

# Formelsamling i medisinsk statistikk

- Dette er en formelsamling til O. O. Aalen (red.): Statistiske metoder i medisin og helsefag, Gyldendal, 2006.
- Merk at boken har en nettside der det er lagt ut rettelser og supplerende stoff, se <http://www.med.uio.no/imb/stat/statbok/>

## Gjennomsnitt

$$\bar{x} = \frac{1}{n}(x_1 + x_2 + x_3 + \cdots + x_n)$$

## Median

Alle observasjoner ordnes i stigende rekkefølge. Ved *ulike* antall observasjoner, er medianen definert som den midterste av dem. Ved *like* antall, er medianen definert som gjennomsnittet av de to midterste.

## Standardavvik

$$s = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$$

## Grupperte data

Intervallmidtpunkter  $m_1, m_2, \dots, m_k$ . Hyppigheter  $f_1, f_2, \dots, f_k$ . Totalt antall observasjoner:  $n$ . Gjennomsnitt og standardavvik er gitt ved:

$$\bar{x} = \frac{1}{n}(m_1 f_1 + m_2 f_2 + \cdots + m_k f_k) = \frac{1}{n} \sum_{j=1}^k m_j f_j$$

$$s = \sqrt{\frac{1}{n-1} \left\{ \sum_{j=1}^k (m_j - \bar{x})^2 f_j \right\}}$$

*Median* og *fraktiler* for grupperte data finnes ved lineær interpolasjon.

## Insidens og prevalens

*Prevalens* angir andelen i befolkningen som har en viss sykdom.

*Insidensraten* beregnes som antall *nye* tilfeller av sykdommen over et tidsintervall, dividert med totalt antall personår under risiko.

## Regneregler for sannsynlighet

Hvis begivenhetene  $A$  og  $B$  er disjunkte has

$$P(A \cup B) = P(A) + P(B)$$

For alle begivenheter  $A$  og  $B$  has

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

Definisjon av *betinget sannsynlighet*

$$P(A | B) = \frac{P(A \cap B)}{P(B)}$$

Begivenhetene  $A$  og  $B$  er *uavhengige* hvis

$$P(A \cap B) = P(A) \cdot P(B)$$

En tilsvarende produktregel er gyldig om vi har flere uavhengige begivenheter.

Regelen om *total sannsynlighet*

$$P(A) = P(A | B) \cdot P(B) + P(A | \bar{B}) \cdot P(\bar{B})$$

*Bayes' lov*

$$P(B | A) = \frac{P(B)P(A | B)}{P(B)P(A | B) + P(\bar{B})P(A | \bar{B})}$$

## Diagnostiske tester

*Sensitivitet:* Sannsynlighet for positiv test gitt at det foreligger sykdom.

*Spesifisitet:* Sannsynlighet for negativ test gitt at det ikke foreligger sykdom.

*Positiv prediktiv verdi:* Sannsynlighet for at det foreligger sykdom gitt positiv test.

*Negativ prediktiv verdi:* Sannsynlighet for at det ikke foreligger sykdom gitt negativ test.

## Kombinatorikk

Trekning av  $s$  kuler fra en boks med  $n$  kuler.

Antall *ordnede* utvalg *med tilbakelegging*

$$n^s$$

Antall *ordnede* utvalg *uten tilbakelegging*

$$n(n-1)(n-2) \cdots (n-s+1)$$

Antall *ikke-ordnede* utvalg *uten tilbakelegging*

$$\binom{n}{s} = \frac{n(n-1)(n-2) \cdots (n-s+1)}{s!}$$

## Forventning og varians for teoretisk fordeling

$$E(X) = \sum_{\text{alle } x_i} x_i P(X = x_i)$$

$$\text{Var}(X) = \sum_{\text{alle } x_i} (x_i - E(X))^2 P(X = x_i)$$

## Regneregler for forventning og varians

$$E(aX + b) = aE(X) + b,$$

$$\text{Var}(aX + b) = a^2 \text{Var}(X), \quad \text{SD}(aX + b) = |a| \text{SD}(X)$$

$$E(X_1 + X_2 + \dots + X_n) = E(X_1) + E(X_2) + \dots + E(X_n)$$

Hvis  $X_1, X_2, \dots, X_n$  er parvis *stokastisk uavhengige* has:

$$\text{Var}(X_1 + X_2 + \dots + X_n) = \text{Var}(X_1) + \text{Var}(X_2) + \dots + \text{Var}(X_n)$$

## Binomisk fordeling

Sannsynligheten for at en begivenhet  $A$  inntreffer  $x$  ganger i løpet av  $n$  binomiske forsøk, er

$$P(X = x) = \binom{n}{x} p^x (1-p)^{n-x}, \quad x = 0, 1, \dots, n$$

Forventning og varians i binomisk fordeling er gitt ved:

$$E(X) = np, \quad \text{Var}(X) = np(1-p)$$

## Poissonfordeling

Sannsynligheten for  $x$  forekomster, når forventning er lik  $\lambda$ , er gitt ved:

$$P(X = x) = \frac{\lambda^x e^{-\lambda}}{x!} \quad \text{for } x = 0, 1, 2, \dots$$

Forventning og varians er gitt ved:

$$E(X) = \lambda \quad \text{og} \quad \text{Var}(X) = \lambda$$

Poissonfordelingen anvendes også ved Poissonprosesser.

## Normalfordeling

En stokastisk variabel  $X$  sies å være normal  $(\mu, \sigma)$  hvis den følger en normalfordeling med forventning (sentrum)  $\mu$  og standardavvik (spredning)  $\sigma$ . Den standardiserte variable  $Y = (X - \mu)/\sigma$  er normal  $(0,1)$ . Sannsynlighetstettheten til normalfordelingen er gitt ved følgende formel:

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

der  $\exp(a)$  er det samme som eksponensialfunksjonen  $e^a$ .

## Formler for gjennomsnitt

La  $\bar{X}$  være gjennomsnittet av de uavhengige variablene  $X_1, X_2, \dots, X_n$ . Da gjelder:

$$E(\bar{X}) = \mu, \quad \text{Var}(\bar{X}) = \frac{\sigma^2}{n}, \quad \text{SD}(\bar{X}) = \frac{\sigma}{\sqrt{n}}, \quad s_{\bar{X}} = \frac{s}{\sqrt{n}}$$

Hvis variablene også er normalfordelte, vil et konfidensintervall være gitt ved

$$\bar{X} \pm c s_{\bar{X}}$$

der  $c$  bestemmes ut fra Studentfordelingen med  $n - 1$  frihetsgrader.

En teststørrelse er gitt ved

$$t = \frac{\bar{X} - a}{s_{\bar{X}}} = \frac{\bar{X} - a}{s} \sqrt{n}$$

og denne er Studentfordelt med  $n - 1$  frihetsgrader når  $H_0: \mu = a$  gjelder.

## Sammenlikning av pardata

Man tar differansen innenfor hvert par og bruker konfidensintervallet og teststørrelsen over med  $a = 0$ . Forutsetningen er at differansene er uavhengige og normalfordelte.

## Sammenlikning av to gjennomsnitt

Vi forutsetter uavhengige og normalfordelte observasjoner. Forøvrig antas gjennomsnittene å komme fra to uavhengige utvalg. Følgende teststørrelse er Studentfordelt med  $n_1 + n_2 - 2$  frihetsgrader når  $H_0$  gjelder

$$t = \frac{\bar{X}_1 - \bar{X}_2}{s_f \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

der  $s_f$  er definert ved

$$s_f = \sqrt{\frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}}$$

Et konfidensintervall er gitt ved

$$\bar{X}_1 - \bar{X}_2 \pm c s_f \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$$

der  $c$  er bestemt av Studentfordelingen med  $n_1 + n_2 - 2$  frihetsgrader.

## Poissonfordeling som tilnærming til binomisk fordeling

Binomisk fordeling kan tilnærmes med en Poissonfordeling hvis:

$$(1) \quad p \leq 0.05 \quad \text{og} \quad (2) \quad n \geq 50$$

## Normalfordeling som tilnærming til binomisk fordeling

Når  $n$  i en binomisk fordeling er så stor at  $np \geq 5$  og  $n(1-p) \geq 5$ , vil den binomiske fordelingen likne mye på en normalfordeling med parametre

$$\mu = np, \quad \sigma = \sqrt{np(1-p)}$$

## Normalfordeling som tilnærming til Poissonfordeling

Når  $\lambda$  i en Poissonfordeling er minst lik 5, vil Poissonfordelingen likne mye på en normalfordeling med parametre

$$\mu = \lambda, \quad \sigma = \sqrt{\lambda}$$

## Estimering av sannsynlighet (andel)

Hvis det er observert  $X$  forekomster ved  $n$  binomiske forsøk, er estimatet for sannsynligheten  $p$  gitt ved  $p^*$ , mens estimert standardfeil er gitt ved  $s_p$

$$p^* = X/n, \quad s_p = \sqrt{\frac{p^*(1-p^*)}{n}}$$

Fordelingen til  $p^*$  er tilnærmet normalfordelt under de samme forutsetninger som for binomisk fordeling, med  $\mu = p$  og  $\sigma = \sqrt{\frac{p(1-p)}{n}}$ .

Et 95% konfidensintervall for  $p^*$  er gitt ved

$$p^* \pm 2s_p$$

## Testing av nullhypotese om en sannsynlighet

$$Z = \frac{X - np_0}{\sqrt{np_0(1-p_0)}}$$

## Teststørrelse for sammenlikning av to sannsynligheter (andeler)

$$Y = \frac{p_1^* - p_2^*}{\sqrt{(\frac{1}{n_1} + \frac{1}{n_2})\bar{p}(1-\bar{p})}}$$

### Konfidensintervall for differanse mellom to andeler

$$p_1^* - p_2^* \pm 1.96 \sqrt{\frac{p_1^*(1-p_1^*)}{n_1} + \frac{p_2^*(1-p_2^*)}{n_2}}$$

### Teststørrelse for sammenlikning av to Poissonvariabler

$$Y = \frac{X_1 - X_2}{\sqrt{X_1 + X_2}}$$

### Konfidensintervall for relativ risiko

Relativ risiko:

$$RR = \frac{a/(a+c)}{b/(b+d)}$$

Hjelpstørrelse:

$$s_{RR} = \sqrt{\frac{1}{a} + \frac{1}{b} - \frac{1}{a+c} - \frac{1}{b+d}}$$

95% konfidensintervall for  $RR$ :

$$(RR \times e^{-1.96 s_{RR}}, RR \times e^{1.96 s_{RR}})$$

### Konfidensintervall for odds-ratio

Odds-ratio:

$$OR = \frac{a/c}{b/d} = \frac{a \cdot d}{b \cdot c}$$

Hjelpstørrelse:

$$s_{OR} = \sqrt{\frac{1}{a} + \frac{1}{b} + \frac{1}{c} + \frac{1}{d}}$$

95% konfidensintervall for  $OR$ :

$$(OR \times e^{-1.96 s_{OR}}, OR \times e^{1.96 s_{OR}})$$

### Kji-kvadrattest

Kji-kvadrattesten for en  $2 \times 2$ -tabell kan beregnes ut fra følgende formel:

$$\chi^2 = \frac{N(ad - bc)^2}{(a+b)(a+c)(b+d)(c+d)}$$

Formelen er basert på oppsettet i tabellen øverst s. 130 i læreboken, der  $N$  er summen av tallene i tabellen. Formelen er ikke gitt i boken, men gir samme svar som beregningen av størrelsen  $S$  på s. 136.

## Regresjonsanalyse

Helningskoeffisienten,  $b$ , og skjæringspunktet med  $y$ -aksen,  $a$ , for minste-kvadraterslinjen er gitt ved

$$\hat{b} = s_{xy}/s_x^2, \quad \hat{a} = \bar{y} - \hat{b}\bar{x}$$

der  $s_x$  og  $s_y$  er standardavvikene til henholdsvis  $x$ - og  $y$ -verdiene, mens  $s_{xy}$  er definert ved

$$s_{xy} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

Minste kvadratsum er gitt ved

$$\text{rest} = \sum_{i=1}^n (y_i - \hat{a} - \hat{b}x_i)^2 = (n-1)(s_y^2 - \hat{b}s_{xy})$$

Standardavvik som måler variasjonen i punktene rundt den beste linjen:

$$s_{\text{reg}} = \sqrt{\frac{\text{rest}}{n-2}}$$

Konfidensintervall for  $\hat{b}$  bestemmes ut fra formelen:

$$\hat{b} \pm c \frac{s_{\text{reg}}}{\sqrt{(n-1)s_x^2}}$$

der  $c$  bestemmes ut fra en studentfordeling med  $n-2$  frihetsgrader.

## Korrelasjon

Korrelasjonskoeffisienten er definert på følgende måte:

$$r = \frac{s_{xy}}{s_x s_y}$$

## Utvalgsstørrelse

Parallellgruppestudie – målevariabler:

$$n = 2 \left( \frac{\sigma}{\Delta} \right)^2 \cdot k$$

Overkrysningsstudie – målevariabler:

$$n = \left( \frac{\sigma_d}{\Delta} \right)^2 \cdot k$$

Utvalgsstørrelse – binomisk responsvariabel:

$$n = \frac{p_1(1-p_1) + p_2(1-p_2)}{(p_2 - p_1)^2} \cdot k$$

$k$  bestemmes av tabellen:

|               |      | Teststyrke |      |      |
|---------------|------|------------|------|------|
|               |      | 0.80       | 0.90 | 0.95 |
| Siginifikans- | 0.10 | 6.2        | 8.6  | 10.8 |
| nivå          | 0.05 | 7.9        | 10.5 | 13.0 |
| (tosidig)     | 0.01 | 11.7       | 14.9 | 17.8 |

## Utvalgsstørrelse basert på presisjon i estimater

Binomisk respons:

$$n = \left( \frac{1.96}{d} \right)^2 p(1-p)$$

Kontinuerlig respons:

$$n = \left( \frac{1.96 \cdot \sigma}{d} \right)^2$$

## Infeksjonsmodellering

Kritisk vaksinasjonsdekning for å unngå en epidemi:

$$p_{\text{crit}} = 1 - \frac{1}{R_0}$$