

Introducerende Statistik og Dataanalyse med R

Gruppespecifik regression

Jens Ledet Jensen



Formulering af generel lineær model

Flere regressionslinjer (gruppespecifik regression)

Teori: uafhængighed og χ^2 -fordeling

n observationer: x_1, x_2, \dots, x_n uafhængige

$$\left. \begin{array}{l} \text{alle normalfordelte} \\ \text{alle samme varians} \end{array} \right\} X_i \sim N(\bullet, \sigma^2)$$

$E(X_i) = \xi_i$: Model for ξ_i -erne

M : $\xi_i = \text{sum af led}$

Led: faktor deler ind i grupper, hver gruppe får sit eget bidrag

Led: regression, $\beta \cdot t_i$

Model i dag: β afhænger af gruppe defineret ved faktor

Den formelle definition af en generel lineær model:

$$\boldsymbol{\xi} = (\xi_1, \dots, \xi_n)^\top \in L \quad L \text{ lineært underrum af } R^n$$

Eksempel 1: $X_i \sim N(\mu, \sigma^2)$

$$\boldsymbol{\xi} = (\mu, \dots, \mu)^\top = \mu \mathbf{e} \in \text{Span}(\mathbf{e}), \quad \mathbf{e} = (1, 1, \dots, 1)^\top$$

$$\text{Span}(\mathbf{v}_1, \dots, \mathbf{v}_k) = \{\boldsymbol{\xi} = z_1 \mathbf{v}_1 + \dots + z_k \mathbf{v}_k \mid z_j \in R, j = 1, \dots, k\}$$

Eksempel 2: $X_i \sim N(\mu_{G_i}, \sigma^2)$, $G = (1, \dots, 1, 2, \dots, 2)$, $n = n_1 + n_2$

$$\boldsymbol{\xi} = (\mu_1, \dots, \mu_1, \mu_2, \dots, \mu_2)^\top = \mu_1 \mathbf{e}_1 + \mu_2 \mathbf{e}_2 \in \text{Span}(\mathbf{e}_1, \mathbf{e}_2)$$

$$\mathbf{e}_1 = (1, \dots, 1, 0, \dots, 0)^\top, \quad \mathbf{e}_2 = (0, \dots, 0, 1, \dots, 1)^\top$$

Eksempel 3: $X_i \sim N(\alpha + \beta t_i, \sigma^2)$

$$\boldsymbol{\xi} = (\alpha + \beta t_1, \dots, \alpha + \beta t_n)^\top = \alpha \mathbf{e} + \beta \mathbf{t} \in \text{Span}(\mathbf{e}, \mathbf{t})$$

$$\mathbf{e} = (1, \dots, 1)^\top, \mathbf{t} = (t_1, \dots, t_n)^\top$$

Generelt: en faktor G , der deler op i k grupper, giver k vektorer $\mathbf{e}_1, \dots, \mathbf{e}_k$, hvor den j 'te har 1 for de observationsnumre, der ligger i gruppe j , og nul ellers

En regressionsvariabel t giver en vektor \mathbf{t}

$$L = \text{Span}(\dots)$$

"plus" i modelformel giver ekstra vektorer i $\text{Span}(\dots)$

Faktor G deler ind i g grupper: vektorer $\mathbf{e}_1(G), \dots, \mathbf{e}_g(G)$

Faktor D deler ind i d grupper: vektorer $\mathbf{e}_1(D), \dots, \mathbf{e}_d(D)$

Underrum L_G og L_D har ikke-trivielt snit:

$$\mathbf{e}_1(G) + \dots + \mathbf{e}_g(G) = \mathbf{e}_1(D) + \dots + \mathbf{e}_d(D)$$

$$\dim(\text{Span}(\mathbf{e}_1(G), \dots, \mathbf{e}_g(G), \mathbf{e}_1(D), \dots, \mathbf{e}_d(D))) = g + d - 1$$

Estimere parametre i middelværdimodel = finde $\xi \in L$

Minimere $\sum_i (x_i - \xi_i(M))^2$ samme som at finde projektion af $\mathbf{x} = (x_1, \dots, x_n)^T$ på L :

$$\sum_i (x_i - \xi_i)^2 = \|\mathbf{x} - \xi\|^2 = \|\mathbf{x} - \mathbf{P}\mathbf{x} + \mathbf{P}\mathbf{x} - \xi\|^2 = \|\mathbf{x} - \mathbf{P}\mathbf{x}\|^2 + \|\mathbf{P}\mathbf{x} - \xi\|^2$$

\mathbf{P} er projektionsmatricen, $\mathbf{x} - \mathbf{P}\mathbf{x}$ er vinkelret på alt i L , $\hat{\xi} = \mathbf{P}\mathbf{x}$

$L = \text{Span}(\mathbf{h}_1, \dots, \mathbf{h}_k)$, $\mathbf{h}_1, \dots, \mathbf{h}_k$ lineært uafhængige

$$\mathbf{P} = \mathbf{H}(\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T$$

$\mathbf{h}_1, \dots, \mathbf{h}_k$ er søjlerne i \mathbf{H}

Vise at $\mathbf{P} = \mathbf{H}(\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T$ er projektionsmatricen:

Typisk vektor i L : $\boldsymbol{\xi} = \mathbf{H}\boldsymbol{\theta}$, $\boldsymbol{\theta} \in \mathbf{R}^k$

vise vinkelret på $\mathbf{x} - \mathbf{P}\mathbf{x}$

$$\begin{aligned}(\mathbf{x} - \mathbf{P}\mathbf{x})^T \mathbf{H}\boldsymbol{\theta} &= \mathbf{x}^T \mathbf{H}\boldsymbol{\theta} - \mathbf{x}^T \mathbf{H}(\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{H}\boldsymbol{\theta} \\ &= \mathbf{x}^T \mathbf{H}\boldsymbol{\theta} - \mathbf{x}^T \mathbf{H}\boldsymbol{\theta} = 0\end{aligned}$$

Eksempel: projektion på 1.aksen i \mathbf{R}^2 , $\mathbf{H} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$

$$\begin{pmatrix} 1 \\ 0 \end{pmatrix} ((1, 0) \begin{pmatrix} 1 \\ 0 \end{pmatrix})^{-1} (1, 0) = \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix}$$

M_2 er en undermodel af M_1 : $L_2 \subset L_1$

Vise: $\text{SSD}(M_2) = \text{SSD}(M_1) + \text{SSD}(M_1, M_2)$

$$\text{SSD}(M_1, M_2) = \sum_i (\hat{\xi}_i(M_1) - \hat{\xi}_i(M_2))^2$$

$\mathbf{P}_1, \mathbf{P}_2$ er projektionsmatricerne

$$\text{SSD}(M_2) = \|\mathbf{x} - \mathbf{P}_2\mathbf{x}\|^2 = \|\mathbf{x} - \mathbf{P}_1\mathbf{x} + \mathbf{P}_1\mathbf{x} - \mathbf{P}_2\mathbf{x}\|^2$$

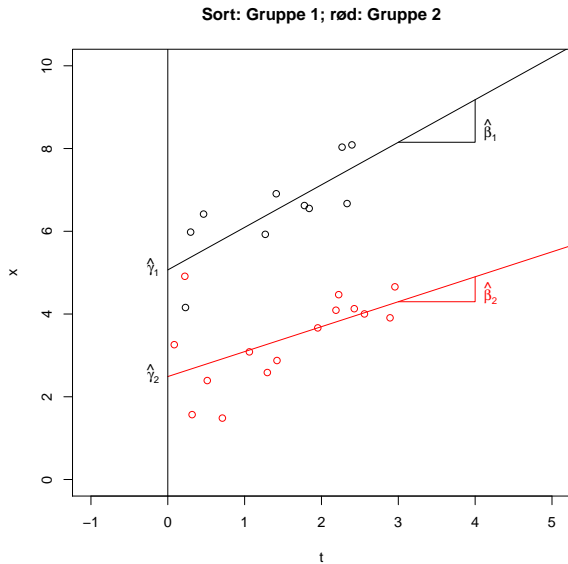
$$= \|\mathbf{x} - \mathbf{P}_1\mathbf{x}\|^2 + \|\mathbf{P}_1\mathbf{x} - \mathbf{P}_2\mathbf{x}\|^2$$

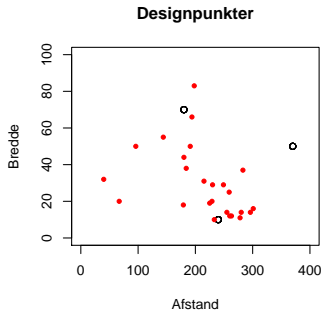
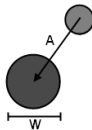
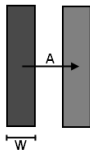
idet $\mathbf{P}_1\mathbf{x} - \mathbf{P}_2\mathbf{x} \in L_1$

Den generelle lineære model er indført via underrum af R^n

Næste: Model for gruppespecifik regression

Regression delt op i grupper





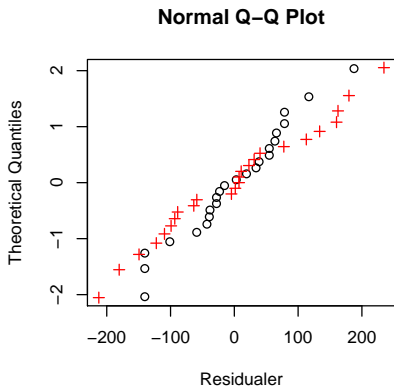
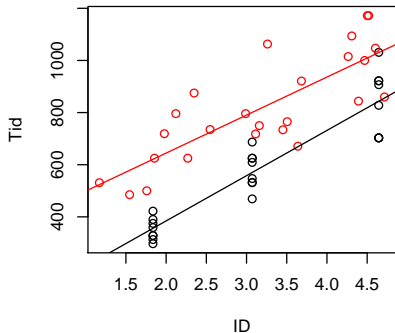
Fitts lov: Tid til at flytte mus er lineær i $ID = \log_2 \left(\frac{A+W}{W} \right)$

Rød: cirkeldesign, sort: rektangeldesign

Data

Eksperiment	Afstand	Bredde	Tid
Rektangel	370	50	625
Rektangel	370	50	625
Rektangel	370	50	547
Rektangel	370	50	531
Rektangel	370	50	469
Rektangel	370	50	609
Rektangel	370	50	532
Rektangel	370	50	687
Rektangel	240	10	703
Rektangel	240	10	922
:			
Cirkel	278	11	860
Cirkel	283	37	718
Cirkel	40	32	531
Cirkel	233	10	1047
Cirkel	191	50	625
Cirkel	179	18	734

Plot af data



- 1) Teste lineær sammenhæng for Rektangel-data (sort) (sidste uge)
- 2) Teste samme varians omkring de to linjer
- 3) Teste samme hældning i de to linjer

TidRekt, IdRekt: målinger fra Rektangeldata

TidCirc, IdCirc: målinger fra Cirkeldata

Tid = (TidRekt, TidCirc)

Id = (IdRekt, IdCirc)

Design=factor(rep(c("Rekt", "Circ"), c(24, 25)), levels=c("Rekt", "Circ"))

levels=c("Rekt", "Circ") for at undgå lexicografisk ordning

Model:

$$\text{TidRekt}_i \sim N(\alpha_1 + \beta_1 \cdot \text{ldRekt}_i, \sigma_1^2)$$

$$\text{TidCirc}_i \sim N(\alpha_2 + \beta_2 \cdot \text{ldCirc}_i, \sigma_2^2)$$

$(\alpha_1, \alpha_2, \beta_1, \beta_2, \sigma_1^2, \sigma_2^2)$ kan variere frit

Teste $\sigma_1^2 = \sigma_2^2$

$$\text{lmUD1} = \text{lm}(\text{TidRekt} \sim \text{ldRekt}); \quad \text{lmUD2} = \text{lm}(\text{TidCirc} \sim \text{ldCirc})$$

$$s_{\text{Rekt}} = 85, \quad s_{\text{Circ}} = 122$$

$$\text{var.test}(\text{lmUD1}, \text{lmUD2}) \quad (\text{bartlett.test}(\text{list}(\text{lmUD1}, \text{lmUD2})))$$

Resultat: $p\text{-værdi} = 0.10$

Konklusion: data strider ikke mod samme varians i de to eksperimenter

Undersøg om der er samme varians omkring linje for 4, 6 eller 8 cylindre

Model: $\text{mpg} \sim N(\alpha_{D_i} + \beta_{D_i} \text{vaegt}_i, \sigma_{D_i}^2)$, hypotese: $\sigma_4^2 = \sigma_6^2 = \sigma_8^2$

```
mpg=mtcars[,1]
```

```
vaegt=mtcars[,6]
```

```
Design=factor(mtcars[,2])
```

```
lmUD1=lm(mpg[Design=="4"]~vaegt[Design=="4"])
```

```
lmUD2=lm(mpg[Design=="6"]~vaegt[Design=="6"])
```

```
lmUD3=lm(mpg[Design=="8"]~vaegt[Design=="8"])
```

Modelkontrol: `plot(vaegt[Design=="4"],lmUD1$residuals)` og
`qqnorm(lmUD1$residuals,datax=TRUE)`

```
bartlett.test(list(lmUD1,lmUD2,lmUD3))
```

Konklusion: Data strider ikke mod samme varians omkring linjen da ...

Tid: alle tidsmålinger

Id: alle Id-værdier

Design: factor med værdier **Rekt** og **Circ**

$$\text{Model: } \text{Tid}_i \sim N(\alpha_{\text{Design}_i} + \beta_{\text{Design}_i} \cdot \text{Id}_i, \sigma^2)$$

Model hvor hvert eksperiment har sin egen lineære sammenhæng

Modelformel, R: $\text{Tid} \sim \text{Design} * \text{Id}$ samme som $\text{Tid} \sim \text{Design} + \text{Design} * \text{Id}$

Design giver α_{Rekt} og α_{Circ}

$\text{Design} * \text{Id}$ giver $\beta_{\text{Rekt}} \cdot \text{Id}_i$ og $\beta_{\text{Circ}} \cdot \text{Id}_i$

Vektoren af middelværdier kan skrives som $\xi = \alpha_1 \mathbf{e}_1 + \alpha_2 \mathbf{e}_2 + \beta_1 \mathbf{t}_1 + \beta_2 \mathbf{t}_2$

$$\mathbf{e}_1 = \begin{cases} 1 & \text{på gruppe 1} \\ 0 & \text{på gruppe 2} \end{cases} \quad \mathbf{e}_2 = \begin{cases} 0 & \text{på gruppe 1} \\ 1 & \text{på gruppe 2} \end{cases}$$

$$\mathbf{t}_1 = \begin{cases} t_i & \text{på gruppe 1} \\ 0 & \text{på gruppe 2} \end{cases} \quad \mathbf{t}_2 = \begin{cases} 0 & \text{på gruppe 1} \\ t_i & \text{på gruppe 2} \end{cases}$$

Generelt: underrum L udspændes af søjlerne i \mathbf{H} , $\xi = \mathbf{H}\theta$

$$\hat{\theta} = (\mathbf{H}^\top \mathbf{H})^{-1} \mathbf{H}^\top \mathbf{x}$$

Antag $\sum_I t_i = \sum_{II} t_i = 0$, hvor I og II refererer til gruppe 1 og gruppe 2
 Lad $\mathbf{H} = (\mathbf{e}_1, \mathbf{e}_2, \mathbf{t}_1, \mathbf{t}_2)$

$$\mathbf{H}^T \mathbf{H} = \begin{pmatrix} n_1 & 0 & 0 & 0 \\ 0 & n_2 & 0 & 0 \\ 0 & 0 & SS_1 & 0 \\ 0 & 0 & 0 & SS_2 \end{pmatrix}, \quad SS_1 = \sum_I t_i^2, \quad SS_2 = \sum_{II} t_i^2$$

$$\mathbf{H}^T \mathbf{x} = (\sum_I x_i, \sum_{II} x_i, \sum_I t_i x_i, \sum_{II} t_i x_i)^T$$

Dette giver estimerne

$$\hat{\alpha}_1 = \frac{\sum_I x_i}{n_1}, \quad \hat{\alpha}_2 = \frac{\sum_{II} x_i}{n_2}, \quad \hat{\beta}_1 = \frac{\sum_I t_i x_i}{SS_1}, \quad \hat{\beta}_2 = \frac{\sum_{II} t_i x_i}{SS_2}$$

Parametrisering i R af modellen: $\text{Design} + \text{Design} * \text{Id}$

$$\text{Intercept} = \alpha_{\text{Rekt}}, \quad \text{DesignCirc} = \alpha_{\text{Circ}} - \alpha_{\text{Rekt}}$$

$$\text{Id} = \beta_{\text{Rekt}}, \quad \text{DesignCirc:Id} = \beta_{\text{Circ}} - \beta_{\text{Rekt}}$$

```
Design=factor(rep(c("Rekt","Circ"),c(24,25)),levels=c("Rekt","Circ"))
```

```
summary(lm(Tid~F+F*Id))
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	37.10	63.32	0.586	0.56089
DesignCirc	315.89	91.46	3.454	0.00122
Id	173.62	18.71	9.280	5.12e-12
DesignCirc:Id	-27.82	27.03	-1.029	0.30888

Residual standard error: 105.3 on 45 degrees of freedom

Bemærk *levels=...* for at undgå lexicografisk ordning

Stiger *Tid* med *index of difficulty* på samme måde i de to eksperimenter?

Hypotese: $\beta_{\text{Circ}} = \beta_{\text{Rekt}}$

Teste reduktion fra model $\text{Tid}_i \sim N(\alpha_{\text{Design}_i} + \beta_{\text{Design}_i} \cdot \text{Id}_i, \sigma^2)$
til $\text{Tid}_i \sim N(\alpha_{\text{Design}_i} + \beta \cdot \text{Id}_i, \sigma^2)$

fra model $\text{Tid} \sim \text{Design} * \text{Id}$ til model $\text{Tid} \sim \text{Design} + \text{Id}$

Fortolkning af model:

Additivitet: uanset eksperiment er der samme forskel i middelværdi mellem to værdier af *Id*

uanset værdien af *Id* er der samme forskel i middelværdi mellem de to eksperimenter

Generelle F -test:

$$F = \frac{(SSD(M_2) - SSD(M_1)) / (df(M_2) - df(M_1))}{SSD(M_1) / df(M_1)} \sim F(df(M_2) - df(M_1), df(M_1))$$

```
anova(lm(Tid~Design+Id),lm(Tid~Design*Id))
```

```
Model 1: Tid ~ Design + Id
```

```
Model 2: Tid ~ Design * Id
```

	Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
1	46	510747.7				
2	45	499001.3	1	11746.4	1.0593	0.3089

Konklusion: data strider ikke mod hypotesen om samme hældning

$\beta_{\text{Rekt}} = \beta_{\text{Circ}}$ (p -værdi = 0.31)

Gå tilbage til parametertabel og find p -værdi der

Vi beskriver data med modellen $Tid_i \sim N(\alpha_{Design_i} + \beta \cdot Id_i, \sigma^2)$

teste reduktion til model $N(\alpha + \beta \cdot Id_i, \sigma^2)$

fortolkning: samme lineære sammenhæng i de to eksperimenter

```
anova(lm(Tid~Id),lm(Tid~Design+Id))
```

Model 1: Tid ~ Id

Model 2: Tid ~ Design + Id

	Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
1	47	1141624.8				
2	46	510747.7	1	630877.1	56.8193	0.0000 ***

Konklusion: data strider kraftigt mod hypotesen om samme skæring

$\alpha_{Rekt} = \alpha_{Circ}$

(Teste samme skæring i fulde model: p -værdi = 0.0012)


```
summary(lm(Tid~Design+Id))
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	79.54	48.09	1.654	0.105
DesignCirc	227.00	30.11	7.538	1.43e-09
Id	160.29	13.51	11.864	1.35e-15

Residual standard error: 105.4 on 46 degrees of freedom

F-test for at fjerne 1 parameter = t-test kvadreret

Fulde model: Residual standard error: 105.3 on 45 degrees of freedom

Konfidensintervaller i model $Tid_i \sim N(\alpha_{Design_i} + \beta \cdot Id_i, \sigma^2)$

```
confint(lm(Tid~Design+Id))
```

	2.5 %	97.5 %
(Intercept)	-17.26535	176.3364
DesignCirc	166.37933	287.6124
Id	133.09496	187.4859

I cirkel-eksperiment er tidsforbruget cirka 200 (ms) over rektangel-eksperiment

Fordobling af sværhedsgrad: stigning på cirka 160 ms

Konfidensinterval for spredning omkring linjen:

Webbog afsnit 2.6 med $s^2 = 105.4^2$ og $df = 46$: [88, 132]

Undersøg om der er samme hældning for de tre designs

Model: ..., Hypotese: ...

```
mpg=mtcars[,1]  
vaegt=mtcars[,6]  
Design=factor(mtcars[,2])
```

Vi bruger det generelle F -test fra afsnit 4.7

```
anova(lm(mpg~Design+vaegt),lm(mpg~Design*vaegt))
```

$F=...$ som vurderes i en $F(.,.)$ -fordeling, og p -værdien for testet er ... Da denne er under 0.05, siger vi, at ...

Vores model er nu ...

Skift til webbog afsnit 5.1

Samme model kan parametriseres på flere måder, for eksempel

$$\alpha_1, \alpha_2 - \alpha_1, \alpha_3 - \alpha_1, \dots, \alpha_k - \alpha_1$$

$$\text{eller: } \alpha_1, \alpha_2, \alpha_3, \dots, \alpha_k$$

Opskrevet med vektorer svarer dette til:

$$\alpha_1 \mathbf{e}_1 + \alpha_2 \mathbf{e}_2 + \dots + \alpha_k \mathbf{e}_k$$

$$= \alpha_1 (\mathbf{e}_1 + \dots + \mathbf{e}_k) + (\alpha_2 - \alpha_1) \mathbf{e}_2 + \dots + (\alpha_k - \alpha_1) \mathbf{e}_k$$

Model med gruppebestemt regressionkoefficient er omtalt

Næste: Forstå fordelingsudsagn og uafhængighedsudsagn

Først: vektor af normalfordelte stokastiske variable

Fitts law for thumbs

Vektor af stokastiske variable: $\mathbf{Z} = (Z_1, \dots, Z_n)^\top$

$$E(\mathbf{Z}) = (E(Z_1), \dots, E(Z_n))^\top$$

$$\text{Var}(\mathbf{Z}) = \begin{pmatrix} \text{Var}(Z_1) & \text{Cov}(Z_1, Z_2) & \cdots & \text{Cov}(Z_1, Z_n) \\ \text{Cov}(Z_2, Z_1) & \text{Var}(Z_2) & \cdots & \text{Cov}(Z_2, Z_n) \\ \vdots & \vdots & & \vdots \\ \text{Cov}(Z_n, Z_1) & \text{Cov}(Z_n, Z_2) & \cdots & \text{Var}(Z_n) \end{pmatrix}$$

$$\text{Var}(Z) = E(Z - E(Z))^2, \quad \text{Cov}(Z_1, Z_2) = E(Z_1 - E(Z_1))(Z_2 - E(Z_2))$$

Kendte regneregler:

$$E\left(\sum_i a_i Z_i\right) = \sum_i a_i E(Z_i)$$

$$\text{Var}\left(\sum_i a_i Z_i\right) = \sum_{i,j} a_i a_j \text{Cov}(Z_i, Z_j)$$

$$\text{Cov}\left(\sum_i a_i Z_i, \sum_i b_i Z_i\right) = \sum_{i,j} a_i b_j \text{Cov}(Z_i, Z_j)$$

Dette giver $E(\mathbf{BZ}) = \mathbf{B}E(\mathbf{Z})$ og $\text{Var}(\mathbf{BZ}) = \mathbf{B}\text{Var}(\mathbf{Z})\mathbf{B}^\top$, $\mathbf{B} : k \times n$

Bevis:

$$(\mathbf{BZ})_i = \sum_j B_{ij} Z_j \rightarrow E(\mathbf{BZ})_i = \sum_j B_{ij} E(Z_j) = (\mathbf{B}E(\mathbf{Z}))_i$$

$$\begin{aligned}\text{Cov}\left(\sum_i B_{ui} Z_i, \sum_j B_{vj} Z_j\right) &= \sum_{i,j} B_{ui} B_{vj} \text{Cov}(Z_i, Z_j) \\ &= \sum_{i,j} B_{ui} \text{Var}(\mathbf{Z})_{ij} \mathbf{B}_{jv}^\top \\ &= (\mathbf{B}\text{Var}(\mathbf{Z})\mathbf{B}^\top)_{uv}\end{aligned}$$

$Z_i \sim N(\mu_i, \sigma^2)$, $i = 1, \dots, n$, uafhængige, skrives som

$$\mathbf{Z} \sim N_n(\boldsymbol{\mu}, \sigma^2 \mathbf{I}), \quad \mathbf{I} \text{ diagonal med } 1 \text{ i diagonalen}$$

Da $E(\mathbf{BZ}) = \mathbf{B}E(\mathbf{Z})$ og $\text{Var}(\mathbf{BZ}) = \mathbf{B}\text{Var}(\mathbf{Z})\mathbf{B}^\top$ skriver vi

$$\mathbf{W} = \mathbf{BZ} \sim N_n(\mathbf{B}\boldsymbol{\mu}, \sigma^2 \mathbf{B}\mathbf{B}^\top),$$

Vi skriver nu $\mathbf{W} \sim N_n(\boldsymbol{\xi}, \boldsymbol{\Sigma})$, $\boldsymbol{\xi} = \mathbf{B}\boldsymbol{\mu}$, $\boldsymbol{\Sigma} = \sigma^2 \mathbf{B}\mathbf{B}^\top$ og får

$$\mathbf{AW} \sim N_n(\mathbf{A}\boldsymbol{\xi}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top)$$

idet $\mathbf{AW} = (\mathbf{AB})\mathbf{Z} \sim N_n((\mathbf{AB})\boldsymbol{\mu}, \sigma^2 (\mathbf{AB})(\mathbf{AB})^\top)$

$$= N_n(\mathbf{A}\boldsymbol{\xi}, \mathbf{A}(\sigma^2 \mathbf{B}\mathbf{B}^\top)\mathbf{A}^\top)$$

Vi har etableret notation for vektor med normalfordelte variable

Næste: fundamentale spaltningssætning (uden bevis)

Opsætning:

$$\mathbf{Y} \sim N_n(\mathbf{0}, \sigma^2 \mathbf{I}), \quad \mathbf{B}_1 : k_1 \times n, \quad \mathbf{B}_2 : k_2 \times n$$

$$\mathbf{Z}_1 = (Z_{11}, \dots, Z_{1,k_1})^\top = \mathbf{B}_1 \mathbf{Y} \quad \text{og} \quad \mathbf{Z}_2 = (Z_{21}, \dots, Z_{2,k_2})^\top = \mathbf{B}_2 \mathbf{Y}$$

Hvis $\mathbf{B}_1 \mathbf{B}_2^\top = \mathbf{0}$ så er \mathbf{Z}_1 og \mathbf{Z}_2 stokastisk uafhængige

En $n \times n$ matrice \mathbf{B} er en ortogonal projektionsmatrice hvis og kun hvis

$$\mathbf{B}\mathbf{B} = \mathbf{B} \text{ og } \mathbf{B}^T = \mathbf{B}$$

Definer $\mathbf{Z} = (Z_1, \dots, Z_n)^T = \mathbf{B}\mathbf{Y}$. Så gælder

$$\sum_{i=1}^n Z_i^2 = \|\mathbf{B}\mathbf{Y}\|^2 = (\mathbf{B}\mathbf{Y})^T(\mathbf{B}\mathbf{Y}) \sim \sigma^2 \chi^2(k)$$

k : rang af \mathbf{B} = dimensionen af rum der projektteres ned på

Spaltningssætningen er omtalt

Bruge sætningen

Et normalfordelt observationssæt

$$X_i \sim N(\mu, \sigma^2), i = 1, \dots, n, \quad Y_i = X_i - \mu \sim N(0, \sigma^2)$$

$$\bar{X} = \frac{X_1}{n} + \dots + \frac{X_n}{n} = \bar{Y} + \mu = \mathbf{B}_1 \mathbf{Y} + \mu$$

$$\mathbf{B}_1 = \mathbf{e}^\top = \left(\frac{1}{n}, \frac{1}{n}, \dots, \frac{1}{n}\right)$$

$$\text{SSD} = \sum_i (X_i - \bar{X}_i)^2 = \sum_i (Y_i - \bar{Y})^2 = \|(\mathbf{I} - \mathbf{E})\mathbf{Y}\|^2 = \|\mathbf{B}_2 \mathbf{Y}\|^2$$

$$\mathbf{E}, n \times n, \text{ alle søjler lig med } \mathbf{e}, \quad \mathbf{B}_2 = \mathbf{I} - \mathbf{E}$$

Betingelse fra spaltningssætningen:

$$\mathbf{B}_1 \mathbf{B}_2^\top = \mathbf{e}^\top (\mathbf{I} - \mathbf{E}) = \left(\frac{1}{n} - \frac{1}{n}\right)(1, 1, \dots, 1) = \mathbf{0}$$

\bar{X} og SSD er uafhængige og $\text{SSD} \sim \sigma^2 \chi^2(n-1)$ idet

$$\mathbf{B}_2 \mathbf{B}_2 = (\mathbf{I} - \mathbf{E})(\mathbf{I} - \mathbf{E}) = \mathbf{I} - 2\mathbf{E} + \mathbf{E}\mathbf{E} = \mathbf{I} - \mathbf{E} = \mathbf{B}_2$$

Uafhængighed mellem skøn over parametre i middelværdi og variansskøn:
webbog afsnit 6.2

Model: $\xi \in L$, L udspændt af søjlerne i H : $\xi = H\theta$

$$\hat{\theta} = B_1 Y + \theta, \quad B_1 = (H^T H)^{-1} H^T, \quad Y = X - \xi$$

Projektionsmatrix: $P = H(H^T H)^{-1} H^T$

$$\text{SSD}(M) = \|B_2 Y\|^2, \quad B_2 = I - P$$

Betingelse i spaltningssætningen:

$$B_1 B_2^T = P(I - P) = P - PP = P - P = 0$$

Spaltningssætningen har givet os de resultater vi har brugt tidligere

Slut for i dag - efter en meget teoritung dag