

Dagens program:

Mandag den 14. februar

Den generelle lineære model

- Kort om likelihoodfunktion
- Flerdimensional normalfordeling
- Generel lineær model
- Estimation, fordeling af residualer mv
- Test for modelreduktion
- Successiv testning

Eksempel, hjerneblodgennemstrømning hos malere

Poul Thyregod, introslide.tex
Specialkursus vid.stat. forår 2005

DTU
IMM

Likelihoodfunktion, begreber :

- **score-funktion:** første afledede af log-likelihood har middelværdien nul
- **observeret information:** anden afledede af log-likelihood (krumning)
- **forventet information:** middelværdien af den observerede information
Og det er netop dispersionsmatricen for scorefunktionen

Poul Thyregod, introslide.tex
Specialkursus vid.stat. forår 2005

DTU
IMM

Likelihoodfunktion :

Kommer fra udtrykket for sandsynlighedstætheden

Sandsynlighedstætheden beskriver en **hypotese** af observationer - for en forelagt værdi af parameteren

likelihoodfunktionen er en funktion af parameteren - for forelagte værdier af observationerne

likelihoodfunktionen beskriver **parameterusikkerhed**, nemlig grad af overensstemmelse mellem data og parameter.

Poul Thyregod, introslide.tex
Specialkursus vid.stat. forår 2005

DTU
IMM

Likelihoodfunktion, transformation af parametre :

Lad $\theta = \theta(\beta) \in \mathbb{R}^k$ for $\beta \in B \subset \mathbb{R}^m$ med $m \leq k$

Da er scorefunktionen (mht. β)

$$l'_\beta(\beta; \mathbf{y}) = \mathbf{J}^T l'_\theta(\theta(\beta); \mathbf{y}) \quad (2.11)$$

hvor

$$\mathbf{J} = \frac{\partial \theta}{\partial \beta}$$

med elementer

$$\mathbf{J}_{ij} = \frac{\partial \theta_i}{\partial \beta_j} \quad i = 1, 2, \dots, k; \quad j = 1, 2, \dots, m$$

og den forventede information er

$$\mathbf{i}_\beta(\beta) = -\mathbb{E} \left[\frac{\partial^2}{\partial \beta \partial \beta^T} \ell(\theta(\beta); \mathbf{Y}) \right] = \mathbf{J}^T \mathbf{i}_\theta(\theta(\beta)) \mathbf{J} \quad (2.12)$$

Poul Thyregod, introslide.tex
Specialkursus vid.stat. forår 2005

DTU
IMM

Normalfordeling :

$$f(\mathbf{y}; \boldsymbol{\mu}, \sigma^2) = \frac{1}{(\sqrt{2\pi})^k \sigma^k \det(\Sigma)} \exp \left[\frac{-1}{2\sigma^2} (\mathbf{y} - \boldsymbol{\mu})' \Sigma^{-1} (\mathbf{y} - \boldsymbol{\mu}) \right] \quad \text{for } \mathbf{y} \in \mathbb{R}^k \quad (3.4)$$

$$\delta_\Sigma(\mathbf{y}_1, \mathbf{y}_2) = \mathbf{y}_1' \Sigma^{-1} \mathbf{y}_2 \quad (3.5)$$

et indre produkt på \mathbb{R}^k .

Tilsvarende **ortogonalitetsbegreb** (svarende til at det indre produkt er nul), og en **norm** (defineret ved $\|\mathbf{y}\|_\Sigma = \sqrt{\delta(\mathbf{y}, \mathbf{y})}$).

Med notationen

$$D(\mathbf{y}; \boldsymbol{\mu}) = \|\mathbf{y} - \boldsymbol{\mu}\|_\Sigma^2 = (\mathbf{y} - \boldsymbol{\mu})' \Sigma^{-1} (\mathbf{y} - \boldsymbol{\mu}) \quad (3.6)$$

bliver tæthedsfunktionen

$$f(\mathbf{y}; \boldsymbol{\mu}, \sigma^2) = \frac{1}{(\sqrt{2\pi})^k \sigma^k \det(\Sigma)} \exp \left[\frac{-1}{2\sigma^2} D(\mathbf{y}; \boldsymbol{\mu}) \right] \quad (3.7)$$



Generel lineær model:

Model (koordinatfri):

$$H_0 : \boldsymbol{\mu} \in L,$$

et lineært (affint) m -dimensionalt underrum af \mathbb{R}^k

(egentlig $\boldsymbol{\mu} - \boldsymbol{\mu}_0 \in L$)

Parametriseret (i koordinater):

$$H_0 : \boldsymbol{\mu}(\boldsymbol{\beta}) = \mathbf{X}\boldsymbol{\beta}$$

Søjlerne i \mathbf{X} udspænder underrummet L (et lokalt koordinatsystem).

$\boldsymbol{\beta}$ angiver punkternes lokale koordinater.

Utallige ækvivalente parametriske fremstillinger af en given model.



Normalfordeling, scorefunktion og information :

log-likelihood

$$\ell_\mu(\boldsymbol{\mu}, \sigma^2; \mathbf{y}) = -(k/2) \log(\sigma^2) - \frac{1}{2\sigma^2} (\mathbf{y} - \boldsymbol{\mu})' \Sigma^{-1} (\mathbf{y} - \boldsymbol{\mu}) \quad (3.8)$$

Scorefunktion (m.h.t. $\boldsymbol{\mu}$)

$$\frac{\partial}{\partial \boldsymbol{\mu}} \ell_\mu(\boldsymbol{\mu}, \sigma^2; \mathbf{y}) = \frac{1}{\sigma^2} [\Sigma^{-1} \mathbf{y} - \Sigma^{-1} \boldsymbol{\mu}] = \frac{1}{\sigma^2} \Sigma^{-1} (\mathbf{y} - \boldsymbol{\mu}) \quad (3.9)$$

Observeret information (m.h.t. $\boldsymbol{\mu}$)

$$\mathbf{j}(\boldsymbol{\mu}; \mathbf{y}) = \frac{1}{\sigma^2} \Sigma^{-1}. \quad (3.10)$$

afhænger ikke af observationerne y , hvorfor

den **forventede information** er

$$\mathbf{i}(\boldsymbol{\mu}) = \frac{1}{\sigma^2} \Sigma^{-1}.$$



Eksempel, sædvanlig endimensional regression:

Ramushøjder, afsnit 3.5

Alder (år)	Højde (mm)
x_i	y_i
8.0	51.2
8.5	53.0
9.0	54.3
9.5	54.5



Eksempel, sædvanlig endimensional regression:

Model:

$$E[Y] = \beta_0 + \beta_1 x_i$$

$$\mu_1 = \beta_0 + \beta_1 x_1$$

$$\mu_2 = \beta_0 + \beta_1 x_2$$

$$\mu_3 = \beta_0 + \beta_1 x_3$$

$$\mu_4 = \beta_0 + \beta_1 x_4$$

dvs

$$\mu = X\beta$$



Eksempel, sædvanlig endimensional regression:

med

$$X = \begin{pmatrix} 1 & 8.0 \\ 1 & 8.5 \\ 1 & 9.0 \\ 1 & 9.5 \end{pmatrix}$$



Generel lineær model, estimation under model:

Koordinatfrit:

maksimum-likelihood estimatet $\hat{\mu}$ for sættet af middelværdier fås ved at minimere $D(y; \mu)$, dvs som den ortogonale (m.h.t δ_Σ) projktion, $p_L(y)$ af punktet y ned på det lineære underrum, L .

Parametrisk:

Vi bruger (2.11) og (3.9) og får scorefunktionen

$$\frac{\partial}{\partial \beta} \ell_\beta(\beta, \sigma^2; y) = \left[\frac{\partial \mu}{\partial \beta} \right]' \frac{\partial}{\partial \mu} \ell_\mu(\mu(\beta), \sigma^2; y) = \frac{1}{\sigma^2} X' [\Sigma^{-1} y - \Sigma^{-1} X \beta] \quad (3.11)$$

hvorfor maksimum-likelihood estimatet for β fås som løsning til **normalligningen**

$$X' \Sigma^{-1} y = X' \Sigma^{-1} X \hat{\beta} \quad (3.12)$$

Altså m ligninger med m ubekendte.



Generel lineær model, estimation under model:

Hvis X har fuld rang (m), er der en entydig løsning, som bestemmes ved

$$\hat{\beta} = (X' \Sigma^{-1} X)^{-1} X' \Sigma^{-1} y \quad (3.16)$$

Hatmatricen

$$H = X [X' \Sigma^{-1} X]^{-1} X'$$

er projektionsmatricen, der overfører y til de fitede værdier $\hat{\mu}$ (nemlig sætter hat på μ 'erne).

idet nemlig

$$\hat{\mu} = p_L(y) = X \hat{\beta} = Hy$$

Residualerne er

$$r = y - X \hat{\beta} = (I - H)y \quad (3.17)$$



Generel lineær model, residualer og fittede værdier:

Opspaltning:

$$y = X\hat{\beta} + r = Hy + (I - H)y$$

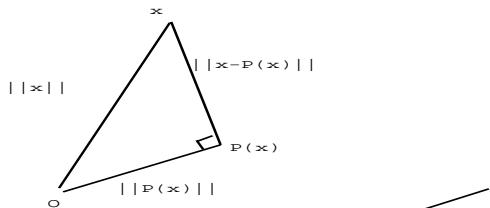
Normalligningen

$$X'\Sigma^{-1}(y - X\hat{\beta}) = 0 \quad (3.19)$$

viser, at residualerne $r = (y - X\hat{\beta})$ er **ortogonale** (mht Σ^{-1}) på underrummet, L , udspændet af X

Generel lineær model, opspaltning :

$$\begin{aligned} D(y; 0) &= D(y; X\hat{\beta}) + D(X\hat{\beta}; 0) \\ &= (y - X\hat{\beta})'\Sigma^{-1}(y - X\hat{\beta}) + (X\hat{\beta})'\Sigma^{-1}(X\hat{\beta}) \end{aligned} \quad (3.20)$$



Generel lineær model, fordeling af estimatorer:

$$E[\hat{\beta}] = \beta \quad (3.21)$$

$$\mathbf{D}[\hat{\beta}] = \sigma^2 (X'\Sigma^{-1}X)^{-1} \quad (3.22)$$

Bruges ved test af enkelte parameterværdier

$$\mathbf{D}[X\hat{\beta}] = \sigma^2 X [X'\Sigma^{-1}X]^{-1} X' = \sigma^2 H$$

$$\mathbf{D}[r] = \sigma^2 (I - H)$$

rang (frihedsgrader) $k - m$.

Bruges til at **standardisere** residualer

Generel lineær model, estimation af σ^2 :

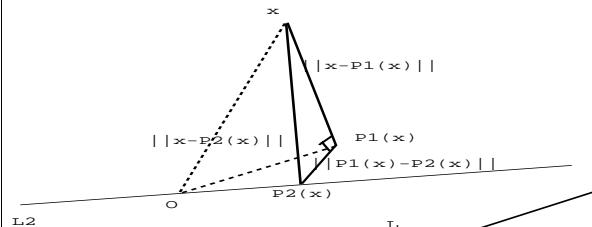
σ^2 estimeres ved

$$s^2 = r'\Sigma^{-1}r / (k - m) \quad (3.24)$$

som er uafhængig af $\hat{\mu}$

Generel lineær model, test af delhypotese:

Deltypotese: $\mu \in L_2 \subset L$



Poul Thyregod, introslide.tex
Specialkursus vid.stat. forår 2005

17

IMM

DTU

DTU

Partielt test, type III opspaltnings:

Antag at H_i gælder.

Vi vil nu undersøge $H_{i+1} \subset H_i$

Kvotienttestet for H_{i+1} under antagelse af at H_i kan opretholdes har teststørrelsen

$$F(\mathbf{y}) = \frac{D(\hat{\mu}, \hat{\mu}) / (m_1 - m_2)}{s_1^2}, \quad (3.33)$$

hvor $\hat{\mu}$ angiver de fittede værdier under H_i og $\hat{\mu}$ angiver de fittede værdier under H_{i+1} og s_1^2 angiver residualvariansen under H_i .

Testet kaldes det partielle kvotienttest - eller et Type III test.

Testet "korrigerer" for alle de effekter, der er i modellen på det pågældende trin.

Poul Thyregod, introslide.tex
Specialkursus vid.stat. forår 2005

19

IMM

DTU

Generel lineær model, kvotienttest af delhypotese:

$$\begin{aligned} q_{2|1}(\mathbf{y}) &= \left[\frac{D(\mathbf{y}; p_2(\mathbf{y}))}{D(\mathbf{y}; p_1(\mathbf{y})) + D(p_1(\mathbf{y}); p_2(\mathbf{y}))} \right]^{k/2} \\ &= \left[\frac{1}{1 + \frac{D(p_1(\mathbf{y}); p_2(\mathbf{y}))}{D(\mathbf{y}; p_1(\mathbf{y}))}} \right]^{k/2} = \left[\frac{1}{1 + \frac{m_1 - m_2}{k - m_1} F} \right]^{k/2} \end{aligned}$$

Kvotienttestet fører således til F-teststørrelsen

$$F = \frac{D(p_1(\mathbf{y}); p_2(\mathbf{y})) / (m_1 - m_2)}{D(\mathbf{y}; p_1(\mathbf{y})) / (k - m_1)}$$

der følger en $F(m_1 - m_2, k - m_1)$ -fordeling (eksakt resultat)

Poul Thyregod, introslide.tex
Specialkursus vid.stat. forår 2005

18

IMM

DTU

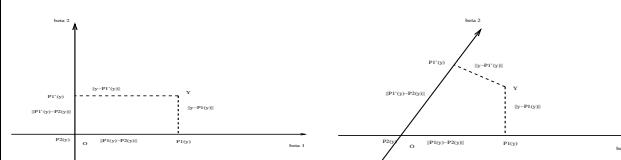
IMM

Successiv testning i hypotesekæde:

$$\mathbb{R} \subseteq L_M \dots \subseteq L_2 \subseteq L_1 \subseteq L_0 \subseteq \mathbb{R}^k, \quad (3.32)$$

Type I opspaltnings udtrykker successive projktioner ned gennem kæden af underrum.

Størrelsen af leddene i Type I opspaltningen kan afhænge af rækkefølgen i modelformlen



Poul Thyregod, introslide.tex
Specialkursus vid.stat. forår 2005

20

Successiv testning i hypotesekæde:

Type III opspaltnings svarer til at man fjerner et led - og vurderer reduktionen i model SSQ ved at fjerne det

så smider man det ind igen, og fjerner et af de andre.

Type III opspaltningen afhænger ikke af rækkefølgen i modelformlen.

I praksis tager man udgangspunkt i en type III opspaltning og fjerner mindst betydende - og mest komplikerede led.

derefter reestimeres - og reduceres

Indskud om multikollinearitet:

Det kan ske, at to eller flere af de forklarende variable følges ad (er næsten lineært afhængige). Konsekvenser af dette er, at modellen er signifikant, men estimatorne for de enkelte koeficienter har stor varians ($X'X$ er næsten singulær). Desuden er fortolkningen af de enkelte koeficienter højst tvivlsom.

Diagnostikker for multikollinearitet udskrives sammen med estimatorne:

$$TOL_i = 1 - R_i^2$$

hvor R_i^2 er den multiple determinationskoeficient fra regressionen af den i 'te forklarende variable på de øvrige. Lille TOL er tegn på multikollinearitet.

$$V[\hat{\beta}_i] \propto 1/TOL_i$$

$$VIF_i = 1/TOL_i$$

Hvis $VIF_i > 10$ grund til bekymring

Diagnostikker:

Oversigt over diagnostiske størrelser samt parametriseringer i notens afsnit 3.9.1 til 3.9.4

Eksempel, Blodgennemstrømningsindeks hos malere:

Eksempel på

- Brug af forklarende regressionsvariabel til korrektion for forskelle i populationer (Kovariansanalyse)
- Brug af modelformler
- Forskel på typel og type III opspaltning

Eksempel, Blodgennemstrømningsindeks hos malere:

Først tegnes boxplot for ISI for de to grupper
Ikke stærke tegn på forskel

Så proves t-test for to uafhængige stikprøver
ved brug af FIT-menuen

Tegn ISI mod alder, marker de to grupper hver for sig
Indikation af sammenhæng mellem ISI og alder

Korriger for alderseffekten ved stor model:
 $ISI = \text{ORG ALDER} + \text{ORG} * \text{ALDER}$
To rette linier
Test for reduktion af model til to parallelle linier:
 $ISI = \text{ORG ALDER}$
Nu er effekten af brugen af oplosningsmidler signifikant.

IMM