

Nynorsk



Kontakt under eksamen: Thiago G. Martins      46 93 74 29

## EKSAMEN I TMA4315 GENERALISERTE LINEÆRE MODELLAR

Torsdag 13. desember, 2012  
Tid: 09:00 – 13:00

Tilletne hjelpemiddel:

Tabeller og formler i statistikk, Tapir Forlag

K. Rottmann: Matematisk formelsamling

Calculator HP30S / CITIZEN SR-270X

Gult, stempla A4-ark med eigne handskrivne notat.

Sensur: 10. januar, 2013

### Oppgåve 1 Nedbør i Trondheim i morgen?

NTNU matematikkstudenten Konrad Kontrollfrik liker å vere forberedt til neste dag, og han vil vite om det blir nedbør eller ikkje i morgen tidleg. Han lagar derfor statistiske modellar for nedbør førekommst, og vurderer tre forklaringsvariable:

1. Mengde nedbør (i  $mm$ ) ifølgje venværslelet i morgen ( $Fore$ ).
2. Ein binær variabel  $ForeBin$  som indikerer om venværslelet seier nedbør ( $ForeBin = 1$ ) eller ikkje nedbør ( $ForeBin = 0$ ).
3. Ein variabel  $OF$  som indikerer korleis gårdsdagens varsel og dagens observasjon stemmer:
  - $OF = 0$ : ingen nedbør observert, og ingen nedbør i varsel.
  - $OF = 1$ : nedbør observert, men ingen nedbør i varsel

- $OF = 2$ : ingen nedbør observert, men nedbør i varsel
- $OF = 3$ : nedbør observert, og nedbør i varsel

Konrad har fått tak i data for 100 påfølgjande dagar. Data for dei første ti er gjeve i Tabell 1.

Han bruker  $R$  for å tilpasse desse modellane (sjå redigert utskrift under); *model 1* gjev **result1**, *model 2* gjev **result2** og *model 3* gjev **result3**.

```
summary(result1)
```

Call:

```
glm(formula = Occur ~ -1 + Fore + ForeBin + as.factor(OF),
family = binomial(link = "logit"), data = OccurData)
```

Deviance Residuals:

Min	1Q	Median	3Q	Max
-2.1714	-0.6614	-0.5975	0.7493	2.2624

Coefficients:

	Estimate	Std. Error	z value	Pr(> z )
Fore	0.4609	0.2294	2.009	0.044494 *
ForeBin	0.8883	0.6559	1.354	0.175641
as.factor(OF)0	-1.6327	0.4349	-3.755	0.000174 ***
as.factor(OF)1	-1.1941	0.9690	-1.232	0.217818
as.factor(OF)2	-2.5144	0.6764	-3.717	0.000202 ***
as.factor(OF)3	-1.7609	0.5983	-2.943	0.003249 **

---

Signif. codes: 0 ‘\*\*\*’ 0.001 ‘\*\*’ 0.01 ‘\*’ 0.05 ‘.’ 0.1 ‘ ’ 1

(Dispersion parameter for binomial family taken to be 1)

```
Null deviance: 138.63 on 100 degrees of freedom
Residual deviance: 101.09 on 94 degrees of freedom
AIC: 113.09
```

```
> summary(result2)
```

Call:

```
glm(formula = Occur ~ Fore + ForeBin + OF, family = binomial(link = "logit"),
```

```
data = OccurData)
```

Coefficients:

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	-1.6438	0.4280	-3.840	0.000123 ***
Fore	0.5204	0.2224	2.340	0.019286 *
ForeBin	0.7344	0.6396	1.148	0.250832
OF	-0.1355	0.2026	-0.669	0.503589

---

Signif. codes: 0 ‘\*\*\*’ 0.001 ‘\*\*’ 0.01 ‘\*’ 0.05 ‘.’ 0.1 ‘ ’ 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 125.37 on 99 degrees of freedom  
 Residual deviance: 103.19 on 96 degrees of freedom  
 AIC: 111.19

```
> summary(result3)
```

Call:

```
glm(formula = Occur ~ Fore, family = binomial(link = "logit"),
  data = OccurData)
```

Deviance Residuals:

Min	1Q	Median	3Q	Max
-2.1246	-0.6984	-0.6146	0.7471	1.8759

Coefficients:

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	-1.5707	0.3152	-4.983	6.27e-07 ***
Fore	0.6683	0.1856	3.601	0.000317 ***

---

Signif. codes: 0 ‘\*\*\*’ 0.001 ‘\*\*’ 0.01 ‘\*’ 0.05 ‘.’ 0.1 ‘ ’ 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 125.37 on 99 degrees of freedom  
 Residual deviance: 104.81 on 98 degrees of freedom  
 AIC: 108.81

Occure	Fore	ForeBin	OF
1	3.0	1	3
0	0.5	1	3
0	0.0	0	2
0	0.0	0	0
0	0.0	0	0
0	0.0	0	0
0	0.7	1	0
0	0.4	0	2
0	0.0	0	0
1	0.4	0	0

Tabell 1: Data frå ti dagar i Konrad sitt datasett om nedbør forkomst, ingen nedbør (Occure=0) eller nedbør (Occure = 1). Og tilgjengeleg; mengde nedbør frå vervesel i mm, binær varsel og *OF* variabelen.

- a) Sett opp matematisk dei generaliserte lineære modellane (GLM) som er brukt. Spesifiser og diskuter antakingar.
- Vidare skriv ut designmatrisa  $X$  for dei første 6 observasjonane for kvar modell.  
 Når det er relevant, spesifiser kva strategi som er brukt for å sikre identifiserbarhet.  
 Beskriv kort forskjellane mellom *model 1* og *model 2*.
- b) Basert på resultata frå R svar på dei følgjande spørsmåla:  
 I følgje *model 1*: Kva er sannsynet for nedbør dersom varselet er på 5mm og *OF* = 0 ?  
 I følgje *model 2*:Kva er sannsynet for nedbør dersom varselet er på 5mm og *OF* = 3 ?  
 I følgje *model 3*: Kva er odds ratio mellom ein dag med varsel 0mm og ein dag med 5mm?
- c) Konrad vil no samanlikne modellar: Sett opp ei hypotese for å samanlikne modell 1 mot modell 3 ved å bruke likelihood ratio testen (dvs basert på deviance), og utfør testen.  
 Kva modell vil du føretrekke, modell 1, modell 2 eller modell 3. Kvifor?
- d) Konrad er og interessert i nedbør førekomst på Trondheim Lufthamn Værnes og på Vassfjellet. Han har observasjonar og varsel også for desse stadane;  $Location \in \{\text{Trondheim}, \text{Vaernes}, \text{Vassfjellet}\}$  for dei same 100 dagane.
- Han tilpassar tre modellar ved å kjøre R kommandoane under:
- ```
res4 = glm(Occur~Fore, family=binomial(link="logit"), data=OccurDataTVV)
res5 = glm(Occur~Fore + Location, family=binomial(link="logit"), data=OccurDataTVV)
res6 = glm(Occur~Fore*Location, family=binomial(link="logit"), data=OccurDataTVV)
```
- Forklar kort dei tre modellane, og lag skisser for å illustrerer.

Sjå på modellen som er brukt for å få `res4`, og diskuter om antakingane for GLMar er oppfylt no. Foreslå ein alternativ modell.

### Oppgåve 2 Nedbør i Trondheim som snø, sludd eller regn?

Gjeve at det blir nedbør i morgen, så er Konrad interessert i kva type nedbør det blir.

Vi kallar denne storleiken  $C$ , og klassifiserer nedbør i tre klassar; snø ( $C = 1$ ), sludd ( $C = 2$ ) og regn ( $C = 3$ ). Vidare vil vi undersøke om temperaturvarselet er ein god forklaringsvariabel.

- a) La  $C_i$  vere klassen av nedbør for dag  $i$ , og la  $t_i$  vere temperaturvarselet som er gyldig for dag  $i$ .

Foreslå ein passande modell for  $C$ . Spesifiser modellen matematisk, og forklar han med ord/figur(ar). Spesielt forklar korleis parametra skal bli tolka.

### Oppgåve 3 Nedbør i Trondheim, mengde

Vi vil no modellere mengda av nedbør i løpet av eit døgn, gjeve at det er nedbør, og kallar denne storleiken  $Y$ . Det er vanleg å modellere  $Y$  som gammafordelt  $Y \sim Ga(\alpha, \beta)$ , med sannsynstettleik;

$$f_Y(y) = \frac{1}{\beta^\alpha \Gamma(\alpha)} y^{\alpha-1} \exp(-y/\beta).$$

I denne oppgåva skal vi anta det er  $N$  observasjonar, kvar gammafordelt  $Y_i \sim Ga(\alpha_i, \mu_i/\alpha_i)$ . Her er  $\alpha_i$ ane antatt kjent, det vil seie at dei er nuisance parameter.

- a) Vis at gammafordelinga er medlem i den eksponensielle familien når  $\mu_i$  er parameteren av interesse.

Bruk dette til å finne uttrykk for forventningverdien og variansen til  $Y_i$ , som funksjon av  $(\alpha_i, \mu_i)$ . Frå dette tolk  $\alpha_i$ .

- b) Forklar kva ein metta (saturated) modell er.

Sett opp log-likelihood funksjonen uttrykt ved  $\mu_i$ , og bruk dette til å finne maximum likelihood estimatorane for  $\mu_i$ -ane for den metta modellen.

Finn deviancen (basert på alle  $N$  observasjonane).

- c) Vi vil no lage ein modell for mengde nedbør (gjeve at det er nedbør) med nedbørvarsle som forklaringsvariabel.

La  $Y_i$  vere mengde nedbør på dag  $i$ , og la  $x_i$  vere nedbørvarslelet som er gyldig for dag  $i$ . Sett opp ein GLM for dette, og begrunn valet ditt for link-funksjon og lineærkomponent.