



## Dagens program:

Mandag den 28. februar

### Generaliserede lineære modeller

- Motiverende eksempel, logistisk regression
- Eksponentielle dispersionsfamilier
- Generaliseret lineær model
- Estimation, fordeling af residualer mv
- Test for modelreduktion
- Eksempler på modeller

Datasæt musefostre og blommer



## Andel døde musefostre som fkt af indgivet konc.:

Hvad er der galt med lineær model og mindste kvadraters metode (GLM) ?

- Andelen (responsen) skal ligge mellem nul og 1
- Variansen på responset er ikke konstant, men afhænger af middelværdien for responset, så observationerne bør ikke tillægges samme vægt

Det er i hvert fald et problem, når vi arbejder med "tællel", dvs for binomialfordeling og Poissonfordeling



## Toxitetsvurdering for kølevæske:

Indeks	Antal døde	Antal fostre	Indgivent koncentration [mg/kg pr. dag]
$i$	$z_i$	$n_i$	$x_i$
1	15	297	0
2	17	242	62.5
3	22	312	125
4	38	299	250
5	144	285	500



## Eksponentielle familier af fordelinger:

**Familie** af fordelinger: En indekseret samling af fordelinger

Eksponentielle familier er sådan at **summen** af uafhængige obs fra fordelingen er **sufficient**, dvs at man kan tillade sig at reducere et observationssæt til summen af obs uden at miste information.

**Exponentiel familie**, fælles grundform:

$$f_X(x; \theta) = c(x) \exp\{\theta x - \kappa(\theta)\} \quad (4.2)$$

Vi udvider til **eksponentielle dispersionsfamilier** med grundformen

$$f_Y(y; \theta, \lambda) = c(y, \lambda) \exp[\lambda\{\theta y - \kappa(\theta)\}] \quad (4.5)$$

## Eksempel, binomialfordeling:



$$g_Y(y; p) = \binom{n}{ny} p^{ny} (1-p)^{n(1-y)} = \binom{n}{x} (1-p)^n \left(\frac{p}{1-p}\right)^{ny} \quad (4.6)$$

for  $y = 0, 1/n, \dots, 1$  og  $0 \leq p \leq 1$ .

Lad

$$\theta = \log\left(\frac{p}{1-p}\right) \quad \text{for } 0 < p < 1$$

så kan frekvensfunktionen udtrykkes som

$$f_Y(y; \theta) = \exp\{n(\theta y - \log(1 + \exp(\theta)))\}$$

med

$$\kappa(\theta) = \log(1 + \exp(\theta))$$

og  $\lambda = n$ .

Altså to ækvivalente parametriseringer, ved  $p$  eller ved **kanonisk parameter**  $\theta$ .

## Middelværdiafbildning :



Funktionen

$$\tau(\theta) = \kappa'(\theta) \quad (4.11)$$

definerer en enetydig afbildning  $\mu = \tau(\theta)$  af parameterrummet,  $\Omega$ , for den kanoniske parameter,  $\theta$ , ind på en delmængde,  $\mathcal{M}$ , af den reelle akse.

Afbildningen kaldes **middelværdiafbildningen**.

Den omvendte afbildning,  $\theta = \tau^{-1}(\mu)$ , der fører middelværdien over i den kanoniske parameter, kaldes den **kanoniske linkfunktion**

## Kumulantfrembringer :



Funktionen  $\kappa(\cdot)$  beskriver hele fordelingen.

Således fås middelværdi og varians for fordelingen svarende til en bestemt værdi af  $\theta$

$$E[Y] = \kappa'(\theta) \quad (4.9)$$

$$V[Y] = \frac{\kappa''(\theta)}{\lambda}, \quad (4.10)$$

Naturligt at kalde  $\lambda$  for præcisionsparameter. Jo større  $\lambda$  desto mindre varians (større præcision)

## Variansfunktion :



(enheds) variansfunktionen udtrykker fordelings varians som en funktion

$$V(\mu) = \kappa''(\tau^{-1}(\mu)) \quad (4.13)$$

af middelværdien.

To ækvivalente parametriseringer af en eksponentiel dispersionsfamilie, nemlig

- ved **kanonisk parameter**,  $\theta$  og **kumulantfrembringer**,  $\kappa(\theta)$  med middelværdiafbildning mv.
- eller ved **middelværdiparameter**,  $\mu$  og **variansfunktion**,  $V(\mu)$ .



Vi indfører **enhedsdeviansen** ved

$$d(y; \mu) = 2 \int_{\mu}^y \frac{y-u}{V(u)} du \quad , \quad (4.14)$$

Tætheden for en observation kan udtrykkes ved enhedsdeviansen (dvs som funktion af  $\mu$ ) som

$$g_Y(y; \mu, \lambda) = a(y, \lambda) \exp \left\{ -\frac{\lambda}{2} d(y; \mu) \right\} \quad , \quad (4.16)$$

hvor  $d(y; \mu)$  er enhedsdeviansen og  $\lambda$  er præcisionsparameteren

[Vi kunne også have indført enhedsdevians som logaritmen til den relative likelihood af  $\mu$  (i forhold til maksimal likelihood)]

$$d(y; \mu) = -2[\ell(\mu; y) - \max_{\mu} \ell(\mu; y)]$$



observationsæt  $Y_1, Y_2, \dots, Y_n$ . Uafhængige samme eksponentielle dispersionsparameterfamilie, kumulantfrembringer  $\kappa(\cdot)$ . Præcisionsparameter for fordeling af  $Y_i$  givet som  $\lambda_i = w_i$  (kendt).

log-likelihood:

$$\ell_{\theta}(\boldsymbol{\theta}; \mathbf{y}) = \sum_{i=1}^k w_i (\theta_i y_i - \kappa(\theta_i)) \quad (4.21)$$

$$\ell_{\mu}(\boldsymbol{\mu}; \mathbf{y}) = -\frac{1}{2} \sum_{i=1}^k w_i d(y_i; \mu_i) = -\frac{1}{2} D(\mathbf{y}; \boldsymbol{\mu}) \quad , \quad (4.22)$$

hvor

$$D(\mathbf{y}; \boldsymbol{\mu}) = \sum_{i=1}^k w_i d(y_i; \mu_i) \quad . \quad (4.23)$$

### Enhedsdevianser svarende til forskellige fordelinger:



Family	$\mathcal{M}$	$V(\mu)$	unit deviance $d(y, \mu)$	$\lambda$	$\theta$
Normal	$(-\infty, \infty)$	1	$(y - \mu)^2$	$1/\sigma^2$	$\mu$
Poisson	$(0, \infty)$	$\mu$	$2\{y \ln(y/\mu) - (y - \mu)\}$	$-\alpha$	$\ln(\mu)$
Gamma	$(0, \infty)$	$\mu^2$	$2 \times \{y/\mu - \ln(y/\mu) - 1\}$	$\alpha^b$	$1/\mu$
Bin	$(0, 1)$	$\mu(1 - \mu)$	$2\{y \ln(y/\mu) + (1 - y) \ln((1 - y)/(1 - \mu))\}$	$n^c$	$\ln(\mu/(1 - \mu))$
Neg Bin	$(0, 1)$	$\mu(1 + \mu)$	$2\{y \ln(y(1 + \mu)/(1 + y\mu)) + \ln((1 + \mu)/(1 + y))\}$	$r^d$	$\ln(\mu)$
I Gauss	$(0, \infty)$	$\mu^3$	$(y - \mu)^2 / (y \times \mu^2)$		$1/\mu^2$
GHS <sup>e</sup>	$(-\infty, \infty)$	$1 + \mu^2$	$2y\{\arctan(y) - \arctan(\mu)\} + \ln((1 + \mu^2)/(1 + y^2))$	$1/\sigma^2$	$\arctan(\mu)$

### Scorefunktion :



I analogi med normalfordelingsituationen ser vi først på scorefunktion og observeret information svarende til fuld model (hver obs sin egen parameter)

$$\ell'_{\theta}(\boldsymbol{\theta}; \mathbf{y}) = \text{diag}(\mathbf{w})(\mathbf{y} - \boldsymbol{\tau}(\boldsymbol{\theta})) \quad (4.25)$$

$$\ell'_{\mu}(\boldsymbol{\mu}; \mathbf{y}) = \text{diag} \left\{ \frac{w_i}{V(\mu_i)} \right\} (\mathbf{y} - \boldsymbol{\mu}) \quad (4.27)$$

Er netop proportional med  $(\mathbf{y} - \boldsymbol{\mu})$

## Observeret information:



Observeret information

$$\mathbf{j}(\boldsymbol{\theta}; \mathbf{y}) = \text{diag}\{w_i V(\tau(\theta_i))\} \quad (4.26)$$

Afhænger af parameterværdien  $\boldsymbol{\theta}$

$$\mathbf{j}(\boldsymbol{\mu}; \mathbf{y}) = \text{diag} \left\{ w_i \left[ \frac{1}{V(\mu_i)} + (y_i - \mu_i) \frac{V'(\mu_i)}{V(\mu_i)^2} \right] \right\}, \quad (4.27)$$

Afhænger også af  $\mathbf{y}$

## Forskel fra den generelle lineære model (normalfordeling):



- fordelingen af observationerne beskrives ved en mere generel klasse af fordelinger, (eksponentielle dispersionsfamilier, hvis støtte ikke nødvendigvis er hele den reelle akse, som for normalfordelingen),
- og derfor måler man overensstemmelse mellem observation og parameterværdi på en måde, der afspejler egenskaberne ved den pågældende fordeling
- dels, at den lineære del ikke vedrører middelværdien,  $\boldsymbol{\mu}$ , men en funktion,  $\boldsymbol{\eta}$ , af middelværdien.
- I normalfordelingsmodeller skriver vi sædvanligvis obs = middelværdi plus støj, hvor støj  $\sim N(0, \sigma^2)$ . Det kan vi ikke for de andre fordelingsmodeller. For dem siger vi  $Y \sim \text{Ford}(\boldsymbol{\mu})$ , og så formulerer vi modeller for  $\boldsymbol{\mu}$  (eller rettere for en funktion  $\boldsymbol{\eta} = g(\boldsymbol{\mu})$ ).

## Generaliseret lineær model:



Ingredienser:

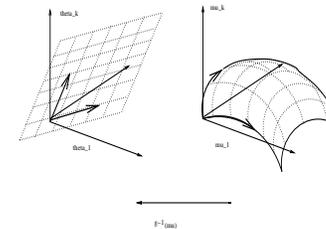
- observationsæt  $Y_1, Y_2, \dots, Y_k$
- eksponentiel dispersionsmodel med variansfunktion  $V(\boldsymbol{\mu})$
- **linkfunktion**  $\boldsymbol{\eta} = g(\boldsymbol{\mu})$  der fører middelværdi over i **lineær prædiktor** (koordinatvis)
- affin hypotese vedrørende  $\boldsymbol{\eta}_1, \boldsymbol{\eta}_2, \dots, \boldsymbol{\eta}_k$

$$H_0 : \boldsymbol{\eta} - \boldsymbol{\eta}_0 \in L, \quad (4.30)$$

dvs  $\boldsymbol{\eta} = \mathbf{X}\boldsymbol{\beta}$

Ideen er at vi laver en lineær model for en funktion (linkfunktionen) af middelværdien

## Middelværdirum og rum for lineær prædiktor:



Lineær prædiktor til venstre, middelværdirum til højre



Vi skal bare minimere **deviansen**

$$D(\mathbf{y}; \boldsymbol{\mu}) = \sum_{i=1}^k w_i d(y_i, \mu_i) .$$

under betingelsen  $\boldsymbol{\eta} = \mathbf{X}\boldsymbol{\beta}$

Scorefunktion mht  $\boldsymbol{\beta}$

$$l'_{\boldsymbol{\beta}}(\boldsymbol{\beta}; \mathbf{y}) = \frac{\partial}{\partial \boldsymbol{\beta}} l_{\boldsymbol{\beta}}(\boldsymbol{\beta}; \mathbf{y}) = [\mathbf{X}(\boldsymbol{\beta})]^T \mathbf{i}_{\boldsymbol{\mu}}(\boldsymbol{\mu}) (\mathbf{y} - \boldsymbol{\mu}(\boldsymbol{\beta})) ,$$

hvor

$$\mathbf{X}(\boldsymbol{\beta}) \stackrel{\text{DEF}}{=} \frac{\partial \boldsymbol{\mu}}{\partial \boldsymbol{\beta}} = \left[ \frac{d\boldsymbol{\mu}}{d\boldsymbol{\eta}} \right]^T \frac{\partial \boldsymbol{\eta}}{\partial \boldsymbol{\beta}} = \text{diag} \left\{ \frac{1}{g'(\mu_i)} \right\} \mathbf{X} \quad (4.35)$$

er den **lokale design matrix**

For kanonisk link bestemmes løsningen fra **Middelværdiligningen**:

$$\mathbf{X}^T \mathbf{y} = \mathbf{X}^T \boldsymbol{\mu}(\boldsymbol{\beta}) \quad (4.43)$$



Middelværdiligning

$$\begin{pmatrix} \sum_{i=1}^k n_i \exp(\alpha + \beta x_i) / \{1 + \exp(\alpha + \beta x_i)\} \\ \sum_{i=1}^k x_i n_i \exp(\alpha + \beta x_i) / \{1 + \exp(\alpha + \beta x_i)\} \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^k z_i \\ \sum_{i=1}^k x_i z_i \end{pmatrix}, \quad (4.53)$$



Højst én løsning pga konveksitet.

Man kan dog risikere at løsningen ligger på randen - og at randen ikke hører med til parameterområdet, nemlig feks en løsning  $\theta = -\infty$  svarende til  $p = 0$  i binomialfordelingssituation..



**law of error propagation**

For  $Z = g(Y)$  med  $E[Y] = \mu$  gælder

$$E[Z] \approx g(\mu), \quad V[Z] \approx (g'(\mu))^2 V[Y]$$

Resultatet fremkommer ved Taylorudvikling (dvs linearisering) af  $g(Y)$  omkring  $Y = \mu$

$$g(Y) = g(\mu) + (Y - \mu)g'(\mu) + o(Y - \mu)$$

Man kan bestemme maksimum-likelihood løsningen ved iterativt genvægtet mindste kvadraters metode.



## Hvad er fidusen ved en generaliseret lineær model ?

**Eksempel:** Estimation af hældning i regressionsmodeller

Data  $(x_1, y_1); (x_2, y_2); \dots; (x_n, y_n)$

Model:

$$E[Y] = \beta x$$

Individuelle estimater for hældningen:

$$\hat{\beta}_i = \frac{y_i}{x_i}$$



## Fordeling af estimatorer

$$\hat{\beta} - \beta \in N_m(\mathbf{0}, \Sigma), \quad (4.44)$$

$$\mathbf{D}[\hat{\beta}] = \Sigma = [\mathbf{X}^T \mathbf{W}(\beta) \mathbf{X}]^{-1} \quad (4.45)$$

med

$$\mathbf{W}(\beta) = \text{diag} \left\{ \frac{w_i}{[g'(\mu_i)]^2 V(\mu_i)} \right\}, \quad (4.46)$$

For kanonisk link bliver den bare

$$\mathbf{W}(\beta) = \text{diag} \{ w_i V(\mu_i) \} \quad (4.47)$$



## Hvad er fidusen ved en generaliseret lineær model ?

Normalfordeling:  $Y \in N(\beta x, \sigma^2)$ .  $V[Y] = \sigma^2$

$$\hat{\beta} = \frac{\sum y_i x_i}{\sum x_i^2} = \frac{\sum x_i^2 \hat{\beta}_i}{\sum x_i^2}$$

Poissonfordeling:  $Y \in P(\beta x)$ .  $V[Y] = \beta x$

$$\hat{\beta} = \frac{\sum y_i}{\sum x_i} = \frac{\sum x_i \hat{\beta}_i}{\sum x_i}$$

Gammafordeling:  $Y \in G(\alpha, \beta x / \alpha)$ .  $V[Y] = (\beta x)^2$

$$\hat{\beta} = \frac{\alpha \sum y_i / x_i}{\sum \alpha} = \frac{\sum \hat{\beta}_i}{n}$$

Eksempel 4.14.6

Vejet gennemsnit af individuelle hældninger, hvor de individuelle hældninger vægtes med deres **præcision**



## Fittede værdier, residualer

Fittede værdier af lineær prædikator,  $\hat{\eta}$  og af middelværdier  $\hat{\mu}$ .

Residualer:

Deviansresidual

$$r_i^D = r_D(y_i; \hat{\mu}_i) \stackrel{\text{def}}{=} \text{sign}(y_i - \hat{\mu}_i) \sqrt{w_i d(y_i, \hat{\mu}_i)} \quad (4.72)$$

Måler i likelihoodværdier.

Pearson-residual

$$r_i^P = r_P(y_i; \hat{\mu}_i) \stackrel{\text{def}}{=} \frac{y_i - \hat{\mu}_i}{\sqrt{V(\hat{\mu}_i) / w_i}} \quad (4.73)$$

Simpel, intuitiv.

Standardisering sikrer samme varians (ved division med  $(1 - h_{ii})$ )



## Test for modeltilpasning, goodness of fit

Sætning 4.8.1 (side 143):

Teststørrelse  $D(\mathbf{y}; \boldsymbol{\mu}(\hat{\boldsymbol{\beta}}))$ .

Under  $H_0: \boldsymbol{\eta} = \mathbf{X}\boldsymbol{\beta}$  vil  $D(\mathbf{y}; \boldsymbol{\mu}(\hat{\boldsymbol{\beta}}))$  approximativt følge en  $\chi^2(k - m)$ -fordeling. Forkast for store værdier.

Evt kan man bruge Pearson teststørrelse.

Bemærk, at i modsætning til normalfordelingsmodellen kan vi her bruge residualdeviansen til at teste modeltilpasning.



## Test for modelreduktion

Foregår ved opspaltning af deviansen:

$$D(\boldsymbol{\mu}(\hat{\boldsymbol{\beta}}); \hat{\boldsymbol{\mu}}_M) = D(\boldsymbol{\mu}(\hat{\boldsymbol{\beta}}); \boldsymbol{\mu}(\hat{\boldsymbol{\gamma}})) + D(\boldsymbol{\mu}(\hat{\boldsymbol{\gamma}}); \hat{\boldsymbol{\mu}}_M)$$

Under hypotesen  $\boldsymbol{\eta} = \mathbf{Z}\boldsymbol{\gamma}$  vil  $D(\boldsymbol{\mu}(\hat{\boldsymbol{\beta}}); \boldsymbol{\mu}(\hat{\boldsymbol{\gamma}}))$  approximativt følge en  $\chi^2(m - r)$ -fordeling.



## Analysis of deviance tabel

Variationskilde	$f$	Devians	middeldevians	Goodness of fit fortolkning
Model $H_0$	$m - 1$	$D(\boldsymbol{\mu}(\hat{\boldsymbol{\beta}}); \hat{\boldsymbol{\mu}}_M)$	$\frac{D(\boldsymbol{\mu}(\hat{\boldsymbol{\beta}}); \hat{\boldsymbol{\mu}}_M)}{m - 1}$	$G^2(H_M H_0)$
Residual (Error)	$k - m$	$D(\mathbf{y}; \boldsymbol{\mu}(\hat{\boldsymbol{\beta}}))$	$\frac{D(\mathbf{y}; \boldsymbol{\mu}(\hat{\boldsymbol{\beta}}))}{k - m}$	$G^2(H_0)$
Korrigeret Total	$k - 1$	$D(\mathbf{y}; \hat{\boldsymbol{\mu}}_M)$		$G^2(H_M)$



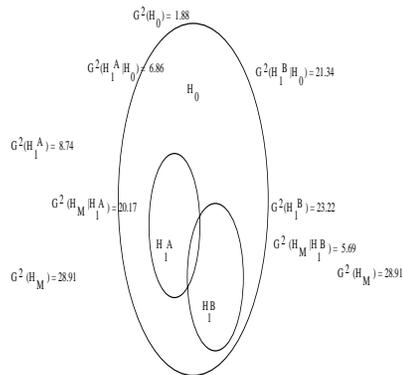
## Eksempel Blommeudplantning

Eksempel 4.11.1 side 152.

Antal overlevende/antal udplantede ved blommeudplantning

Længde	tykkelse		
	MID	TYK	TYN
KO	10/20	11/20	4/20
LA	14/20	18/20	6/20

Tofaktor forsøg med binomial respons



- Vi har udvidet til andre responsfordelinger, end normalfordelingen
- Disse familier kan beskrives ved kun én parameter, fx middelværdien. Variansen afhænger af middelværdien
- Middelværdien kan ikke variere på hele den reelle akse, derfor betragter vi en funktion af middelværdien, **linkfunktionen**, og laver lineære modeller for denne
- **Kanonisk link** fører middelværdien over i kanonisk parameter (grundform)
- Vi kan måle afvigelsen mellem en observation og en fittet middelværdi ved **deviansen**, og deviansen tager netop hensyn til at variansen afhænger af middelværdien
- og så er resten næsten lige som normalfordelingen.  
Vi estimerer ved at minimere devians. Devianser er approximativt  $\chi^2$ -fordelte, og vi kan **spalte devianser**.