



Dagens program:

Mandag den 25. april

Hierarkiske normalfordelingsmodeller

- Resumé af ensidet variansanalysemodel med tilfældig effekt estimation af tilfældige effekter, fortolkning som BLUP estimation, empirisk Bayes-estimation, Steins sætning
- Tilfældige regressionslinier
- Hierarkiske flerdimensionale observationer
- Generelle lineære mixed models
- SAS PROC Mixed
- Begyndelse af hierarkiske modeller for eksponentielle familier

Modellens parametre



Modellen har parametrene $\mu_0, \sigma^2, \sigma_b^2$.

Vi sætter

$$\gamma = \sigma_b^2 / \sigma^2$$

(Signal støj- forhold)

Hierarkisk model, normalfordeling



Først udtages en balle tilfældigt. Derefter udtages fire prøver tilfældigt fra denne balle.

Y_{ij} j'te prøve fra i'te balle.

For fastholdt balle, uafhængige gentagelser.

$$Y_{ij} | \mu_i \sim N(\mu_i, \sigma^2), \\ Y_i | \mu_i \sim N(\mu_i, \sigma^2 I_4)$$

Ballerenheden μ_i varierer fra balle til balle:

$$\mu_i \sim N(\mu_0, \sigma_b^2)$$

indbyrdes uafhængige

Fordeling af i'te sæt



Lad \mathbf{Y}_i angive observationer fra i'te udtagne balle.

$\mathbf{Y}_1, \mathbf{Y}_2, \dots, \mathbf{Y}_q$ er uafhængige observationssæt

$$\mathbf{Y}_i \sim N_{n_i}(\boldsymbol{\mu}_0, \sigma^2 \mathbf{I}_{n_i} + \sigma_b^2 \mathbf{J}_{n_i})$$

hvor \mathbf{J}_n angiver en $n \times n$ matrix med etaller.

Dispersionsmatricen for \mathbf{Y}_i er

$$\mathbf{V}_i = \mathbf{D}[\mathbf{Y}_i] = E[(\mathbf{Y}_i - \boldsymbol{\mu}_0)(\mathbf{Y}_i - \boldsymbol{\mu}_0)^T]$$

$$= \begin{pmatrix} \sigma_b^2 + \sigma^2 & \sigma_b^2 & \dots & \sigma_b^2 \\ \sigma_b^2 & \sigma_b^2 + \sigma^2 & \dots & \sigma_b^2 \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_b^2 & \sigma_b^2 & \dots & \sigma_b^2 + \sigma^2 \end{pmatrix}$$

Fordeling af hele observationssættet



\mathbf{Y} består af q uafhængige gentagelser med fordeling som \mathbf{Y}_i

$$\mathbf{D}[\mathbf{Y}] = \mathbf{V} = \text{Blok diag}\{\mathbf{V}_i\}$$

Estimation (prædiktion) af enkelte ballemiddelværdier

Det hierarkiske maximum likelihood estimat er BLUP-estimat (best linear unbiased predictor), et vejet gennemsnit mellem stikprøveresultat og fælles middelværdi

Shrinkage (krympning) mod fælles gennemsnit.

Empirisk Bayes fortolkning (estimeret apriorifordeling).



Estimation (prædiktion) af enkelte ballemiddelværdier

Løsningen tilfredsstiller

$$\hat{\mu}_i = (1 - w_i(\hat{\gamma}))\bar{y}_i + w_i(\hat{\gamma})\hat{\mu}_0$$

et vejet gennemsnit mellem de individuelle gennemsnit \bar{y}_i og det fælles gennemsnit $\hat{\mu}_0$ med vægtene $(1 - w_i(\hat{\gamma}))$ og $w_i(\hat{\gamma})$, hvor

$$1 - w = \frac{V[E[\hat{\theta}|\theta]]}{E[V[\hat{\theta}|\theta]] + V[E[\hat{\theta}|\theta]]}$$

$$w = \frac{E[V[\hat{\theta}|\theta]]}{E[V[\hat{\theta}|\theta]] + V[E[\hat{\theta}|\theta]]}$$

netop er de to bidrag (variationen indenfor balle, og variationen mellem baller) fra opspaltningen af den totale varians

$$V[\hat{\theta}] = E[V[\hat{\theta}|\theta]] + V[E[\hat{\theta}|\theta]]$$

Empirisk Bayes estimat har mindre estimationsfejl end ML-estimat i systematisk model

Stein's sætning vedrørende simultan estimation:
(Bemærkning 1, side 246.)

Ved en **systematisk** (fixed effects) model for de enkelte balleniveauer estimeres hver balle ved sit gennemsnit, \bar{y}_i , med varians σ^2/n .

Betrægt den **gennemsnitlige estimationsfejl**

$$R(\boldsymbol{\mu}, d(\cdot)) = \frac{1}{k} E \left[\sum (d_i(\mathbf{y}) - \mu_i)^2 \right]$$

hvor $d_i(\mathbf{y})$ angiver estimatoren $\hat{\mu}_i$ for det i 'te balleniveau.

Stein's sætning

Betrægt opspaltingen af den gennemsnitlige estimationsfejl i varians og bias:

$$\begin{aligned} R(\boldsymbol{\mu}, d(\cdot)) &= \frac{1}{k} E \left[\sum (d_i(\mathbf{y}) - \mu_i)^2 \right] \\ &= \frac{1}{k} E \left[\sum \left\{ (d_i(\mathbf{y}) - E[d_i])^2 + (E[d_i] - \mu_i)^2 \right\} \right] = \frac{1}{k} \sum \left(V[d_i] + \delta_i^2 \right) \end{aligned}$$

hvor $\delta_i = E[d_i] - \mu_i$ angiver **bias** for estimatoren $d_i(\mathbf{y})$

Såfremt man vælger maximum-likelihood estimatoren $d_i^{ML}(\mathbf{y}) = \bar{y}_i$, finder man

$$R(\boldsymbol{\mu}, d^{ML}(\cdot)) = \frac{1}{k} \sum V[\bar{y}_i] = \sigma^2/n$$

idet estimatoren er **central** (ikke har nogen bias).

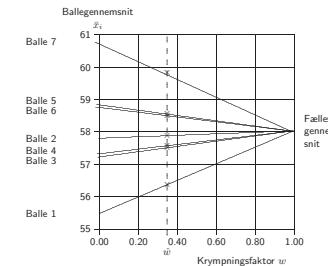
Man kan vise (Stein 1955), at såfremt man benytter empirisk Bayes estimatoren

$$d_i^{EB}(\mathbf{y}) = w(\mathbf{y})\bar{y}_i + (1 - w(\mathbf{y}))\bar{y}_i$$

gælder for et vilkårligt parametersæt $\boldsymbol{\mu}$, at

$$R(\boldsymbol{\mu}, d^{EB}(\cdot)) < R(\boldsymbol{\mu}, d^{ML}(\cdot))$$

Empirisk Bayes estimat som shrinkage estimat (vejet gennemsnit mellem fælles middelværdi og individuelle gennemsnit)



Intuitivt

For $w < 1$ vil $E[w\bar{y}_i] = w\mu_i < \mu_i$ og $V[w\bar{y}_i] = w^2V[\bar{y}_i] < \sigma^2/n$

Ved brug af empirisk Bayes-estimatet får vi altså reduceret variansen på bekostning af en lille bias $\delta_i = (1 - w)\mu_i$.

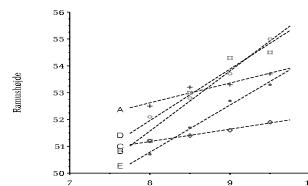
$$R(\boldsymbol{\mu}, d^{EB}(\cdot)) \approx \frac{1}{k} \sum [w^2\sigma^2/n + (1 - w^2)\mu_i^2] < \frac{1}{k} \sum \sigma^2/n$$

Andre eksempler på hierarkiske modeller, og bestemmelse af aposteriorifordelinger

1. Tilfældige regressionslinier
2. Todimensionale målinger

Varierende regressionslinier

Ramushøjder for 5 drenge ved fire forskellige aldre.



Systematiske modeller: Vi kunne lave én fælles regressionslinie for alle 5 drenge, eller vi kunne estimere 5 individuelle regressionslinier (en for hver dreng),

Aposteriorifordeling af regressionskoefficienter

Y en $n \times 1$ dimensional vektor af observationer, og X en $n \times p$ dimensional matrix af kendte koefficienter.

Lad $Y | \beta \sim N_n(X\beta, \sigma^2 V)$ og lad apriorifordelingen af β være

$$\beta \sim N_p(\beta_0, \sigma^2 \Lambda)$$

Aposteriorifordelingen af β efter observation af $Y = y$ er

$$\beta | Y = y \sim N_p(\beta_1, \sigma^2 \Gamma_1)$$

hvor

$$\beta_1 = W\beta_0 + (I - W)\hat{\beta}$$

med $W = (I + \Gamma)^{-1}$ og $\hat{\beta}$ er den sædvanlige mindste kvadraters estimator



Varierende regressionslinier, tilfældig (hierarkisk) model

Data opfattes som en **stikprøve** af drenge

$$Y_i = X\beta_i + \epsilon_i, \quad i = 1, 2, \dots, k$$

hvor

$$\beta_i \sim N_2(\beta_0, \sigma^2 \Gamma), \quad \epsilon_i \sim N_n(0, \sigma^2 I_n),$$

β_i, β_j er indbyrdes uafhængige for $i \neq j$,

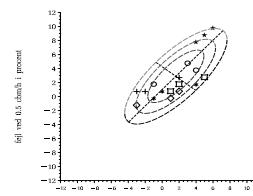
Kovariansmatricen Γ angiver kovariansmatricen i populationsfordelingen af interceper og hældninger (med målestøjen σ^2 trukket ud som en faktor).

Flerdimensional normalfordeling

3 gentagne kalibreringer af 6 flowmålere, der er udtaget af en større målerpopulation. De 6 målere blev hver kalibreret ved de samme to flow, henholdsvis $0.1 [m^3/h]$ og $0.5 [m^3/h]$.

Sædvanlige overflader af registreret fejl ved et flow for de 3 gentagne prøver på hver af de 6 flowmålere

gentagelser på samme måler er markeret med samme symbol





Flerdimensional normalfordeling, hierarkisk model

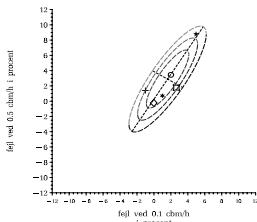
$$\mathbf{X}_{ij} = \boldsymbol{\mu}_0 + \boldsymbol{\alpha}_i + \boldsymbol{\epsilon}_{ij}, \quad i = 41, 42, \dots, 46; \quad j = 1, 2, 3.$$

hvor $\boldsymbol{\alpha}_i$ er uafhængige, $\boldsymbol{\alpha}_i \sim N_2(0, \Sigma_0)$



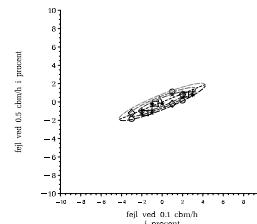
Flerdimensional normalfordeling, variation mellem målere

Sammeværende værdier af estimerede målerøjls ved to flow
for stikprøve bestående af 8 flommålere



Flerdimensional normalfordeling, variation indenfor målere

Sammeværende værdier af kalibreringsrøjls ved to flow
for 3 genetiske prøver på hver af 8 flommålere
måderens spalt(er) er eliminert



Flerdimensional normalfordeling, aposteriorifordeling

Lad $\mathbf{X} | \boldsymbol{\mu} \sim N_p(\boldsymbol{\mu}, \Sigma)$, og lad $\boldsymbol{\mu} \sim N_p(\boldsymbol{\mu}_0, \Sigma_0)$,

Da er aposteriorifordelingen af $\boldsymbol{\mu}$ efter observation af $\mathbf{X} = \mathbf{x}$ givet ved

$$\boldsymbol{\mu} | \mathbf{X} = \mathbf{x} \sim N_p(\mathbf{W}\boldsymbol{\mu}_0 + (\mathbf{I} - \mathbf{W})\mathbf{x}, (\mathbf{I} - \mathbf{W})\Sigma)$$

med

$$\mathbf{W} = \Sigma_0(\Sigma_0 + \Sigma)^{-1} \quad \text{og} \quad \mathbf{I} - \mathbf{W} = \Sigma_0(\Sigma_0 + \Sigma)^{-1}$$

Lad $\boldsymbol{\Gamma} = \Sigma_0\Sigma^{-1}$ betegne den generaliserede kvotient mellem variansen mellem grupper og variansen inden for grupper, da kan matricerne \mathbf{W} og $\mathbf{I} - \mathbf{W}$ udtrykkes ved

$$\mathbf{W} = (\mathbf{I} + \boldsymbol{\Gamma})^{-1} \quad \text{og} \quad \mathbf{I} - \mathbf{W} = (\mathbf{I} + \boldsymbol{\Gamma})^{-1}\boldsymbol{\Gamma}$$



Generel lineær mixed model



$$\mathbf{Y}_{n,1} = \mathbf{X}_{n,k}\boldsymbol{\beta}_{k,1} + \mathbf{Z}_{n,m}\mathbf{U}_{m,1} + \boldsymbol{\epsilon}_{n,1}$$

Fixed effects $\boldsymbol{\beta}$ (designmatrix \mathbf{X})

Random effects \mathbf{U} ("designmatrix" \mathbf{Z})

$\boldsymbol{\epsilon} \sim N_n(\mathbf{0}, \mathbf{R})$ og $\mathbf{U} \sim N_m(\mathbf{0}, \mathbf{D})$ uafhængige

$$E \begin{pmatrix} \mathbf{u} \\ \boldsymbol{\epsilon} \end{pmatrix} = \begin{pmatrix} \mathbf{0} \\ \mathbf{0} \end{pmatrix}$$

$$\mathbf{D} \begin{pmatrix} \mathbf{u} \\ \boldsymbol{\epsilon} \end{pmatrix} = \begin{pmatrix} \mathbf{D} & \mathbf{0} \\ \mathbf{0} & \mathbf{R} \end{pmatrix}$$

\mathbf{R} og \mathbf{D} kendt struktur, men eventuelt med ukendte parametre.



De hierarkiske modeller som generelle lineære mixede modeller

Uldballer:

$$\mathbf{X}_{28,1} = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}, \quad \mathbf{Z}_{28,7} = \begin{pmatrix} 1 & 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & 0 & 0 & \cdots & 0 \\ 0 & 1 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & 1 \end{pmatrix}$$

$$\mathbf{R} = \sigma^2 \mathbf{I}_{28}, \quad \mathbf{D} = \sigma_0^2 \mathbf{I}_7$$

idet nemlig

$$\mathbf{u} = \begin{pmatrix} u_1 \\ u_2 \\ \vdots \\ u_7 \end{pmatrix}$$

angiver det tilfældige bidrag fra hvert af de syv balleniveauer.



Ramushøjder:

$$\mathbf{X}_{20,2} = \begin{pmatrix} 1 & 8.0 \\ 1 & 8.5 \\ 1 & 9.0 \\ 1 & 9.5 \\ 1 & 8.0 \\ \vdots & \vdots \\ 1 & 9.5 \end{pmatrix}, \quad \mathbf{Z}_{20,5} = \begin{pmatrix} 1 & 8.0 & 0 & 0 & \cdots & 0 \\ 1 & 8.5 & 0 & 0 & \cdots & 0 \\ 1 & 9.0 & 0 & 0 & \cdots & 0 \\ 1 & 9.5 & 0 & 0 & \cdots & 0 \\ 0 & 0 & 1 & 8.0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & 9.5 \end{pmatrix}$$

$$\mathbf{R} = \sigma^2 \mathbf{I}_{20}, \quad \mathbf{D} = \Lambda_{2,2} \otimes \mathbf{I}_5 = \text{Blok diag } \Lambda_{2,2}$$

idet nemlig

$$\mathbf{u} = \begin{pmatrix} u_1 \\ u_2 \\ \vdots \\ u_{10} \end{pmatrix}$$

angiver det tilfældige bidrag til intercept, $u_{2\nu-1}$, og hældning, $u_{2\nu}$, fra hver af de fem drenge, $\nu = 1, 2, \dots, 5$.



Marginal fordeling



Marginal fordeling af \mathbf{Y}

$$E[\mathbf{Y}] = \mathbf{X}\boldsymbol{\beta}$$

$$\mathbf{D}[\mathbf{Y}] = \mathbf{V} = \mathbf{R} + \mathbf{ZDZ}'$$



Generel lineær mixed model



Er en meget rig klasse af lineære normalfordelingsmodeller

Forstyrrende tilfældige effekter (blokke i forsøgsplanlægning)
Interessante tilfældige effekter (variation mellem personer)

Kan også løses i ubalancede situationer.



PROC MIXED



Eksempel tilfældige regressionslinier

```
PROC MIXED DATA=ramus ;
CLASS DRENG;
MODEL ramus = ald /solution OUTP= pred ;
RANDOM int ald /subject = dreng solution type=un;
run;
```

Udfører estimation af varians-kovarians matrix for intercept og hældning, samt
BLUP-estimation af de tilfældige effekter.



PROC MIXED



SUBJECT = specifiserer de uafhængige blokke

TYPE = specifiserer strukturen af R og D

Mange muligheder for matricerne R og D

$$\begin{pmatrix} X'R^{-1}X & X'R^{-1}Z \\ Z'R^{-1}X & Z'R^{-1}Z + D^{-1} \end{pmatrix} \begin{pmatrix} \beta \\ u \end{pmatrix} = \begin{pmatrix} X'R^{-1}y \\ Z'R^{-1}y \end{pmatrix}$$

løses ved iteration - og iterativ bestemmelse af parametre i R og D

Estimaterne for u er BLUP-estimatorer