

Introduktion til Statistik

Forelæsning 10: Inferens for andele

Peder Bacher

DTU Compute, Dynamiske Systemer
Bygning 303B, Rum 010
Danmarks Tekniske Universitet
2800 Lyngby – Danmark
e-mail: pbac@dtu.dk

Forår 2021

Kapitel 7: Inferens for andele

Statistik for andele:

- Andel: $p = \frac{x}{n}$ (x *successer ud af n observationer*)
- Specifikke metoder, én, to og $k > 2$ grupper
 - Binær/kategorisk respons

Specifikke metoder:

- Estimation og konfidensintervaller for andele
 - Metoder korrektion ved små stikprøver
- Hypoteser for én andel (p)
- Hypoteser for to andele
- Analyse af antalstabeller (χ^2 -test) (alle forventede antal > 5)

Chapter 7: Inferences for Proportions

Statistics for proportions:

- Proportion: $p = \frac{x}{n}$ (*x successes out of n observations*)
- Specific methods: one, two and $k > 2$ samples:
 - Binary/categorical response

Specific methods:

- Estimation and confidence interval of proportions
 - Methods for correction for small samples
- Hypotheses for one proportion
- Hypotheses for two proportions
- Analysis of contingency tables (χ^2 -test) (all expected > 5)

Overview

- 1 Intro
- 2 Konfidensinterval for én andel
 - Eksempel 1
- 3 Hypotesetest for én andel
 - Eksempel 1 - fortsat
- 4 Konfidensinterval og hypotesetest for forskel på to andele
 - Eksempel 2
- 5 Hypotesetest for flere andele
 - Eksempel 2 - fortsat
- 6 Analyse af antalstabeller

Forskellige analyse/data-situationer

Hypotesetests og konfidensintervaller for:

- Én middelværdi (*one-sample, i.e. one group/population*)
- To middelværdier (*two-sample, i.e. two groups/populations*)
- Næste uge: For flere middelværdier (*k-sample, i.e. k groups/populations*)

I dag: Hypotesetests og konfidensintervaller for:

- Én andel
- To andele
- Flere andele (kun hypotesetest)
- Flere "multi-categorical" andele (kun hypotesetest)

Estimation af andele

Estimation af andele fås ved at observere antal gange x en hændelse har indtruffet ud af n forsøg:

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

$$\hat{p} \in [0; 1]$$

Spørgsmål om andel (socrative.com, ROOM: pbac)

Hvilken kan ikke en være en andel?

- A: $103/900$
- B: $12/80$
- C: 0.957
- D: $202/154$
- E: 0.224

Konfidensinterval for én andel

Method 7.3

Såfremt der haves en stor stikprøve, fås et $(1 - \alpha)\%$ konfidensinterval for p

$$\left[\hat{p} - z_{1-\alpha/2} \cdot \hat{\sigma}_{\hat{p}} , \quad \hat{p} + z_{1-\alpha/2} \cdot \hat{\sigma}_{\hat{p}} \right] \quad \left[\hat{p} - z_{1-\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} , \quad \hat{p} + z_{1-\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \right]$$

(Vi siger: Med stor sikkerhed vælger vi at tro at p i dette interval)

Hvordan?

Følger af at approximere binomialfordelingen med normalfordelingen

As a rule of thumb

The normal distribution gives a good approximation of the binomial distribution if np and $n(1-p)$ are both greater than 15

Konfidensinterval for én andel

Middelværdi og varians i binomialfordelingen, kapitel 2:

$$\begin{aligned}E(X) &= np \\ \text{Var}(X) &= np(1 - p)\end{aligned}$$

Derfor får man

$$\begin{aligned}E(\hat{p}) &= E\left(\frac{X}{n}\right) = \frac{np}{n} = p \\ \text{Var}(\hat{p}) &= \sigma_{\hat{p}}^2 = \text{Var}\left(\frac{X}{n}\right) = \frac{1}{n^2} \text{Var}(X) = \frac{p(1 - p)}{n}\end{aligned}$$

Eksempel 1

Venstrehåndede:

$p =$ Andelen af venstrehåndede i Danmark

eller:

Kvindelige ingeniørstuderende:

$p =$ Andelen af kvindelige ingeniørstuderende

Eksempel 1

Venstrehåndede ($x = 10$ ud af $n = 100$):

$$\hat{\sigma}_{\hat{p}} = \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} = \sqrt{\frac{10/100(1-10/100)}{100}} = 0.03$$
$$0.10 \pm 1.96 \cdot 0.03 \Leftrightarrow 0.10 \pm 0.06 \Leftrightarrow [0.04, 0.16]$$

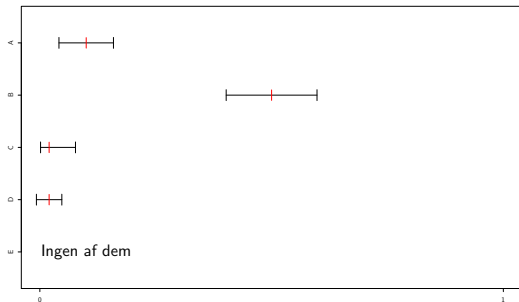
Bedre "small sample" metode - "plus 2-approach" (Remark 7.7):

Anvend samme formel på $\tilde{x} = 10 + 2 = 12$ og $\tilde{n} = 104$:

$$\sqrt{\frac{\tilde{p}(1-\tilde{p})}{\tilde{n}}} = \sqrt{\frac{12/104(1-12/104)}{104}} = 0.0313$$
$$0.115 \pm 1.96 \cdot 0.0313 \Leftrightarrow 0.115 \pm 0.061 \Leftrightarrow [0.054, 0.18]$$

Spørgsmål om plus 2-approach (socrative.com, ROOM: pbac)

Hvilket af følgende intervaller er med plus 2-approach?



Trin ved Hypotesetest

Trin ved Hypotesetest:

1. Opstil hypoteser og vælg signifikansniveau α
 2. Beregn teststørrelse
 3. Beregn p -værdi (eller kritisk værdi)
 4. Fortolk p -værdi og/eller sammenlign p -værdi og signifikansniveau, og derefter drag en konklusion
- (Alternativ 4. Sammenlign teststørrelse og kritisk værdi og drag en konklusion)

Hypotesetest for én andel

Vi betragter en nul- og alternativ hypotese for én andel p :

$$H_0 : p = p_0$$

$$H_1 : p \neq p_0$$

Man vælger som sædvanligt enten at acceptere H_0 eller at forkaste H_0

Theorem 7.10 og Method 7.11

Såfremt stikprøven er tilstrækkelig stor ($np_0 > 15$ og $n(1 - p_0) > 15$) bruges teststørrelsen:

$$z_{\text{obs}} = \frac{x - np_0}{\sqrt{np_0(1 - p_0)}}$$

Under nulhypotesen gælder at den tilsvarende tilfældige variabel Z følger en standard normalfordeling, dvs. $Z \sim N(0, 1^2)$

Test ved brug af p -værdi (Method 7.11)

Find p -værdien (bevis mod nulhypotesen):

- We only use two-sided: $2P(Z > |z_{\text{obs}}|)$ in exercises and exams
- Remark 7.9 om one-sided "less" og "greater"

Kritiske værdier

Alternativ hypotese	Afvis nulhypotese hvis
$p \neq p_0$	$z_{\text{obs}} < -z_{1-\alpha/2}$ eller $z_{\text{obs}} > z_{1-\alpha/2}$

Eksempel 1 - fortsat

Er halvdelen af alle danskere venstrehåndede?

$$H_0 : p = 0.5, \quad H_1 : p \neq 0.5$$

Teststørrelse:

$$z_{\text{obs}} = \frac{x - np_0}{\sqrt{np_0(1-p_0)}} = \frac{10 - 100 \cdot 0.5}{\sqrt{100 \cdot 0.5(1-0.5)}} = -8$$

Er p -værdien under 0.05? (dvs. skal nulhypotesen forkastes ved $\alpha = 0.05$)

A: Ja B: Nej C: Ved ikke

R: prop.test - een andel

```
## Single proportion  
  
## Testing the probability = 0.5 with a two-sided alternative  
## We have observed 518 out of 1154  
## Without continuity corrections  
  
prop.test(x=518, n=1154, p = 0.5, correct = FALSE)
```

Konfidensinterval for forskel på to andele

Method 7.15

$$(\hat{p}_1 - \hat{p}_2) \pm z_{1-\alpha/2} \cdot \hat{\sigma}_{\hat{p}_1 - \hat{p}_2}$$

hvor

$$\hat{\sigma}_{\hat{p}_1 - \hat{p}_2} = \sqrt{\frac{\hat{p}_1(1 - \hat{p}_1)}{n_1} + \frac{\hat{p}_2(1 - \hat{p}_2)}{n_2}}$$

Rule of thumb:

Både $n_i p_i \geq 10$ and $n_i(1 - p_i) \geq 10$ for $i = 1, 2$

Hypotesetest for forskel på to andele, Method 7.18

Two sample proportions hypothesis test

Såfremt man ønsker at sammenligne to andele (her vist for et tosidet alternativ)

$$H_0 : p_1 = p_2$$

$$H_1 : p_1 \neq p_2$$

Fås teststørrelsen:

$$z_{\text{obs}} = \frac{\hat{p}_1 - \hat{p}_2}{\sqrt{\hat{p}(1 - \hat{p})\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}}, \quad \text{hvor} \quad \hat{p} = \frac{x_1 + x_2}{n_1 + n_2}$$

Og for passende store stikprøver:

Brug standardnormalfordelingen igen

Eksempel 2

Sammenhæng mellem brug af p-piller og risikoen for blodprob i hjertet (hjerteinfarkt)

I et studie (USA, 1975) undersøgte man dette. Fra et hospital havde man indsamlet følgende to stikprøver

	p-piller	Ikke p-piller
Blodprob	23	35
Ikke blodprob	34	132

Er der sammenhæng mellem brug af p-piller og sygdomsrisiko

Udfør et test for om der er sammenhæng mellem brug af p-piller og risiko for blodprob i hjertet. Anvend signifikansniveau $\alpha = 5\%$.

Eksempel 2

Sammenhæng mellem brug af p-piller og risikoen for blodprob i hjertet

	p-piller	Ikke p-piller
Blodprob	$x_1 = 23$	$x_2 = 35$
Ikke blodprob	34	132
Sum	$n_1 = 57$	$n_2 = 167$

Estimater i hver stikprøve

$$\hat{p}_1 = \frac{x_1}{n_1} = \frac{23}{57} = 0.40, \quad \hat{p}_2 = \frac{x_2}{n_2} = \frac{35}{167} = 0.21$$

R: prop.test - to andele

```
## Pill study: two proportions
```

```
## Reading the table into R
```

```
pill.study <- matrix(c(23, 34, 35, 132), ncol = 2)  
rownames(pill.study) <- c("Blood Clot", "No Clot")  
colnames(pill.study) <- c("Pill", "No pill")
```

```
## Testing that the probabilities for the two groups are equal  
prop.test(t(pill.study), correct = FALSE)
```

```
## Or simply directly by
```

```
prop.test(x=c(23,35), n=c(57,167), correct = FALSE)
```

Hypotesetest for flere andele

Sammenligning af c andele

I nogle tilfælde kan man være interesseret i at vurdere om to eller flere binomialfordlinger har den samme parameter p , dvs. man er interesseret i at teste nulhypotesen

$$H_0: p_1 = p_2 = \dots = p_c = p$$

mod en alternativ hypotese at disse andele ikke er ens

Hypotesetest for flere andele

Tabel af observerede antal for c stikprøver:

	stikprøve 1	stikprøve 2	...	stikprøve c	Total
Succes	x_1	x_2	...	x_c	x
Fiasko	$n_1 - x_1$	$n_2 - x_2$...	$n_c - x_c$	$n - x$
Total	n_1	n_2	...	n_c	n

Fælles (gennemsnitlig) estimat:

Under nulhypotesen fås et estimat for p

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

Hypotesetest for flere andele

Fælles (gennemsnitlig) estimat:

Under nulhypotesen fås et estimat for p

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

"Brug" dette fælles estimat i hver gruppe:

såfremt nulhypotesen gælder, vil vi forvente at den j 'te gruppe har e_{1j} successer og e_{2j} fiaskoer, hvor

$$e_{1j} = n_j \cdot \hat{p} = n_j \cdot \frac{x}{n}$$

$$e_{2j} = n_j(1 - \hat{p}) = n_j \cdot \frac{n - x}{n}$$

Hypotesetest for flere andele

Generel formel for beregning af forventede værdier i antalstabeller:

$$e_{ij} = (j\text{'th column total}) \cdot \frac{(i\text{'th row total})}{(\text{total})}$$

Beregning af teststørrelse - Method 7.20

Teststørrelsen bliver

$$\chi_{\text{obs}}^2 = \sum_{i=1}^2 \sum_{j=1}^c \frac{(o_{ij} - e_{ij})^2}{e_{ij}}$$

hvor o_{ij} er observeret antal i celle (i,j) og e_{ij} er forventet antal i celle (i,j)

Find p -værdi eller brug kritisk værdi - Method 7.20

Stikprøvefordeling for test-størrelse:

χ^2 -fordeling med $(c - 1)$ frihedsgrader

Kritisk værdi metode

Såfremt $\chi_{\text{obs}}^2 > \chi_{1-\alpha}^2(c - 1)$ forkastes nulhypotesen

Rule of thumb for validity of the test:

Alle forventede værdier $e_{ij} \geq 5$

Eksempel 2 - fortsat

De OBSERVEREDE værdier o_{ij}

	p-piller	Ikke p-piller	Total
Blodprob	23	35	
Ikke blodprob	34	132	

Eksempel 2 - fortsat

Beregn de FORVENTEDE værdier e_{ij} (altså forventede *under* H_0)

	p-piller	Ikke p-piller	Total
Blodprob			$x = 58$
Ikke blodprob			
	$n_1 = 57$	$n_2 = 167$	$n = 224$

Eksempel 2 - fortsat

Beregn de FORVENTEDE værdier e_{ij} (altså forventede *under* H_0)

	p-piller	Ikke p-piller	Total
Blodprob	14.76	43.24	$x = 58$
Ikke blodprob	42.24	123.76	
	$n_1 = 57$	$n_2 = 167$	$n = 224$

Brug "reglen" for forventede værdier fire gange, f.eks. :

$$e_{12} = 167 \cdot \frac{58}{224} = 43.24$$

Eksempel 2 - fortsat

Teststørrelsen:

$$\chi_{\text{obs}}^2 = \frac{(o_{11} - e_{11})^2}{e_{11}} + \frac{(o_{12} - e_{12})^2}{e_{12}} + \frac{(o_{21} - e_{21})^2}{e_{21}} + \frac{(o_{22} - e_{22})^2}{e_{22}}$$
$$=$$

$$\chi_{\text{obs}}^2 = \frac{(23 - 14.76)^2}{14.76} + \frac{(35 - 43.24)^2}{43.24} + \frac{(34 - 42.24)^2}{42.24} + \frac{(132 - 123.76)^2}{123.76}$$
$$= 8.33$$

Kritisk værdi og p -værdi:

```
## Kritisk værdi  
qchisq(0.95, 1)  
  
## [1] 3.8
```

```
## p-værdi  
1 - pchisq(8.33, df=1)  
  
## [1] 0.0039
```


R: chisq.test - to andele

```
## Pill study: two proportions, chi-square test  
  
## Chi2 test for testing the probabilities for the two groups are equal  
chisq.test(pill.study, correct = FALSE)  
## If we want the expected numbers save the test in an object  
chi <- chisq.test(pill.study, correct = FALSE)  
## The expected values  
chi$expected
```

Antalstabeller

Antalstabel

- Flere end 2 kategorier (f.eks. fire.: rød, grøn, blå, sort)
- Beregningerne er ens for begge følgende setups

To mulige setups

- Setup 1: c stikprøver med r kategorier:
 - Test om der er forskel i fordelingen mellem kategorierne for hver stikprøve
- Setup 2: To kategoriske variabel (r kategorier) målt på samme individer (parret setup):
 - Test om der er forskel i fordelingen mellem de to grupper

Setup 1: c stikprøver med r kategorier

En 3×3 tabel - 3 stikprøver, 3-kategori udfald

	4 uger før	2 uger før	1 uge før
Kandidat I	79	91	93
Kandidat II	84	66	60
ved ikke	37	43	47
	$n_1 = 200$	$n_2 = 200$	$n_3 = 200$

Er stemmefordelingen ens?

$$H_0: p_{i1} = p_{i2} = p_{i3}, i = 1, 2, 3$$

Setup 2: To kategoriske variabel (r kategorier) målt på samme individer (parret setup)

En 3×3 tabel - 1 stikprøve, to stk. 3-kategori variable:

	dårlig	middel	god
dårlig	23	60	29
middel	28	79	60
god	9	49	63

Er der uafhængighed mellem inddelingskriterier?

$$H_0 : p_{ij} = p_{i \cdot} p_{\cdot j}$$

f.eks. er der sammenhæng mellem den måde elever klarer sig i matematik som i dansk?

Beregning af teststørrelse – uanset type af tabel

I en antalstable med r rækker og c søjler, fås teststørrelsen

$$\chi^2_{\text{obs}} = \sum_{i=1}^r \sum_{j=1}^c \frac{(o_{ij} - e_{ij})^2}{e_{ij}}$$

hvor o_{ij} er observeret antal i celle (i,j) og e_{ij} er forventet antal i celle (i,j)

Generel formel for beregning af forventede værdier i antalstabeller:

$$e_{ij} = (j\text{'th column total}) \cdot \frac{(i\text{'th row total})}{(\text{total})}$$

Spørgsmål (socrative.com, ROOM: pbac)

En 3×4 tabel - 4 stikprøver, 3-kategori udfald

	Gruppe A	Gruppe B	Gruppe C	Gruppe D	n_j
Han	3	3	2	2	10
Hun	3	3	5	2	13
Tvekøn	4	4	3	6	17
n_i	10	10	10	10	40

Hvad er e_{23} ? (forventning af hunner i gruppe C under H_0)

- A: $10 \cdot 10/40$
- B: 3
- C: $10 \cdot 13/40$
- D: $17 \cdot 4/40$
- E: Ved ikke

Find p -værdi eller brug kritisk værdi – Method 7.22

Stikprøvefordeling for test-størrelse:

χ^2 -fordeling med $(r-1)(c-1)$ frihedsgrader

Kritisk værdi metode

Såfremt $\chi_{\text{obs}}^2 > \chi_{1-\alpha}^2$ med $(r-1)(c-1)$ frihedsgrader forkastes nulhypotesen

Rule of thumb for validity of the test:

Alle forventede værdier $e_{ij} \geq 5$

R: chisq.test - antalstabeller

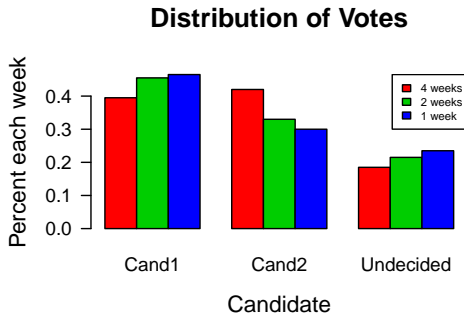
```
## Poll study: contingency table, chi-square test

## Reading the table into r
poll <- matrix(c(79, 91, 93, 84, 66, 60, 37, 43, 47), ncol = 3, byrow = TRUE)
colnames(poll) <- c("4 weeks", "2 weeks", "1 week")
rownames(poll) <- c("Cand1", "Cand2", "Undecided")

## Column percentages
colpercent <- prop.table(poll, 2)
colpercent
```


R: chisq.test - antalstabeller

```
barplot(t(colpercent), beside = TRUE, col = 2:4, las = 1,
        ylab = "Percent each week", xlab = "Candidate",
        main = "Distribution of Votes")
legend( legend = colnames(poll), fill = 2:4,"topright", cex = 0.5)
par(mar=c(5,4,4,2)+0.1)
```



R: chisq.test - antalstabeller

```
## Testing same distribution in the three populations  
chi <- chisq.test(poll, correct = FALSE)  
chi  
## Expected values  
chi$expected
```