} = \frac{(911 - 779.4)^2}{779.4} + \dots$$

$$= 49.0155. \quad \text{Jämfors med } \chi^2_{1-\alpha, (r-1)(c-1)}$$

$$\chi^2_{0.999, 3} = 16.266. \quad \text{Eftersom } 49.0155 > 16.266$$

så kan vi förkasta H_0 på nivån $\alpha = 0.001$

Provtagningstider – problemställningar från labbet

Vid vårt kemilaboratorium analyserar vi bland annat en rad prover från den näraliggande provtagningscentralen. Det har framkommit starka önskemål om att vissa patienter som tar prover ska kunna träffa en läkare vid samma besök och inte behöva boka in olika dagar för provtagning och läkarbesök. För att detta önskemål ska uppfyllas måste vi förstås veta hur lång tid det tar från det att provet tagits på patienten till analyssvaret är klart och undersöka om det är rimligt att låta patienter vänta på provsvar samma dag. Det finns flera moment att ta hänsyn till: provet kan få vänta på provtagningscentralen tills det blir hämtat till vårt laboratorium, det behövs en viss manuell handläggningstid av provet och slutligen har vi själva processtiden i maskinen. Dessutom har vi lite olika hanteringstider beroende på vilken dag i veckan det är och om det är för- eller eftermiddag.

I datafilerna `proverfm.Rdata` och `proverem.Rdata` finns det tider (minuter) som det tog ”från patientarm till analysvar”. Provet är på så kallad ”allmän kemi” och, som ni märker, har vi har delat upp data i två filer, en för prover tagna på förmiddagen och en för prover tagna på eftermiddagen.

Nu till våra frågor: När vi gjorde ett histogram på data slogs vi av att det inte alls liknade en normalfördelning. Det kan väl i och för sig vara rimligt, men kan vi dra några slutsatser då? Finns det andra fördelningar som kan användas för att modellera tiden? Mer specifikt:

- Hur sannolikt är det att en förmiddagspatient får vänta mer än två timmar på analysvar?
- Vi vill kunna säga att ”95 % av förmiddagspatienterna kommer att ha sitt provsvar snabbare än x minuter”. Vad är då x ? Om vi vill göra samma uttalande för eftermiddagspatienterna, vad är x då?
- Vi beräknar att ca 40 % av proverna kommer på förmiddagen och resten på eftermiddagen. Hur troligt är det att ett patientprov, taget någon gång under dagen, tar mer än två timmar att analysera?
- Om vi har 50 patientprover på förmiddagen, vad är sannolikheten att **genomsnittstiden** för dessa 50 prov överstiger en timme?

Tips på arbetsgång

- Titta på data (histogram, empirisk fördelningsfunktion). beräkna enkla mått (medelvärde, standardavvikelse). Jämför förmiddags- och eftermiddagstider.
- Anpassa en lämplig standardfördelning till förmiddagstiderna och skatta parametrarna i den fördelningen.
- Gör samma sak för eftermiddagstiderna.
- Svara på frågorna (a)–(c) genom att utnyttja de anpassade fördelningarna.
- Fundera på vad medelvärdet av 50 förmiddagstider har för fördelning. Använd de ”enkla mått” du beräknade tidigare för att hitta rätt parametrar i fördelningen.