

At træffe sine valg i en usikker verden - eller den statistiske modellerings rolle.

Af E. Susanne Christensen. Lektor i statistik.

Institut for Matematiske Fag. Aalborg Universitet.

I mange tilfælde og i mange forskellige faglige sammenhænge må man træffe en afgørelse eller basere en overbevisning på et ikke fuldstændigt informationsgrundlag. Dette er fx tilfældet, når man ønsker at forudsige udfaldet af et kommende folketingsvalg og kun har en opinionsundersøgelse til rådighed, eller når man ønsker at vide, om antallet af rygere er stigende blandt unge mennesker, men ikke har mulighed for at spørge alle unge, om de ryger eller ej. Resultaterne vil i sådanne tilfælde bære præg af, hvilken **stikprøve**¹ man lægger til grund for til sin undersøgelse.

Eksempel:

Vi vil undersøge, om holdningen til skattelettelser i Danmark er den samme for gymnasieelever som for gymnasielærere. For at finde ud af det, må vi indsamle nogle data, som vi kan basere vores konklusioner på: Vi skal lave en stikprøve!

Hvis vi vidste, det var så heldigt, at alle gymnasieelever i landet var enige med hinanden, og at også alle gymnasielærere var indbyrdes enige i spørgsmålet om skattelettelser, så kunne vi hurtigt blive færdige. Vi skulle bare spørge én person fra hver gruppe, hvad de mente om sagen, og så fastslå, om de to grupper også var enige; Og vi ville således være helt sikre på, at vi havde draget den rigtige konklusion. Men så nemt går det sjældent.

Der kan være ret stor forskel på meningene inden for en gruppe personer. Og selvom det faktisk skulle forholde sig sådan, at de fleste gymnasieelever går ind for skattelettelser, så kan vi jo godt være uheldige – til vores stikprøve - at udvælge de elever, der er imod ... og lige så uheldige kunne vi være med vores udvælgelse blandt lærerne. Hvis det er tilfældet, så vil vi ende op med den forkerte konklusion om gruppernes mening om spørgsmålet.

Det, vi kan gøre for at mindske risikoen for at lave forkerte konklusioner, er at tage "fornuftige" stikprøver, der er "store nok" til, at risikoen for at komme frem til en forkert konklusion bliver "ac-

¹ Ord markeret med rød farve indikerer, at der findes en specifik matematisk definition af ordet. Når det ikke defineres her skyldes det, at den almindelige brug og opfattelse af ordet normalt får en til at agere i overensstemmelse med den korrekte matematiske definition. I dette tilfælde kan finde den korrekte definition i [2] (stikprøve=sample).

ceptabel". Statistik går (blandt andet) ud på at præcisere alle de begreber, der her står i anførselstegn. Hvad sandsynligheden er for at drage en forkert konklusion, kan beregnes. I hvert tilfælde hvis man følger nogle enkle regler, når man indsamler data. Hvilke udvælgelsesmåder man kan bruge i praksis, er beskrevet nærmere i [2] og [3].

I denne note skal vi se på, hvordan man kan sætte tal på ens usikkerhed i et par specifikke tilfælde.

Statistisk test for uafhængighed mellem to inddelingskriterier.

Hvis vi vil undersøge, om det at bruge "mange penge på tøj" er lige udbredt blandt unge kvinder og unge mænd, er den mest objektive måde at forholde sig på at lave en empirisk undersøgelse. Dvs. man indsamler data og drager sine slutninger på baggrund af dem. Allerførst er der et par ting, der skal præciseres!

Hvad er det for nogle unge, vi interesserer os for? I den statistiske terminologi hedder det at fastlægge *populationen*. Lad os sige, at det er unge mellem 15-20 år og bosiddende i Danmark, vi er interesseret i.

Så skal vi have præciseret, hvad vi egentligt forstår ved at "bruge mange penge på tøj". I den statistiske terminologi siger man, at vi skal have formuleret *modellen og hypoteserne*.

Modellen kan her være, at andelen af kvinder, der bruger mere end 1500 kr. om måneden på tøj er p_k og den tilsvarende for mændene er p_m . (Andelen p_k svarer til **sandsynligheden**² for at en kvinde **tilfældigt udvalgt**³ fra populationen bruger mere end 1500 kr. på tøj om måneden). Grænsen for, hvornår man "bruger mange penge på tøj", er jo her sat subjektivt og kan selvfølgelig gøres til genstand for diskussion ☺. Vi kommer tilbage til hvilke *matematiske krav*, der skal stilles til denne grænse.

Hypotesen er, at $p_m = p_k$, eller sagt i ord: Andelen af unge mænd, der bruger mange penge på tøj, er den samme som den tilsvarende andel for unge kvinder.

For at undersøge vores hypotese kan vi fx gennemføre følgende forsøg:

Vi udvælger et antal unge mellem 15 og 20 år tilfældigt og spørger dem om, hvor mange penge, de bruger på tøj om måneden. Så tæller vi op, hvor mange kvinder, der bruger mere end 1500 kr., og hvor mange mænd.

(En anden statistisk undersøgelse kunne foretages med udgangspunkt i svarene om størrelsen af de brugte beløb og fx undersøge, om det gennemsnitlige forbrug for kvinder er større end det gennemsnitlige forbrug for mænd. I så fald skulle man bruge den statistiske metode, der hedder sammenligning af to middelværdier. Det kan man læse mere om i fx [2])

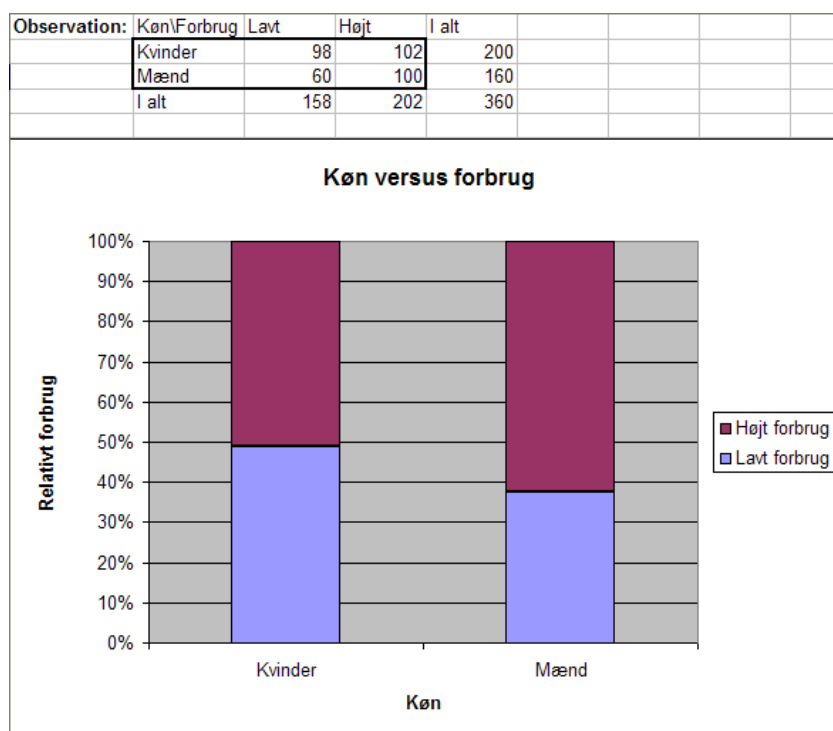
² For en introduktion til basal sandsynlighedsregning se fx [1] eller [2].

³ Tilfældigt udvalgt betyder, at alle individer i populationen har samme sandsynlighed for at blive udvalgt. Vore beregninger forudsætter, at stikprøven er udvalgt på denne måde.

Vores resultat af undersøgelsen kan vi organisere i en tabel som nedenfor:
 (Tallene er rent gætværk - ikke resultat af en virkelig undersøgelse. Sådan én kan I jo lave!).

	<1500 kr./måned	≥ 1500 kr./måned	I alt
Kvinder	98	102	200
Mænd	60	100	160
I alt	158	202	360

Vi kan eventuelt fremstille tallene fra stikprøven grafisk som nedenfor:



Vi kan nu gætte på andelen af unge kvinder, der bruger mange penge på tøj, simpelthen ved at regne ud, hvor stor andelen er i stikprøven. Vi siger, vi *estimerer* p_k . I dette tilfælde får vi *et estimat* på $102/200 = 0.51$, hvilket vil sige, at vi tror, at 51 % af de unge kvinder bruger mange penge på tøj. Tilsvarende estimerer vi andelen af mænd til $100/160 = 0.625$. Vi tror altså, at der var 62.5 % af de unge mænd, der bruger mange penge på tøj.

I den **stikprøve** vi har, er andelen af kvinder, der bruger mange penge på tøj, altså mindre end den tilsvarende andel for mænd.

Umiddelbart kan det altså se ud som om, at vores hypotese om samme andel for de to køn af personer, der bruger mange penge på tøj, IKKE holder stik. Resten af øvelsen går ud på at afgøre, om dette blot er en tilfældig følge af en uheldig stikprøve, ELLER om det vi har set, er så markant, at vi

tør tage det som udtryk for en forskel mellem de to køn generelt, altså noget vi tror, der gælder for hele populationen.

For at kunne regne på sagen og lave et statistisk test er det matematisk set vigtigt, at de enkelte svar på, hvor mange penge man bruger, er uafhængige⁴ af hinanden. Hvis man fx i sin stikprøve har valgt en gruppe venner med stor indbyrdes påvirkning, så vil denne gruppe sagtens kunne have en adfærd, som er atypisk for populationen som helhed, og derved påvirke undersøgelsens resultat uhensigtsmæssigt.

Statistisk hypotesetest minder logisk set om det, du måske kender fra din matematikundervisning som et modstridsargument. Man antager en ting, gennemfører en række logiske argumenter og ender op med en konklusion, der klart er forkert. Heraf slutter man, at den oprindelige antagelse IKKE kan være rigtig. I statistik tager man hensyn til, at verden ikke er deterministisk, så hér kan man ikke konkludere, at udgangsantagelsen ikke er sand, men man kan eventuelt slutte, at det, man har set i sit forsøg vil være USANDSYNLIGT, hvis udgangsantagelsen er sand. Dermed tyder forsøget på, at antagelsen ikke er rigtig.

Vores antagelse om, at forbrugsmønstret er ens for de to køn formuleres som vores udgangshypotese. Hvis vores test kan afvise den hypotese, så har vi et vist belæg for at påstå, at der er en **signifikant forskel**⁵ mellem de to køn. Vi har således en udgangshypotese, som per tradition kaldes H_0 og en alternativ hypotese H_1 givet som:

$$H_0 \quad p_k = p_m$$

$$H_1 \quad p_k \neq p_m$$

En anden måde at udtrykke H_0 på er, at der er uafhængighed mellem det at bruge mange penge på tøj og ens køn. H_1 svarer så til, at der er afhængighed mellem de to inddelingskriterier - forbrug og køn, dvs.

H_0 *Der er uafhængighed mellem de to kriterier.*

H_1 *Der er ikke uafhængighed mellem de to kriterier.*

Vi starter med at antage, at H_0 udtrykker den sande tilstand af verden. I så fald kan vi estimere andelen af unge, der bruger mange penge på tøj, uden hensyntagen til, om de tilhører det ene eller det andet køn. Andelen af "storforbrugere" estimeres så til $202/360=0.5611$, altså 56.11 %. Så hvis vi har en gruppe på 200 unge, vil vi forvente, at $200*0.5611=112.22$ af dem er storforbrugere, og $200*(1-0.5611)=89.78$ af dem var ikke-storforbrugere, uanset hvilket køn de har.

Her har vi brugt regelen, at hvis sandsynligheden for at være storforbruger er givet ved p , så er sandsynligheden for det modsatte, nemlig at være ikke-storforbruger, givet ved $(1-p)$. Ud over at være logisk er dette også en regneregul fra den basale sandsynlighedsteori.

⁴ At to hændelser A og B er uafhængige betyder at $P(A \cap B) = P(A) * P(B)$. For regneregler for sandsynligheder se [2] eller [3].

⁵ Begrebet statistisk signifikans er relateret til statistisk testteori. Se fx [2].

Ved at regne på den måde kan vi udfylde skemaet med de værdier, som vi ville have forventet at se, hvis verden opførte sig som vores H_0 foreskriver.

Forventede værdier under antagelse af at der er uafhængighed:

	<1500 kr./måned	≥1500 kr./måned	Ialt
Kvinder	$\frac{158}{360} * 200 = 87.78$	$\frac{202}{360} * 200 = 112.22$	200
Mænd	$\frac{158}{360} * 160 = 70.22$	$\frac{202}{360} * 160 = 89.78$	160
Ialt	158	202	360

Afviselserne mellem det resultat, vi fik i forsøget, og de hér udregnede forventede værdier er et udtryk for, hvor langt forsøgets virkelighed er fra den verden, der er modelleret i H_0 .

Imidlertid er det sådan, at summen af afvigelserne $(98-87.77)+(102-112.23)+(60-70.23)+(100-89.77) = 0$, og sådan vil det altid være. Så at lægge afvigelserne sammen giver os ikke noget samlet billede af, hvor stor afvigelsen er. I stedet viser det sig at være smart at udregne en χ^2 **teststørrelse**. Man udregner differens mellem det observerede antal og det forventede antal i hver celle, sætter denne differens i anden og dividerer med det forventede antal. Til sidst summeres disse tal for alle celler, altså:

$$\chi^2 = \sum \frac{(\text{obs. antal} - \text{forv. antal})^2}{\text{forv. antal}}$$

En stor værdi af teststørrelsen tyder i denne sammenhæng på, at udgangshypotesen om uafhængighed IKKE er opfyldt. Altså: store værdier af χ^2 teststørrelsen får os til at tro mere på H_1 . Vi har så bare det problem tilbage, at vi skal afgøre, HVOR stor en teststørrelse skal være, før vi mener, den er så stor, at vi ikke vil tro på H_0 . Til det brug skal vi vide, hvor store værdier teststørrelsen normalt vil antage, når H_0 er sand.

Hvis hypotesen om uafhængighed er rigtig, og hvis man har en stor nok stikprøve (sådan at alle de forventede værdier er større end 5), så ved man – takket være nogle matematikeres arbejde – hvilke værdier denne teststørrelse ville antage, hvis man lavede en uendelig række af forsøg som det skitserede. Den statistiske terminologi er, at man kender teststørrelsens **fordeling**⁶ under H_0 , idet den nemlig vil følge det, der hedder en χ^2 fordeling med 1 frihedsgrad. (Udtales "ki i anden"-fordelingen.)

(Og det er her, vi for en kort stund kan vende tilbage til vor subjektivt fastsatte grænse for, hvornår man bruger mange penge på tøj. Havde vi sat den grænse så højt eller så lavt, at der i en af cellerne med de forventede værdier var kommet et tal mindre end 5, så skulle vi enten have lavet grænsen om, eller være gået over til en anden statistik metode.⁷)

⁶ Begrebet fordeling kræver introduktion af begrebet stokastisk variabel for at formaliseres. Se [2] eller [3].

⁷ Fx Fishers eksakte test.

Matematisk kan man vise, at i en verden, hvor køn og forbrugsmønster er uafhængige størrelser, så vil man i 5% af de gange, hvor man udvælger en stikprøve på 360 personer, få en teststørrelse, der er større end 3.84. I 1% af ville man få en teststørrelse, der er større end 6.63.⁸

Disse tal –også kaldet kritiske værdier - kan findes i tabeller, på moderne lommeregnerne, i Excel og i statistiske værktøjsprogrammer. Med Excel ser det fx således ud:

CHIFORDELING		=chiinv(0.05;1)			
	A	B	C	D	E
1	signifikansniveau = 1%	6.634897			
2	signifikansniveau = 5%	3.841459			
3		=chiinv(0.05;1)			
4		CHIINV(sandsynlighed; frihedsgrader)			
5					

Teststørrelsen fra vores stikprøve bliver $\chi^2 = (98-87.78)^2/87.78+(102-112.22)^2/112.22+(60-70.22)^2/70.22+(100-89.78)^2/89.78= 1.19+0.93+1.49+1.16= 4.77$. Den er jo altså er større end 3.84. Så HVIS antagelsen om uafhængighed mellem køn og forbrug skal holde stik, så har vi hér set et forsøg, der vil optræde med en sandsynlighed, der er betydeligt mindre end 0.05. Den statistiske terminologi er, *at testsandsynligheden er mindre end 5 %*. Det er vist lettere at tro på, at antagelsen IKKE holder.

Vi forkaster vores udgangshypotese og siger: "Forsøget har påvist en sammenhæng mellem køn og forbrug på tøj, der er signifikant på 5% niveau."

Men vores teststørrelse er IKKE større end de 6.63. Det betyder at *testsandsynligheden er større end 1 %*. Vi kan derfor ikke påvise en signifikant sammenhæng på 1% niveau.

Hvis man, som det ofte er tilfældet, har en fast grænse for, hvornår man vil vælge at forkaste sin udgangshypotese, fx når testsandsynligheden er mindre end 5 %, siger man, at man arbejder med et *signifikansniveau* på 5%. Ved at bruge et fast signifikansniveau på fx 5% i en hel række af forsøg og test ved man altså, at man i 5% af testene fejlagtigt vil forkaste en sand udgangshypotese. Vi ved, hvor sikker metoden er, men vi ved ikke, om den enkelte beslutning om at tro på udgangshypotesen er rigtig eller ej.

På din lommeregner kan du få den præcise sandsynlighed for at få en χ^2 -teststørrelse, der er større end de 4.78, selvom H_0 er sand. (Ved at slå værdien 4.78 op i en χ^2 fordeling med 1 frihedsgrad.) Denne sandsynlighed kaldes *p-værdien* eller *testsandsynligheden* for testet.

En *p-værdi* er altså sandsynligheden for at få en teststørrelse, der får os til at tvivle mindst lige så meget på H_0 , som den, vi lige har set, selvom H_0 faktisk er den rigtig hypotese. I det aktuelle eksempel får vi en p-værdi på $p = 0.029$.

Også p-værdien kan vi finde ved hjælp af Excel eller lignende hjælpemidler.

⁸ De konkrete tal er fraktiler fra χ^2 -fordelingen med 1 frihedsgrad, der er den approximative fordeling af teststørrelsen. For en intuitiv forklaring af frihedsgrader, se [4].

Med Excel ser det således ud:

CHIFORDELING		=chifordeling(4.773531;1)		
	A	B	C	D
1	Teststørrelse =	4.773531		
2	Testsandsynlighed =	0.028901		
3		=chifordeling(4.773531;1)		

Hvad gør vi, hvis vi har flere niveauer på hvert af inddelingskriterierne?

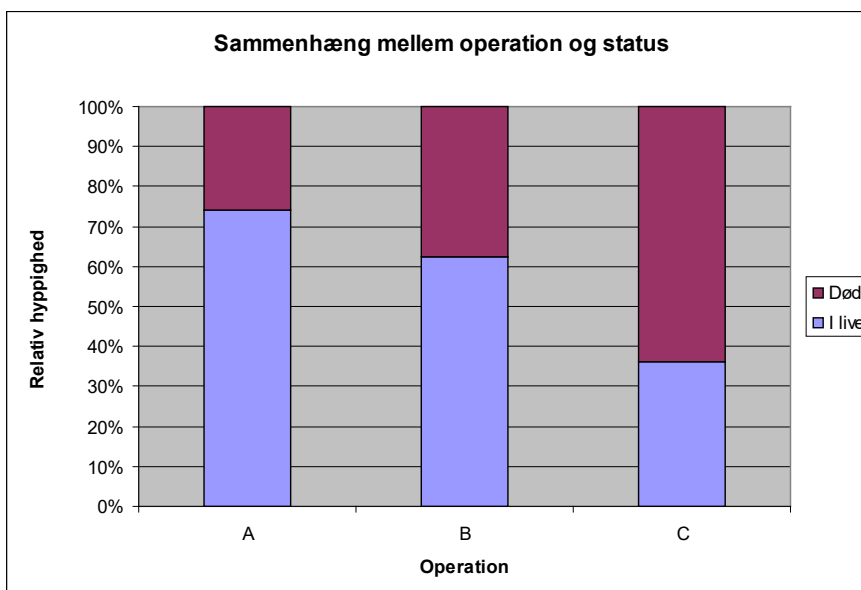
Et andet eksempel: En hjerneforsker undersøger forskellige måder at operere for en bestemt form for hjernetumor. Der er tre forskellige operationstyper: Type A, hvor man kun fjerner selve tumoren; Type B, hvor man også tager lidt af det nærmeste omkringliggende væv bort, og Type C, hvor såvel tumor som en større del af det omkringliggende væv fjernes.

Lægen ønsker at vide, hvordan operationstypen indvirker på chancen for at overleve et halvt år efter operationen. Over en årrække har lægen indsamlet følgende data over resultaterne af operationerne:

(data er fiktive, problemstillingen er autentisk)

	I live	Død	Ialt
Operation A	40	14	54
Operation B	10	6	16
Operation C	9	16	25
Ialt	59	36	95

En grafisk fremstilling af tallene kunne se sådan ud:



I **stikprøven** er der altså forskel på andelen af overlevende efter de forskellige typer af operation. Vi skal forsøge at undersøge, om det er en så stor forskel, at man kan sige, resultaterne kan generaliseres ud over denne stikprøve, altså: er det statistisk signifikant? Vi skal altså forsøge at regne ud, om den sete stikprøve er meget ekstrem, hvis vi antager, at der ikke er forskel på overlevelseschancerne efter de forskellige operationer.

Vi lader p_A være sandsynligheden for at være i live et halv år efter at have gennemgået en operation af type A, og tilsvarende for p_B og p_C .

Udgangshypotesen er, at der er samme sandsynlighed for overlevelse efter alle tre typer operation, dvs

$$H_0 \quad p_A = p_B = p_C$$

H_1 *Ikke alle tre sandsynligheder ens.*

Hvis udgangshypotesen holder, estimerer vi sandsynligheden for overlevelse ved $\frac{59}{95} = 0.6211$, og vi kan udregne de forventede værdier efter samme princip som før:

Forventede værdier	I live	Død	Ialt
Operation A	$0.6211 * 54 = 33.54$	$(1 - 0.6211) * 54 = 20.46$	54
Operation B	$0.6211 * 16 = 9.94$	$(1 - 0.6211) * 16 = 6.06$	16
Operation C	$0.6211 * 25 = 15.52$	$(1 - 0.6211) * 25 = 9.48$	25
Ialt	59	36	95

Alle de forventede størrelser er større end 5, så vi kan bruge χ^2 testet igen, nu bare med nogle lidt andre tal.⁹

Teststørrelsen udregnes som før ved at summe $\frac{(obs.antal - forv.antal)^2}{forv.antal}$ over alle celler.

Det vil sige, at vi her får

$$\frac{(40-33.54)^2}{33.54} + \frac{(14-20.46)^2}{20.46} + \frac{(10-9.94)^2}{9.94} + \frac{(6-6.06)^2}{6.06} + \frac{(9-15.52)^2}{15.52} + \frac{(16-9.48)^2}{9.48} = 10.5.$$

Denne gang skal vi bruge en χ^2 fordeling med $(2-1)*(3-1)=(antal_rækker - 1)*(antal_søjler - 1)=2$ frihedsgrader. Matematikeren, der er i besiddelse af de relevante tabeller, kan her fortælle os, at vi med en sandsynlighed på 5% vil få en teststørrelse større end 5.99, og med sandsynlighed 1% en teststørrelse større end 9.81, NÅR udgangshypotesen er sand.

Vi kan også selv finde en p-værdi via tabel, lommeregner eller i Excel med kommandoen CHIFORDELING(10.5;2), og vi får en testsandsynlighed på $p = 0.0051$.

⁹ Den approximerende χ^2 fordeling har denne gang 2 frihedsgrader. Antallet af frihedsgrader er $(antal_rækker - 1)*(antal_søjler - 1)$.

Så i dette tilfælde er vores valgmuligheder:

Vi fastholder troen på, at de tre operationstyper giver samme overlevelseschance, og vi har set et forsøg, der har mindre end 1 % sandsynlighed for at indtræffe.

ELLER

Vi forkaster hypotesen om ens chancer for overlevelse og siger:

Der er fundet en sammenhæng mellem overlevelseschance og operationstype, der er statistisk signifikant på 1% niveau.

Imidlertid skal man tænke sig om, inden man foreslår operationstype C forbudt. Hvis der er en tredje faktor der influerer billedet, så kan det give misvisende konklusioner, når man kun tager to af dem i betragtning. Vi kigger lidt nærmere på tallene fra før:

	I live	Død	Ialt
Operation A	40	14	54
Operation B	10	6	16
Operation C	9	16	25
Ialt	59	36	95

Eller i procenter:

	I live	Død	Ialt
Operation A	74%	26%	100%
Operation B	63%	37%	100%
Operation C	36%	64%	100%

Efter en nærmere inspektion af journalerne viser det sig, at også patientens alder er noteret, og ved at opdele i to grupper efter alder får vi billedet:

50 år eller derunder:

	I live	Død	Ialt
Operation A	27	1	28
Operation B	2	0	2
Operation C	1	0	1
Ialt	30	1	31

eller i procent

	I live	Død	Ialt
Operation A	96%	4%	100%
Operation B	100%	0%	100%
Operation C	100%	0%	100%

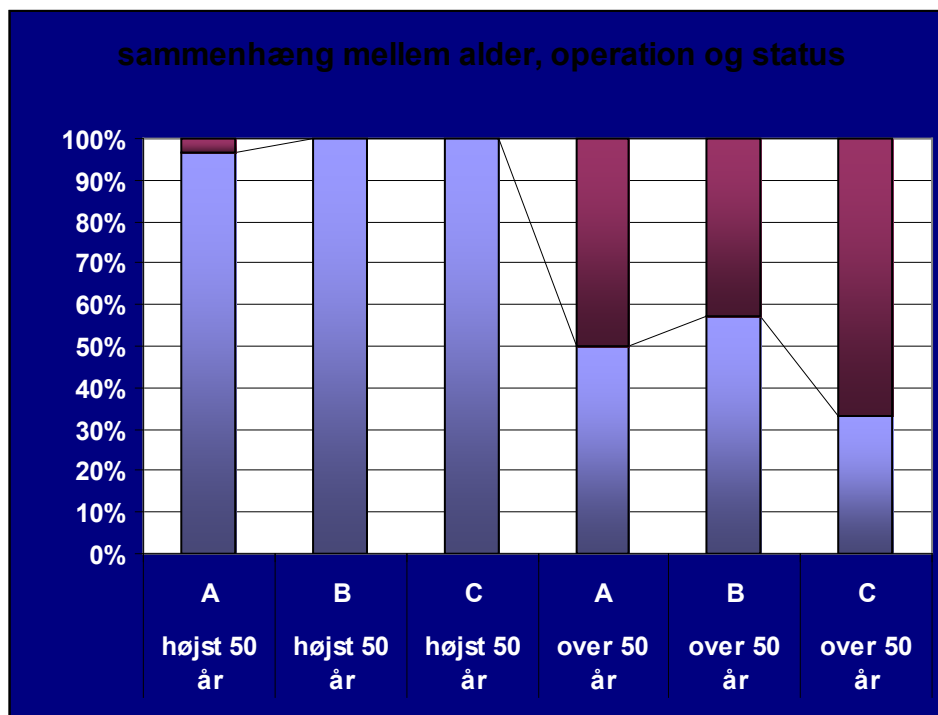
Over 50 år:

	I live	Død	Ialt
Operation A	13	13	26
Operation B	8	6	14
Operation C	8	16	24
Ialt	29	35	64

eller i procent

	I live	Død	Ialt
Operation A	50%	50%	100%
Operation B	57%	43%	100%
Operation C	33%	66%	100%

Grafisk ser det nu sådan ud:



Prøv at lave et nyt test for uafhængighed mellem overlevelse og operationsform for aldersgruppen over 50!

Når man tager alderen med i betragtning, er operation C pludselig ikke længere så stor en skurk.

Her sker der det, at operationstypen og alderen ikke er uafhængige af hinanden eller af overlevelseschancen. Når der er flere vigtige faktorer, der spiller ind på en gang, så bør de alle tages med i analysen, der så bliver noget mere kompliceret. En statistisk metode, man kan anvende, hedder loglineære modeller. Men den sag vil vi ikke komme ind på hér - det må vente til universitetet 😊.

Opgaver:

Opgave 1: En amerikansk undersøgelse af bilisters brug af sikkerhedsseler resulterede i følgende stikprøve:

Køn	Brug af sikkerheds sele			
	Altid	Som regel	Af og til	Aldrig
Mænd	37	60	54	64
Kvinder	39	58	49	39

Spørgsmål 1: Opstil den relevante nulhypotese og den alternative hypotese for at undersøge, om der er uafhængighed mellem køn og brug af sikkerhedssele?

Spørgsmål 2: Udregn tabellen med de forventede værdier og χ^2 teststørrelsen

Spørgsmål 3: Vil det være rimeligt at bruge en χ^2 -fordeling til at vurdere teststørrelsen her?

Spørgsmål 4: Giver stikprøven grundlag for at sige, at der er forskel på de to køns brug af sikkerhedsseler?

Opgave 2: En forretningskæde vil undersøge, om farven på indpakningen af nye kartofler påvirker salget. Butikken sælger derfor i en periode poser med samme slags kartofler, alle med 2.5 kg/pose og til samme pris, men i poser med forskellig farve.

Der bliver i alt sendt 600 poser kartofler ud i butikkerne, hvoraf 520 poser bliver solgt. Af de solgte poser er de 375 gule, og der er 55 gule poser tilbage. De øvrige poser er blå.

Undersøg, om der er grundlag for at påstå, at farven på posen påvirker salget af kartofler. Undervejs skal du formulere de relevante hypoteser, kommentere på begreber som signifikansniveau og/eller p-værdi og forklare den anvendte metode.

Opgave 3: En dyrlæge har fået en mistanke om, at hunde af racen labrador har en større tendens til at udvikle allergi end andre hunderacer. (Bemærk! Fiktivt eksempel)

Gennem det sidste år er der i klinikken blevet registret følgende undersøgelser og resultater af allergi hos hunde:

Race \ allergi	Ingen allergi	Mild allergi	Allergi
Labrador	25	2	10
Schæfer	32	0	7
Puddel	28	3	4
Andet	64	7	8

Spørgsmål 1: Opstil den relevante nulhypotese og den alternative hypotese for at undersøge, om der er uafhængighed mellem hunderace og udvikling af allergi?

Spørgsmål 2: Udregn tabellen med de forventede værdier.

Spørgsmål 3: Vil det være rimeligt at bruge en χ^2 -fordeling til at vurdere teststørrelsen her? Hvis ikke, hvordan kan man så komme videre med undersøgelsen?

En anden anvendelse af χ^2 -testet.

Danmarks statistiks opgørelse af indkomstfordelingen for personer over 15 år i Danmark år 2007 viser følgende billede:

I=Indkomst i 1000 kr.	I<50	50≤I<100	100≤I<150	150≤I<200	200≤I<300	300≤I<400	400≤I<500	500≤I
% af be-folkning	6.4	9.3	17.8	12.3	24.3	18.0	6.6	5.3

En markedsanalytiker har foretaget en undersøgelse af 1000 personers kendskab til et særdeles kostbart fladskærmsprodukt, men efterfølgende er der opstået tvivl om udvælgelsen af stikprøven, der er forgået som interviewundersøgelse over et par dage i et lokalt supermarked. Det frygtes, at stikprøven har fået for mange respondenter med i de lavere indkomstklasser. Heldigvis er der blevet spurgt om folks indkomst, så man kan lave et test for, om indkomstfordelingen i stikprøven synes at komme fra et specielt segment af befolkningen og altså dermed ikke at have den samme fordeling som indkomstfordelingen i Danmark. Hvis det er tilfældet, kan man nemlig ikke generalisere undersøgelsens resultat til hele befolkningen.

Indkomstfordelingen i stikprøven var:

Observerede antal:

I=Indkomst i 1000 kr.	I<50	50≤I<100	100≤I<150	150≤I<200	200≤I<300	300≤I<400	400≤I<500	500≤I
Antal i stikprøven	98	88	199	136	210	179	52	38

Modellen er følgende: Sandsynligheden for, at en tilfældig udvalgt person over 15 år tilhører en given indkomstklasse, er givet ved den procentdel af befolkningen, der tilhører denne klasse ifølge Danmarks statistiks indkomstfordeling. Sandsynligheden for, at en tilfældig udvalgt person har en indkomst på mindre end 50.000 kr. om året er fx 0.064. *Hypotesen* H_0 er, at vores stikprøve har en indkomstfordeling, der er den samme som den danske befolknings.¹⁰ Hvis vores stikprøve skal repræsentere en tilfældig stikprøve på den danske befolkning, så vil vi derfor forvente, at $1000 \cdot 0.064$ personer har en indkomst på mindre end 50.000 kr. På den måde kan vi som før beregne, hvor mange vi vil forvente i hver af de 8 indkomst kategorier:

Forventede antal under udgangshypotesen:

I=Indkomst i 1000 kr.	I<50	50≤I<100	100≤I<150	150≤I<200	200≤I<300	300≤I<400	400≤I<500	500≤I
Antal i stikprøven	64	93	178	123	243	180	66	53

Hypoteserne kan skrives.

H_0 : Indkomstfordelingen i stikprøven adskiller sig ikke signifikant fra indkomstfordelingen i populationen.

H_1 : Indkomstfordelingen i stikprøven er signifikant anderledes end indkomstfordelingen i populationen.

Teststørrelsen udregnes på samme vis som før ved at summere $\frac{(obs.antal - forv.antal)^2}{forv.antal}$ over alle indkomstgrupper.

Dvs. vi får:

$$\chi^2 = \frac{(98-64)^2}{64} + \frac{(88-93)^2}{93} + \frac{(199-178)^2}{178} + \frac{(136-123)^2}{123} + \frac{(210-243)^2}{243} + \frac{(179-180)^2}{180} + \frac{(52-66)^2}{66} + \frac{(38-53)^2}{53} = 33.87.$$

¹⁰ Det skal endnu engang understreges, at for at sikre repræsentativitet af en stikprøve så skal man følge de korrekte udvælgelsesmetoder og altså undgå denne form for "convenience sampling".

Teststørrelsen er denne gang χ^2 fordelt med 7 frihedsgrader. (Beregnes som antallet af indkomstgrupper minus 1). Vores matematiker kan fortælle os, at med et signifikansniveau på 5 % skal vi forkaste udgangshypotesen, når teststørrelsen er større end 14.07, og på 1% signifikansniveau, når den er større end 18.48.

P-værdien for testet er for alle praktiske formåls skyld 0.

Vi kan altså konkludere, at indkomstfordelingen i stikprøven afviger signifikant fra den generelle indkomstfordeling i Danmark, og at markedsanalytikeren skulle have fulgt reglerne for indsamling af repræsentative stikprøver.

Opgaver:

Opgave 1: Du har en mistanke om, at en af dine venner har en "falsk" terning.

Derfor har du i al hemmelighed noteret udfaldet af alle vedkommendes kast med terningen gennem en hel aftens spil. Dine optegnelser viser, at terningen er endt på "1" i alt 5 gange, "2" i alt 4 gange, "3" i alt 5 gange, "4" i alt 6 gange, "5" i alt 5 gange og "6" i alt 13 gange.

Giver dine observationer anledning til at din mistanke bestyrkes? Du forventes at argumentere ud fra statistiske hypoteser og test, med berøring af begreber som signifikans og/eller p-værdi. Desuden forventes du at kommentere, hvilke konsekvenser du vil lade den statistiske analyse få i den konkrete problemstilling.

Opgave 2: En mindre restaurant med et menukort bestående af 5 forskellige, men faste menuer plejer at have følgende ordrefordeling på disse:

menu 1: 30 %, menu 2: 25 %, menu 3: 20 %, menu 4: 15 % og menu 5: 10 %.

Restauranten foretager sine indkøb for at imødegå en efterspørgsel, der følger dette mønster. Imidlertid er man flere gange i den seneste tid løbet tør for menu 5, og man ønsker at afgøre, om det er en tilfældighed, eller om man skal revidere indkøbsplanerne.

I den seneste uge har man haft 543 gæster. Af disse bestilte 152 menu 1, 101 bestilte menu 2, 110 bestilte menu 3, 91 bestilte menu 4 og 89 bestilte menu 5.

Skal man revidere indkøbsplanerne?

Litteratur:

1. J. Burt & G. Barber. Elementary Statistics for Geographers. The Guildford press.
2. P. Newbold. Statistics for Business and Economics. Prentice Hall International Editions.
3. P. Mortensen. Repræsentative undersøgelser. Systime.
4. H.J. Beck, H.C. Hansen, A. Jørgensen, L.Ø. Petersen, P. Bollerslev. Matematik i læreruddannelsen. Gyldendals Uddannelse.