

Kontingenstabeller

A frequency table.

	B_1	B_2	\dots	B_c	Total
A_1	y_{11}	y_{12}	\dots	y_{1c}	y_{1+}
A_2	y_{21}	y_{22}	\dots	y_{2c}	y_{2+}
\vdots	\vdots	\vdots	\ddots	\vdots	\vdots
A_r	y_{r1}	y_{r2}	\dots	y_{rc}	y_{r+}
Total	y_{+1}	y_{+2}	\dots	y_{+c}	$y_{++} = N$

To faktorer, A og B



Fraktkonst af data

(1) Poissonfordelte observationer

$$Y_{hk} \sim p(\mu_{hk}) \text{ uafh.}$$

Simultant:

$$P(Y=y) = \prod_{h=1}^H \prod_{k=1}^{n_h} \frac{e^{-\mu_{hk}} \mu_{hk}^{y_{hk}}}{y_{hk}!}$$

(2) med $\sum_h \sum_k y_{hk} = N$ fast

$$\begin{aligned}
 & P(Y=y) \mid \sum_i Y_i = N \quad (\quad Y_i = Y_{hk} \quad) \\
 & = \frac{P(Y=y)}{P(\sum_i Y_i = N)} \quad \begin{matrix} / \\ i=1, \dots, N \end{matrix} \quad \begin{matrix} \backslash \\ h=1, \dots, H \\ k=1, \dots, c \end{matrix} \\
 & = \left[\prod_{h=1}^H \prod_{k=1}^{n_h} \frac{e^{-\mu_{hk}} \mu_{hk}^{y_{hk}}}{y_{hk}!} \right] \frac{N!}{e^{-\mu_{..}} \mu_{..}^N}
 \end{aligned}$$

$$\begin{aligned}
 &= N! \prod_{h,k} \frac{1}{y_{hk}!} \left(\frac{\mu_{hk}}{\mu_{..}} \right)^{y_{hk}} \\
 &= (y_{11}, \dots, y_{rc}) \prod_{h,k} \pi_{hk}^{y_{hk}}, \quad \pi_{hk} = \frac{\mu_{hk}}{\mu_{..}}
 \end{aligned}$$

dvs.

$$(Y_{11}, \dots, Y_{rc}) \sim m(N; \pi_{11}, \dots, \pi_{rc})$$

↑
multinomialfordelt

(3) med $\sum_h y_{hk} = y_{.k}$ fast for alle k

$$\begin{aligned}
 &P(Y=y | Y_{.k} = y_{.k}, k=1, \dots, c) \\
 &= \left[\prod_{h,k} \frac{e^{-\mu_{hk}} \mu_{hk}^{y_{hk}}}{y_{hk}!} \right] \prod_k \frac{y_{.k}!}{e^{-\mu_{.k}} \mu_{.k}^{y_{.k}}} \\
 &= \prod_k \left[\left(\frac{y_{.k}}{y_{1k}, \dots, y_{ck}} \right) \prod_h \pi_{hk}^{y_{hk}} \right], \quad \pi_{hk} = \frac{\mu_{hk}}{\mu_{.k}}
 \end{aligned}$$

produkt af c multinomialsands.,

dvs.

$$(Y_{1k}, \dots, Y_{ck}) \sim m(y_{.k}; \pi_{11k}, \dots, \pi_{1ck}),$$

$k=1, \dots, c$
nafte

Log-lineare modeller

En log-lineær model er en GLM med logaritmisk link, dvs. $y_i = \ln E[Y_i] = \ln \mu_i$

$$(y_i = x_i^\top \beta)$$

$g(\mu) = \ln \mu$ er harmonisk link for Poissonfordelingen. Log-lineær modeller er derfor egnete til analyse af kontingenstabeller.

Betrægt $Y_i \sim p(\mu_i)$

$$E[Y_i] = \mu_i = \mu_{..} \frac{\mu_i}{\mu_{..}} = \mu_{..} \pi_i$$

Test af uafhængighed mellem faktor A og faktor B i kontingenstabell.

Hypothesen H_0 : $\pi_{hk} = \pi_{h..} \pi_{..k}$ for alle h, k

$$\text{Under } H_0: E[Y_i] = \mu_{..} \pi_{h..} \pi_{..k}$$

$$\eta_{hk} = \ln \mu_{..} + \ln \pi_{h..} + \ln \pi_{..k} \quad (*)$$

$$= \alpha + \gamma_h^A + \gamma_k^B \quad \text{med}$$

$$\begin{aligned} \sum_h \pi_{h..} &= 1 \\ \sum_k \pi_{..k} &= 1 \end{aligned} \quad \left. \right\} \text{lineare bånd}$$

$$\text{antal parametere: } 1 + (r-1) + (c-1) \\ = r+c-1$$

Alternativ parametrisering:

$$\alpha = \ln \mu_{..} + \frac{1}{r} \sum_h \ln \pi_{h..} + \frac{1}{c} \sum_k \ln \pi_{..k}$$

$$\gamma_h^A = \ln \pi_{h..} - \frac{1}{r} \sum_h \ln \pi_{h..}$$

$$\gamma_k^B = \ln \pi_{..k} - \frac{1}{c} \sum_k \ln \pi_{..k}$$

$$\text{med } \sum_h \gamma_h^A = 0$$

$$\sum_k \gamma_k^B = 0$$

Referencemodel (den mattede model)

$$\text{Eneste bånd: } \sum_n \sum_k \pi_{nk} = 1$$

Den til (*) analoge log-lineare model

$$\eta_{hk} = \ln \mu_{..} + \ln \pi_{h.} + \ln \pi_{.k} + \ln \frac{\pi_{hk}}{\pi_{h.} \pi_{.k}}$$

$$= \lambda + \lambda_h^A + \lambda_k^B + \lambda_{hk}^{AB} \quad \text{med}$$

$$\text{yderligere bånd } \sum_n \lambda_{hk}^{AB} = 0$$

(fx)

$$\sum_k \lambda_{hk}^{AB} = 0$$

antal parametre:

$$1 + (r-1) + (c-1) + (r-1)(c-1) = rc$$

Betræk ogå $Y_i | N \sim m(N; \pi_1, \dots, \pi_{rc})$

$$E[Y_i | N] = N \pi_i = N \pi_{hi}$$

$$= N \pi_{h.} \pi_{.k} \text{ under } H_0$$

Alle ovenstående regninger i forb. med Poissonmodellen kan oprettholdes, men blot $\mu_{..}$ skiftes ud med N . Betragt dog, at der er få parametere mindre i multinomialmodellen.

Likelihood i log-lineare modeller

ℓ_p loglikelihood i Poissonmodel

ℓ_m loglikelihood i tilsvarende multinomialmodel

$$\begin{aligned}
 l_p &= \sum_i (y_i \ln \mu_i - \mu_i) + c, \text{ jf. AA s. 227} \\
 l_m &= \sum_i y_i \ln \pi_i + c, \text{ jf. AA s. 136} \\
 &= \sum_i y_i (\ln \mu_i - \ln \mu_{..}) + c, \text{ idet } \pi_i = \frac{\mu_i}{\mu_{..}} \\
 &= \sum_i (y_i \ln \mu_i - \mu_i - y_i \ln \mu_{..} + \mu_{..}) + c \\
 &= \sum_i y_i \ln \mu_i - \mu_{..} - (N \ln \mu_{..} - \mu_{..}) + c \\
 &= l_p - l_p(\mu_{..}) + c
 \end{aligned}$$

Når μ_i indeholder et konstant led, så
 gælder $\hat{\mu}_{..} = \sum_i \hat{\mu}_i = \sum_i y_i = N$, jf. AA s. 236, dvs.
 $\hat{\mu}_i$ vestent i Poissonmodellen har også
 gyldighed i multinomialmodellen.

Udbygrende argumentation

$$y_i = \ln \mu_i = \alpha + z_i^T \beta \quad (\text{konstantledet er skilt ud})$$

$$\begin{aligned}
 l_p &= \sum_i (y_i \ln \mu_i - \mu_i) \\
 &= \sum_i y_i (\alpha + z_i^T \beta) - \mu_{..} \\
 &= N \alpha + \sum_i y_i (z_i^T \beta) - \mu_{..} \\
 &= \sum_i y_i (z_i^T \beta) - N \ln \sum_i \exp z_i^T \beta \\
 &\quad + N \alpha + N \ln \sum_i \exp z_i^T \beta - \mu_{..}
 \end{aligned}$$

$$\begin{aligned}
 &\alpha + \ln \sum_i \exp z_i^T \beta \\
 &= \ln \exp (\alpha + \sum_i z_i^T \beta) \\
 &= \ln \mu_{..}
 \end{aligned}$$

$$\begin{aligned}
 \frac{\exp z_i^T \beta}{\sum_j \exp z_j^T \beta} &= \frac{\exp (\alpha + z_i^T \beta)}{\exp (\alpha + \sum_j z_j^T \beta)} \\
 &= \frac{\mu_i}{\mu_{..}} = \pi_i
 \end{aligned}$$

$$\begin{aligned}
 &= \sum_i y_i \ln (\exp z_i^T \beta) - \sum_i y_i \ln \sum_j \exp z_j^T \beta \\
 &\quad + N (\alpha + \ln \sum_i \exp z_i^T \beta) - \mu_{..} \\
 &= \sum_i y_i \ln \frac{\exp z_i^T \beta}{\sum_j \exp z_j^T \beta} + N \ln \mu_{..} - \mu_{..}
 \end{aligned}$$

$$= \sum_i y_i \ln \bar{\pi}_i + N \ln p_{..} - \mu_{..}$$

dvs.

$$\ell_p(\mu_{..}, p) = \ell_m(p) + \ell_p(\mu_{..})$$

ugift. af p

$\ell_p(\mu_{..}, p)$ og $\ell_m(p)$ maksimeres for samme \hat{p} .

Hvis rækkesummerne er faste, eller hvis spissummerne er faste, kan tilsvarende ekvivalens mellem Poissonfordelingers likelihood og de udgående multinomialfordelingers likelihood eftervises.

Se tilføjelse
side 10

Odds

En haandelse med sands. p har

$$\text{odds } \frac{p}{1-p}$$

Bemerk, at ln odds = logit

Betrægt schema

A 2×2 probability table

Exposure to risk	Occurrence of event		Total
	\bar{E}	E	
\bar{X}	π_{00}	π_{01}	π_{0+}
X	π_{10}	π_{11}	π_{1+}
Total	π_{+0}	π_{+1}	1

Prospektiv studie

X og \tilde{X} valgt på forhånd

ford. på E og \tilde{E} observeres

$$\text{logit } P(E|\tilde{X}) - \text{logit } P(E|X)$$

$$= \ln \frac{\frac{\pi_{01}}{\pi_{00}}}{\frac{\pi_{10}}{\pi_{00}}} - \ln \frac{\frac{\pi_{11}}{\pi_{10}}}{\frac{\pi_{01}}{\pi_{10}}} = \ln \frac{\pi_{01} \pi_{10}}{\pi_{00} \pi_{11}}$$

$$= -\ln w$$

Retrospektiv studie

Der foretages obs. af E og \tilde{E}

der klassificeres efter X og \tilde{X}

$$\text{Logit } P(X|\tilde{E}) - \text{logit } P(X|E)$$

$$= \ln \frac{\frac{\pi_{10}}{\pi_{11}}}{\frac{\pi_{00}}{\pi_{10}}} - \ln \frac{\frac{\pi_{11}}{\pi_{01}}}{\frac{\pi_{00}}{\pi_{01}}} = \ln \frac{\pi_{01} \pi_{10}}{\pi_{11} \pi_{00}}$$

$$= -\ln w$$

'vægtsproduktforholdet' er altså ens

$w = 1$ svarer til uafhængighed

$w > 1$ - - - positivt sammenspiel
af faktorer

$w < 1$ - - - negativt sammenspiel
af faktorer

$$\begin{aligned}
 -\ln w &= \ln \left(\frac{\pi_{01}}{\pi_{00} + \pi_{01}} \cdot \frac{\pi_{00} + \pi_{10}}{\pi_{00}} \cdot \frac{\pi_{10}}{\pi_{10} + \pi_{11}} \cdot \frac{\pi_{10} + \pi_{11}}{\pi_{11}} \right) \\
 &= \lambda_{01}^{xE} - \lambda_{00}^{xE} + \lambda_{10}^{xE} - \lambda_{11}^{xE}
 \end{aligned}$$

Värd mellan parametrarne

$$\lambda_{00}^{xE} + \lambda_{10}^{xE} = 0$$

$$\lambda_{01}^{xE} + \lambda_{11}^{xE} = 0$$

$$\lambda_{00}^{xE} + \lambda_{01}^{xE} = 0$$

$$\lambda_{10}^{xE} + \lambda_{11}^{xE} = 0$$

Heraf

$$\lambda_{10}^{xE} = -\lambda_{00}^{xE}$$

$$\lambda_{11}^{xE} = -\lambda_{01}^{xE}$$

$$\lambda_{01}^{xE} = -\lambda_{00}^{xE}$$

$$\lambda_{11}^{xE} = -(-\lambda_{00}^{xE})$$

$$\ln w = \lambda_{00}^{xE} + \lambda_{00}^{xE} + \lambda_{00}^{xE} + \lambda_{00}^{xE} = 4\lambda_{00}^{xE}$$

Kvarn likelihood

$$\text{Antag } E[Y_i] = \mu_i$$

$$\text{Var}[Y_i] = \alpha V(\mu_i)$$

$$\text{Sät } u = u(Y_i, \mu, \alpha) = \frac{Y - \mu}{\sqrt{V(\mu)}} \quad (\text{index } i \text{ undelat})$$

$$\text{Heraf } E[u] = 0$$

$$\text{Var}[u] = \frac{1}{(\alpha V(\mu))^2} \text{Var} Y = \frac{1}{\alpha^2 V(\mu)}$$

$$\frac{\partial u}{\partial \mu} = \frac{\alpha V(\mu)(-1) - (Y - \mu)\alpha V'(\mu)}{(\alpha V(\mu))^2}$$

$$\Rightarrow -E\left[\frac{\partial u}{\partial \mu}\right] = \frac{1}{\alpha V(\mu)}$$

des. u upphovs sig analogt
til en scoringspunkt

Kvazi likelihood defineres derfor som

$$\begin{aligned} Q(\mu; y) &= \int_y^{\mu} u(t, y) dt \\ &= \int_y^{\mu} \frac{y-t}{\sqrt{V(t)}} dt \end{aligned}$$

Når flere uafhængige observationer, så

$$Q(\mu; y) = \sum_i Q(\mu_i; y_i)$$

Når Q er en egentlig loglikelihoodfkt.,
dvs. når der findes en funktion ℓ , hvor

$$\frac{\partial \ell}{\partial \mu} = \frac{y-\mu}{\sqrt{V(\mu)}}, \quad E[Y] = \mu, \quad \text{Var}[Y] = \sqrt{V(\mu)},$$

så foreligger der en eksponentiel familie
(afvis).

I en regressionsmodel er μ_i en fkt. af β .

Kvazi likelihoodsligningerne

$$\frac{\partial Q}{\partial \mu_i} \frac{\partial \mu_i}{\partial \beta_j} = 0, \quad j=1, \dots, p$$

ses at være

$$\sum_i \frac{y_i - \mu_i}{\sqrt{V(\mu_i)}} \frac{\partial \mu_i}{\partial \beta_j} = 0, \quad j=1, \dots, p$$

og afhænger ikke af dispersionsparametren σ^2 .

Kvazi devians

$$\begin{aligned} D(y; \hat{\mu}) &= -2 \nabla Q(\hat{\mu}; y) \\ &= 2 \int_{\hat{\mu}}^y \frac{y-t}{\sqrt{V(t)}} dt \geq 0 \end{aligned}$$

Når uafhængige observations

$$D(y; \hat{\mu}) = \sum_i D(y_i; \hat{\mu}_i)$$

Eksempel. Overdispersion i Poissonmodellen

vedr stofdata, lf. AA s. 225-226

$$D_i \sim n(p_i; L_i)$$

$$\text{GLM: } E[D_i] = \mu_i L_i$$

$$\text{Var}[D_i] = \gamma E[D_i] = \gamma \mu_i L_i$$



Faktoren γ modelles
for overdispersion

Estimation af γ (jf. AA s. 240):

$$\begin{aligned}\hat{\gamma} &= \frac{1}{n-p} \sum_i \frac{w_i (y_i - \hat{\mu}_i)^2}{V(\hat{\mu}_i)} \\ &= \frac{1}{n-1} \sum_i \frac{(D_i - \hat{\beta} L_i)^2}{\hat{\beta} L_i}\end{aligned}$$

Fortsat fra side 6:

I det maksimum likelihood estimation er man i hvert den ene Poissonmodel, i multinomialmodellen og i tilfældet med uafhængige multinomialmodeller, kan vi, hvadenten vi ønsker at teste for uafhængighed mellem to faktorer, eller vi ønsker at teste for homogenitet mellem en faktors niveauer, i alle tilfælde benytte Pearson χ^2 som teststatistikk, jf.

AA s. 137 og s. 140 samt note 7 s. 10 og s. 11.

Som af: AA
s. 224-225
s. 226-227
lf. 4-5
lf. 6-7
lf. 8-9