

# TMA4315 Generalized linear models H2018

## Module 6: Categorical regression

Mette Langaas, Department of Mathematical Sciences, NTNU,  
with contributions from Ingeborg G. Hem

25.10.2017 [PL], 26.10.2017 [IL]

(Latest changes: 24.10: added IL + corrected typos, 23.10.2018:  
first version)

# Overview

## Learning material

This topic is *new* on the reading list this year.

- ▶ Textbook: Fahrmeir et al (2013): Chapter 6 (not p 344-345 nominal models and latent utility models, not 6.3.2 Sequential model, and not category specific variables on page 344-3458).
- ▶ Classnotes 25.10.2018

## Topics

- ▶ multinomial random component
- ▶ nominal vs. ordinal response
- ▶ ungrouped and grouped data
- ▶ multivariate exponential family
- ▶ nominal response and logit models
- ▶ ordinal response and logit models - based on a latent model
- ▶ likelihood inference

Jump to interactive.

## Categorical random component

We consider a situation where our random variable (response) is given as one of  $c + 1$  possible categories (where we will look at category  $c + 1$  as the reference category).

The categories will either be

- ▶ Unordered: *nominal response variable*. Example: food types in alligator example.
- ▶ Ordered: *ordered response variable*. Example: degrees of mental impairment.

## Assumptions:

- ▶ *Independent* observation pairs  $(\mathbf{Y}_i, \mathbf{x}_i)$ .
- ▶  $\pi_{ir}$ : probability that the response is category  $r$  for subject  $i$ .
- ▶  $\sum_{s=1}^{c+1} \pi_{is} = 1$  for all  $i$ , so that  $\pi_{i,c+1} = 1 - \sum_{s=1}^c \pi_{is}$ . So, we have  $c$  probabilities to estimate.
- ▶ Further, the covariate vector  $\mathbf{x}_i$  consists of the same measurements for each response category (that is, not different covariate types that are measured for each response category - which in our textbook is written as *independent of the response category*).

When coding the response variable we use a dummy variable coding with  $c$  elements (the  $c + 1$  category is the reference level). This means that if we have that  $\pi_{ir} = 1$  then  $\mathbf{y}_i = (0, 0, \dots, 0, 1, 0, \dots, 0)$  with a value of 1 in the  $r$ th element of  $\mathbf{y}_i$ . If observation  $i$  comes from category  $c + 1$  we have  $\mathbf{y}_i = (0, 0, \dots, 0)$ .

## Categorical regression

is modelling and estimating the probabilities

$\pi_{ir} = P(Y_i = r) = P(Y_{ir} = 1)$  as a function of the covariates  $\mathbf{x}_i$ .

The modelling is done differently for nominal (unordered) and ordered categories, but both rely upon the multinomial distribution. For unordered categories, a Poisson distribution can also be used



## The multinomial distribution

Probability mass function for one observation:

$$f(\mathbf{y}) = \pi_1^{y_1} \pi_2^{y_2} \cdots \pi_c^{y_c} (1 - \pi_1 - \pi_2 - \cdots - \pi_c)^{1 - y_1 - y_2 - \cdots - y_c}$$

where then  $\mathbf{y} = (y_1, y_2, \dots, y_c)$  and  $y_r = 1$  if the observation comes from the  $r$ th category.

If we then have  $m$  independent trials then  $\mathbf{y} = (y_1, y_2, \dots, y_c)$  is summed over our  $m$  responses, so that  $y_r$  is the number of observations where the response is from the  $r$ th category.

$$f(\mathbf{y}) = \frac{m!}{y_1! \cdots y_c! (m - y_1 - \cdots - y_c)!} \pi_1^{y_1} \pi_2^{y_2} \cdots \pi_c^{y_c} (1 - \pi_1 - \pi_2 - \cdots - \pi_c)^{m - y_1 - y_2 - \cdots - y_c}$$

The mean and the covariance matrix of the random vector  $\mathbf{Y}$  are given by:

$$\mathbf{E}(\mathbf{Y}) = m\boldsymbol{\pi} = \begin{pmatrix} m\pi_1 \\ m\pi_2 \\ \vdots \\ m\pi_c \end{pmatrix}$$

$$\text{Cov}(\mathbf{Y}) = m \begin{pmatrix} \pi_1(1 - \pi_1) & -\pi_1\pi_2 & \cdots & -\pi_1\pi_c \\ -\pi_2\pi_1 & \pi_2(1 - \pi_2) & \cdots & -\pi_2\pi_c \\ \vdots & \vdots & \ddots & \vdots \\ -\pi_c\pi_1 & -\pi_c\pi_2 & \cdots & \pi_c(1 - \pi_c) \end{pmatrix}$$

**Q:** what about  $\mathbf{E}(Y_{c+1})$  and  $\text{Cov}(Y_1, Y_{c+1})$ ?

Finally, if we look at  $\bar{Y}_r = \frac{1}{m}Y_r$  then  $\bar{\mathbf{Y}} = \frac{1}{m}\mathbf{Y}$  follows a scaled multinomial distribution  $\bar{\mathbf{Y}} \sim \frac{1}{m}M(m, \pi)$  with  $E(\bar{\mathbf{Y}}) = \pi$  and  $\text{Cov}(\bar{\mathbf{Y}}) = \frac{1}{m^2}\text{Cov}(\mathbf{Y})$ .

# Data

## Ungrouped data

$$\mathbf{Y} = \begin{pmatrix} Y_{11} & Y_{12} & \cdots & Y_{1c} \\ Y_{21} & Y_{22} & \cdots & Y_{2c} \\ \vdots & \vdots & \ddots & \vdots \\ Y_{n1} & Y_{n2} & \cdots & Y_{nc} \end{pmatrix}$$

and  $\mathbf{X}$  is an  $n \times p$  matrix as usual.

## Grouped data

As for the binomial case we look at the number of occurrences with a group - that is, one covariate pattern.

$$\mathbf{Y} = \begin{pmatrix} Y_{11} & Y_{12} & \cdots & Y_{1c} \\ Y_{21} & Y_{22} & \cdots & Y_{2c} \\ \vdots & \vdots & \ddots & \vdots \\ Y_{G1} & Y_{G2} & \cdots & Y_{Gc} \end{pmatrix}$$

The notation here is that we have  $n_i$  observation for each covariate pattern (group)  $i$  for  $i = 1, \dots, G$ . This will replace the  $m$  used for the multinomial distribution above.

## Regression with nominal responses

nominal=unordered

Agresti (2015, p203): “The model treats the response variable as nominal scale in the following sense: if the model holds and the outcome categories are permuted in any way, the model still holds with the corresponding permutation of the effects.”

This is a generalization of the binary logit model with  $P(Y = 1)$  vs  $P(Y = 0)$ , to  $c$  models of  $\pi_{ir}$  vs  $\pi_{i,c+1}$  for  $r = 1, \dots, c$ .

The models can be written using log ratios:

$$\ln\left(\frac{\pi_{ir}}{\pi_{i,c+1}}\right) = \mathbf{x}_i^T \beta_r$$

Remark:  $\beta_r$  is the  $p \times 1$  coefficient vector for the  $r$ th response

Using this we may also look at the log ratio for any two probabilities  $\pi_{ia}$  and  $\pi_{ib}$ :

$$\ln\left(\frac{\pi_{ia}}{\pi_{ib}}\right) = \ln\left(\frac{\pi_{ia}}{\pi_{i,c+1}}\right) - \ln\left(\frac{\pi_{ib}}{\pi_{i,c+1}}\right) = \mathbf{x}_i^T (\beta_a - \beta_b)$$



Alternatively, we may write out the model for the probabilities:

$$P(Y_i = r) = \pi_{ir} = \frac{\exp(\mathbf{x}_i^T \beta_r)}{1 + \sum_{s=1}^c \exp(\mathbf{x}_i^T \beta_s)}$$

$$P(Y_i = c+1) = \pi_{i,c+1} = 1 - \pi_{i1} - \dots - \pi_{ic} = \frac{1}{1 + \sum_{s=1}^c \exp(\mathbf{x}_i^T \beta_s)}$$

## Multivariate GLM

This is a multivariate GLM and the multinomial distribution is a *multivariate exponential family*.

$$f(\mathbf{y}_i, \boldsymbol{\theta}_i, \phi) = \exp\left(\frac{\mathbf{y}_i^T \boldsymbol{\theta}_i - b(\boldsymbol{\theta}_i)w_i}{\phi} + c(\mathbf{y}_i, \phi, w_i)\right)$$

where  $\boldsymbol{\theta}$  has dimension  $c$ .

## Multivariate GLM-set-up

1.  $\mathbf{Y}_i$  is multinomial with

$$\mu_i = \mathbf{E}(\mathbf{Y}_i) = \pi_i = \begin{pmatrix} \pi_{i1} \\ \pi_{i2} \\ \vdots \\ \pi_{i,c+1} \end{pmatrix}$$

Remark: if grouped data we instead look at  $\bar{\mathbf{Y}}_i \sim \frac{1}{n_i} M(n_i, \pi_i)$  so that the mean is  $\pi_i$

2. Linear predictor is now a  $c \times 1$  vector:

$$\eta_i = \begin{pmatrix} \eta_{i1} \\ \eta_{i2} \\ \vdots \\ \eta_{i,c} \end{pmatrix} = \begin{pmatrix} \mathbf{x}_i^T \beta_1 \\ \mathbf{x}_i^T \beta_2 \\ \vdots \\ \mathbf{x}_i^T \beta_c \end{pmatrix}$$

3. Link functions ( $c$  of those):  $\mathbf{g}(\boldsymbol{\mu}_i) = \boldsymbol{\eta}_i$  where for the *nominal logit data model* element  $r$  (for  $r = 1, \dots, c$ ) of  $\mathbf{g}$  is

$$g_r(\boldsymbol{\mu}_i) = \ln\left(\frac{\mu_{ir}}{1 - \mu_{i1} - \dots - \mu_{ic}}\right) = \ln\left(\frac{\pi_{ir}}{1 - \pi_{i1} - \dots - \pi_{ic}}\right)$$

We also define response functions ( $\mathbf{h}$ ) with elements  $h_r$  given by  $\pi_{ir} = h_r(\eta_{i1}, \eta_{i2}, \dots, \eta_{ic})$ , and we have for the *nominal data model*

$$\pi_{ir} = h_r((\eta_{i1}, \eta_{i2}, \dots, \eta_{ic})) = \frac{\exp(\eta_{ir})}{1 + \sum_{s=1}^c \exp(\eta_{is})}$$

It turns out that the reference category logits are the canonical links for the multinomial distribution GLM.

In this case, as for the univariate exponential family GLM the loglikelihood is concave with an unique maximum (if it exists) and the expected and observed Fisher information matrices are equal.

As before, we find maximum likelihood parameter estimates from the Fisher scoring or Newton Raphson method.

Remember: now we have  $p \times c$  parameters to estimate —  $p$  for each category  $c$ . All of these coefficients may either be put into a long vector (length  $p \cdot c$ ) — which might be easiest to understand for the estimation, or into a matrix of dimension  $p \times c$  — might be easier for viewing.

## Likelihood

(grouped data)

With the notation that  $\beta$  is a long vector with the coefficients for the  $c$  categories stacked upon each other.

$$L(\beta) = \prod_{i=1}^G f(\mathbf{y}_i | \pi)$$

where  $f$  is the multinomial distribution function.

## Loglikelihood

$$l(\beta) \propto \sum_{i=1}^G \sum_{s=1}^{c+1} y_{is} \ln(\pi_{is})$$

where we remember that  $y_{i,c+1} = n_i - y_{i1} - \dots - y_{ic}$ , and  $1 - \pi_{i1} - \dots - \pi_{ic}$ .

(This formula is also correct for the ordinal model of the next section.) General formulas for the score function and expected Fisher information matrix follow later.



## Deviance

The derivation used for the binary GLM model generalizes directly to the multinomial GLM. The fitted probabilities are  $\hat{\pi}_{ij}$  (group  $i$  and category  $j$ ) and the saturated model (grouped data) is

$$n_i \tilde{\pi}_{ij} = y_{ij}.$$

$$D = 2 \sum_{i=1}^G \sum_{s=1}^{c+1} y_{is} \ln\left(\frac{y_{is}}{n_i \hat{\pi}_{is}}\right)$$

The asymptotic distribution is as before  $\chi^2$  with “the number of groups times number of categories minus 1 ( $Gc$ )” minus “the number of covariates ( $cp$ )”, giving  $Gc - cp = c(G - p)$  degrees of freedom.

The deviance can be used for model check with grouped data ( $G$  groups with  $n_i$  observations), but can be used to compare nested unsaturated models also for individual (ungrouped) data, with again an asymptotic  $\chi^2$  distribution with the difference of number of parameters between the two models.

This formula is also correct for the ordinal model of the next section, except that the number of parameters estimated differ.

## Alligators example

Example and data are taken from Agresti (2015, pages 217-219).  
Research question: what is the factors influencing the primary food choice of alligators?

Data are from 219 captured alligators from four lakes in Florida, where the stomach contents of the alligators were investigated. The weight of different types of food was measured, and then the primary food choice (highest weight) was noted. The primary choice is given as  $y_1:y_5$  below. In addition the size of the alligator (non-adult or adult) was registered.

- ▶ lake: each of the 4 lakes in Florida (1:4)
- ▶ size: non-adult=the size of the alligator (0: 2.3 meters or smaller) and adult=(1: larger than 2.3 meters)
- ▶ y1: fish
- ▶ y2: invertebrate
- ▶ y3: reptile
- ▶ y4: bird
- ▶ y5: other

These data are grouped, and we let y1:fish be the reference category.

```
# data from Agresti (2015), section 6, with use of the VGAM pack  
data = "http://www.stat.ufl.edu/~aa/glm/data/Alligators.dat"  
ali = read.table(data, header = T)  
ali  
attach(ali)
```

```
##   lake size y1 y2 y3 y4 y5  
## 1    1    1 23  4  2  2  8  
## 2    1    0  7  0  1  3  5  
## 3    2    1  5 11  1  0  3  
## 4    2    0 13  8  6  1  0  
## 5    3    1  5 11  2  1  5  
## 6    3    0  8  7  6  3  5  
## 7    4    1 16 19  1  2  3  
## 8    4    0 17  1  0  1  3
```

```
y.data = cbind(y2, y3, y4, y5, y1)
y.data
dim(y.data)
x.data = model.matrix(~size + factor(lake), data = ali)
x.data
dim(x.data)
```

```
##      y2 y3 y4 y5 y1
## [1,]  4  2  2  8 23
## [2,]  0  1  3  5  7
## [3,] 11  1  0  3  5
## [4,]  8  6  1  0 13
## [5,] 11  2  1  5  5
## [6,]  7  6  3  5  8
## [7,] 19  1  2  3 16
## [8,]  1  0  1  3 17
## [1] 8 5
## (Intercept) size factor(lake)2 factor(lake)3 factor(la
## 1           1     1           0           0
## 2           1     0           0           0
```

```
# We use library VGAM:
```

```
library(VGAM)
```

```
# We fit a multinomial logit model with fish (y1) as the r
```

```
fit.main = vglm(cbind(y2, y3, y4, y5, y1) ~ size + factor(1
```

```
  data = ali)
```

```
# summary(fit.main)
```

```
pchisq(deviance(fit.main), df.residual(fit.main), lower.tail
```

```
## [1] 0.1466189
```

Q:

- ▶ Why is the number of degrees of freedom for the residual deviance 12? Hint: there are 8 covariate patterns, and we have 5 response categories.
- ▶ How can you interpret the coefficient for invertebrate (y2) and size? Hint: we have y2,y3,y4,y5 as 1:4.

```
exp(coefficients(fit.main))
```

```
## (Intercept):1 (Intercept):2 (Intercept):3 (Intercept):4
## 0.0404626 0.1259644 0.2471007 0.3402404
## size:2 size:3 size:4 factor(lake)2
## 0.7037987 0.5322405 1.3931262 13.4043304
## factor(lake)2:3 factor(lake)2:4 factor(lake)3:1 factor(lake)3:2
## 0.2596748 0.4401925 16.1245576 5.4329104
## factor(lake)3:4 factor(lake)4:1 factor(lake)4:2 factor(lake)4:3
## 1.9940595 5.2506853 0.2885818 0.4990104
```



Testing out other models, and comparing with LRT-test - by using deviances for different models.

```
# Fit model with only lake:
```

```
fit.lake = vglm(cbind(y2, y3, y4, y5, y1) ~ factor(lake), f
```

```
# Test effect of size (no anova command is available)
```

```
(G2 = deviance(fit.lake) - deviance(fit.main))
```

```
(df.diff = df.residual(fit.lake) - df.residual(fit.main))
```

```
1 - pchisq(G2, df.diff)
```

```
# Size has a significant effect
```

```
## [1] 21.08741
```

```
## [1] 4
```

```
## [1] 0.0003042796
```

```
# Fit model with only size:
fit.size = vglm(cbind(y2, y3, y4, y5, y1) ~ size, family =
# Test effect of lake
(G2 = deviance(fit.size) - deviance(fit.main))
(df.diff = df.residual(fit.size) - df.residual(fit.main))
1 - pchisq(G2, df.diff)
# Lake has a significant effect
```

```
## [1] 49.13308
## [1] 12
## [1] 1.982524e-06
```

Q: explain what is presented below, in particular “what is the probability that the main food source is fish given size=0 and lake=1”?

```
library(knitr)
# Fitted values for main effect model 'fit.main':
fitted = data.frame(fitted(fit.main), lake = ali$lake, size = ali$size)
kable(fitted)
```

y2	y3	y4	y5	y1	lake	size
0.0930988	0.0474566	0.0704015	0.2537396	0.5353035	1	1
0.0230717	0.0718246	0.1408963	0.1940096	0.5701978	1	0
0.6018967	0.0772276	0.0088175	0.0538721	0.2581861	2	1
0.2486452	0.1948374	0.0294161	0.0686628	0.4584385	2	0
0.5168385	0.0887672	0.0358947	0.1742005	0.1842990	3	1
0.1929612	0.2023995	0.1082251	0.2006616	0.2957525	3	0
0.4128558	0.0115665	0.0296712	0.0938024	0.4521040	4	1
0.1396778	0.0238987	0.0810674	0.0979136	0.6574425	4	0

## Regression with ordinal responses

(we will only consider cumulative models - and not sequential models)

An unobservable latent variable  $U_i$  drives the observed category  $Y_i$ .

$$Y_i = r \Leftrightarrow \theta_{r-1} \leq U_i \leq \theta_r$$

where these  $\theta$ s are our unobservable thresholds, and the thresholds are monotonely increasing,  $-\infty = \theta_0 < \theta_1 < \dots < \theta_{c+1} = \infty$ .

We further assume that the latent variables are dependent on our covariates through

$$U_i = -\mathbf{x}_i^T \beta + \varepsilon_i$$

where we have a new random variable that has cumulative distribution function (cdf)  $F$ . No intercept is included due to identifiability issue (shift in intercept would produce the same effect as negative shift in threshold).

We get rid of the latent variable  $U_i$  by considering

$$\begin{aligned} P(Y_i \leq r) &= P(U_i \leq \theta_r) = P(-\mathbf{x}_i^T \beta + \varepsilon_i \leq \theta_r) \\ &= P(\varepsilon_i \leq \theta_r + \mathbf{x}_i^T \beta) = F(\theta_r + \mathbf{x}_i^T \beta) \end{aligned}$$

Observe that the final expression does not include the latent variable  $U_i$ , but includes the unknown threshold and  $k$  regression parameters.

Different choices of  $F$  will give different models, and we will only consider  $F$  to be the cdf for the logistic distribution. (Another popular choice is the cdf of the standard normal distribution.)

$$P(Y_i \leq r) = \frac{\exp(\theta_r + \mathbf{x}_i^T \beta)}{1 + \exp(\theta_r + \mathbf{x}_i^T \beta)}$$

which also can be written as

$$\ln\left(\frac{P(Y_i \leq r)}{P(Y_i > r)}\right) = \theta_r + \mathbf{x}_i^T \beta$$

Our model is a proportional odds model, in the sense that the cumulative odds are proportional across categories

$$\frac{\frac{P(Y \leq r | \mathbf{x}_i)}{P(Y > r | \mathbf{x}_i)}}{\frac{P(Y \leq r | \mathbf{x}_i^*)}{P(Y > r | \mathbf{x}_i^*)}} = \exp((\mathbf{x}_i - \mathbf{x}_i^*)^T \beta)$$

Observe that this is independent of  $r$ .

## Response function

What is the response function here?

$$\pi_{i1} = F(\eta_{i1})$$

$$\pi_{ir} = F(\eta_{ir}) - F(\eta_{i,r-1})$$

where  $\eta_{ir} = \theta_r + \mathbf{x}_i^T \beta$ , and  $F$  is the logistic cdf.



## Plotting the cumulative probabilities

We have five categories, where the fifth is the reference category.

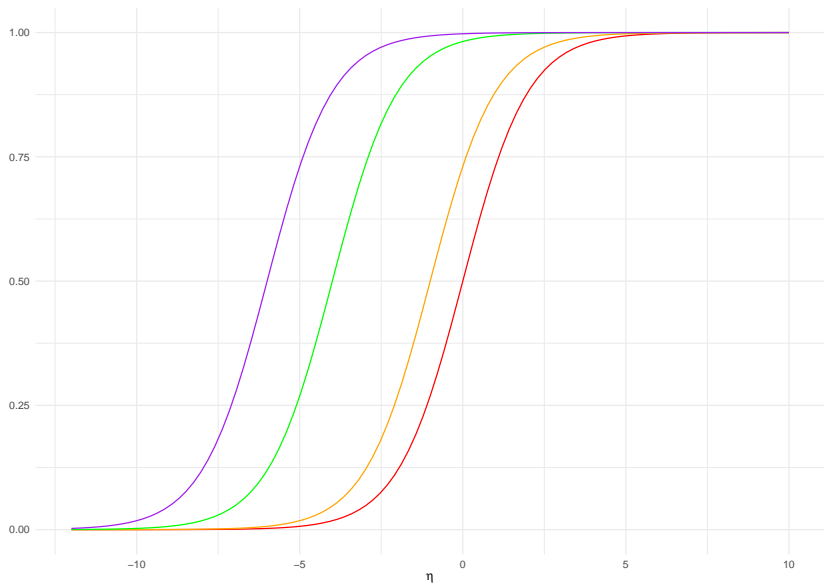
True parameters

- ▶  $\theta_1 = 0$ ,  $\theta_2 = 1$ ,  $\theta_3 = 4$  and  $\theta_4 = 6$  and
- ▶ one covariate with parameter  $\beta = 1$ .

The graph shows the cumulative probability  $P(Y \leq r)$  for  $r = 1$  (red),  $r = 2$  (orange),  $r = 3$  (green),  $r = 4$  (purple).

Observe the parallel lines.

What would  $P(Y \leq 5)$  be? Why is this missing from the plot?



## Mental health data example

Example and data are taken from Agresti (2015, pages 219-223).

Research question: understand mental health issues.

The data comes from a random sample of size 40 of adult residents of Alachua County, Florida, USA.

- ▶ Mental impairment  $Y$ : 1=well, 2=mild symptom formation, 3=moderate symptom formation, 4=impaired.
- ▶ Life event index ( $x_1$ ): composite measure of the number and severity of important life events within the last three years (birth, new job, divorce, death in the family, ...)
- ▶ SES ( $x_2$ ): socioeconomic index, 1=high, 0=low.

These data are ungrouped (but could be grouped). In the original study several other explanatory variables were studied.

```
# Read mental health data from the web:  
library(knitr)  
data = "http://www.stat.ufl.edu/~aa/glm/data/Mental.dat"  
mental = read.table(data, header = T)  
colnames(mental)  
apply(mental, 2, table)  
# kable(mental)
```

```
## [1] "impair" "ses"      "life"  
## $impair  
##  
##  1  2  3  4  
## 12 12  7  9  
##  
## $ses  
##  
##  0  1  
## 18 22  
##  
## $life  
##
```

```
library(VGAM)
```

```
# We fit a cumulative logit model with main effects of 'ses'  
fit.imp = vglm(impair ~ life + ses, family = cumulative(par  
# parallel=T gives proportional odds structure - only into  
summary(fit.imp)
```

```
##
```

```
## Call:
```

```
## vglm(formula = impair ~ life + ses, family = cumulative  
##       data = mental)
```

```
##
```

```
## Coefficients:
```

```
##           Estimate Std. Error z value Pr(>|z|)  
## (Intercept):1 -0.2819    0.6231  -0.452  0.65096  
## (Intercept):2  1.2128    0.6511   1.863  0.06251 .  
## (Intercept):3  2.2094    0.7171   3.081  0.00206 **  
## life          -0.3189    0.1194  -2.670  0.00759 **  
## ses           1.1112    0.6143   1.809  0.07045 .
```

```
## ---
```

```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1
```

The ML fit for this model can be written as

$$\text{logit}(\hat{P}(y_i \leq r)) = \hat{\theta}_r + 0.319x_{i1} + 1.111x_{i2}$$

**Q:** give an interpretation of this model!

Remember:

- ▶ Life event index ( $x_1$ ): composite measure of the number and severity of important life events within the last three years (birth, new job, divorce, death in the family, ...)
- ▶ SES ( $x_2$ ): socioeconomic index, 1=high, 0=low.

**Q:** How can you interpret the last line below? Why is it  $\exp(\text{CI}(\beta))$  and not  $\text{CI}(\exp(\beta))$ ?

```
exp(confint(fit.imp))
```

```
##                2.5 %      97.5 %  
## (Intercept):1 0.2224503  2.5581328  
## (Intercept):2 0.9385968 12.0489557  
## (Intercept):3 2.2342109 37.1467162  
## life          0.5752574  0.9187045  
## ses           0.9114465 10.1266209
```

**Q:** How are these predictions calculated? What is the interpretation?

```
fitted = data.frame(fitted(fit.imp), ses = mental$ses, life = mental$life)
fitted[c(6, 18, 10), ] #0,7 not fitted

xs = cbind(c(2, 7, 2, 7), c(0, 0, 1, 1))
coeff = coefficients(fit.imp)
linpreds = cbind(coeff[1] + xs %*% coeff[4:5], coeff[2] + xs %*% coeff[4:5])
(cprobs = exp(linpreds)/(1 + exp(linpreds)))
(pprobs = cbind(cprobs[, 1], cprobs[, 2] - cprobs[, 1], cprobs[, 3] - cprobs[, 1],
  1 - cprobs[, 3]))
```

```
##           X1           X2           X3           X4 ses life
## 6  0.2850362 0.3548973 0.18808559 0.17198084    0    2
## 18 0.5477558 0.2959810 0.09227184 0.06399141    1    2
## 10 0.1973858 0.3255923 0.22513134 0.25189056    1    7
##           [,1]           [,2]           [,3]
## [1,] 0.28503623 0.6399336 0.8280192
## [2,] 0.07488691 0.2651744 0.4943332
```



Q: What do you see here, and what is the formula for this matrix?

```
vcov(fit.imp)
```

```
##                (Intercept):1 (Intercept):2 (Intercept):3
## (Intercept):1    0.38819709    0.32992954    0.32615019
## (Intercept):2    0.32992954    0.42395851    0.40844529
## (Intercept):3    0.32615019    0.40844529    0.51423495
## life              -0.04231112   -0.05427393   -0.06185757
## ses                -0.15615761   -0.11440293   -0.09055667
```

```
# We consider a model with interaction between 'ses' and 'life'  
fit.int = vglm(impair ~ life + ses + life:ses, family = cur  
  data = mental)  
summary(fit.int)
```

```
##
```

```
## Call:
```

```
## vglm(formula = impair ~ life + ses + life:ses, family =  
##   data = mental)
```

```
##
```

```
## Coefficients:
```

```
##           Estimate Std. Error z value Pr(>|z|)  
## (Intercept):1  0.09807    0.81102   0.121  0.90375  
## (Intercept):2  1.59248    0.83717   1.902  0.05714 .  
## (Intercept):3  2.60660    0.90966   2.865  0.00416 **  
## life           -0.42045    0.19031  -2.209  0.02715 *  
## ses            0.37090    1.13022   0.328  0.74279  
## life:ses       0.18131    0.23611   0.768  0.44255
```

```
## ---
```

```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1
```

```
# And test if there is a significant effect of interaction
G2 = deviance(fit.imp) - deviance(fit.int)
df.diff = df.residual(fit.imp) - df.residual(fit.int)
1 - pchisq(G2, df.diff)
# The effect of interaction is not significant
```

```
## [1] 0.4410848
```

```
# We consider a model where the effect of the covariates mo
# cumulative logits - so not parallell lines for the cdfs
fit.nopar = vglm(impair ~ life + ses, family = cumulative,
summary(fit.nopar)
```

```
# The change in the deviance compared to the model 'fit.imp
# 99.0979-96.7486=2.3493 with df.diff=115-111=4, which is 1

# So model 'fit.imp' seems fine.
```

```
##
```

```
## Call:
```

```
## vglm(formula = impair ~ life + ses, family = cumulative,
```

## Why not use MLR instead of ordinal regression?

Based on Agresti (2015, p 214-216)

To use MLR the ordinal categories need to be replaced with numerical values, and we then need to assume a normal error structure. The following are questions to be answered and possible limitation to be assumed for using MLR instead of ordinal regression:

- ▶ how to translate ordered categories into numerical scores?
- ▶ is it better with an ordinal variable with some range than a single numerical number?
- ▶ MLR will not give probabilities for each response category
- ▶ variability in the response may be dependent on the category, for MLR we assume homoscedasticity

## Likelihood inference

We use the notation that  $\beta$  is a long vector with all regression parameters. The content of this vector is slightly different for our two models, with intercept and  $k$  covariate effects for each response category for the nominal model - and with  $c$  thresholds but the same  $k$ -dimensional  $\beta$  vector for all categories.

Full matrix versions (over all  $i$ ) can be found in our textbook, page 345-346.

### Loglikelihood

We have seen that the loglikelihood is:

$$l(\beta) \propto \sum_{i=1}^n \sum_{s=1}^{c+1} y_{is} \ln(\pi_{is})$$

where we remember that  $y_{i,c+1} = n_i - y_{i1} - \dots - y_{ic}$ , and  $1 - \pi_{i1} - \dots - \pi_{ic}$ .

## Design matrix and coefficient vector

The design matrix  $\mathbf{X}$  and coefficient vector are different for our nominal logit model and our ordinal cumulative model.

### Nominal logit model

$$\mathbf{X}_i = \text{diag}(\mathbf{x}_i^T) = \begin{pmatrix} \mathbf{x}_i^T & 0 & \cdots & 0 \\ 0 & \mathbf{x}_i^T & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \mathbf{x}_i^T \end{pmatrix}$$

where the 0s are  $1 \times p$  vectors. The dimension of the design matrix for covariate pattern  $i$  is  $c \times c \cdot p$ .

The vector of coefficients has dimension  $c \cdot p \times 1$ .

$$\beta = \begin{pmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_c \end{pmatrix}$$

## Ordinal cumulative model

$$\mathbf{X}_i = \begin{pmatrix} 1 & 0 & \cdots & 0 & \mathbf{x}_i^T \\ 0 & 1 & \cdots & 0 & \mathbf{x}_i^T \\ \vdots & \vdots & \ddots & \vdots & \\ 0 & 0 & \cdots & 1 & \mathbf{x}_i^T \end{pmatrix}$$

The dimension of the design matrix for covariate pattern  $i$  is  $c \times (c + k)$

The vector of coefficients has dimension  $(c + k) \times 1$  (where  $p = k + 1$ ), and now the thresholds replace the intercept and are put first in the vector, and the effects of the covariates are the same for all categories.

$$\beta = \begin{pmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_c \\ \beta \end{pmatrix}$$



## Score function

$$\mathbf{s}(\boldsymbol{\beta}) = \sum_{i=1}^G \mathbf{X}_i^T \mathbf{D}_i \boldsymbol{\Sigma}_i^{-1} (\mathbf{y}_i - n_i \boldsymbol{\pi}_i)$$

where

- ▶  $\mathbf{D}_i = \left. \frac{\partial h(\boldsymbol{\eta})}{\partial \boldsymbol{\eta}} \right|_{\boldsymbol{\eta}=\boldsymbol{\eta}_i}$  has dimension  $c \times c$
- ▶  $\boldsymbol{\Sigma}_i = \text{Cov}(\mathbf{Y}_i)$

## Fisher information

The dimension of the matrix is  $cp \times cp$  for the nominal case and  $(c + k) \times (c + k)$  for the ordinal case studied.

$$F(\beta) = \sum_{i=1}^G \mathbf{X}_i^T \mathbf{W}_i \mathbf{X}_i$$

where  $\mathbf{W}_i$  is given as  $\mathbf{D}_i \Sigma_i^{-1} \mathbf{D}_i^T$ .

## Finding the ML estimate

As in modules 1-5 we find the ML estimate by the Fisher scoring or Newton Raphson method.

## Asymptotic distribution

As in modules 1-5 the ML estimator  $\hat{\beta}$  asymptotically follows a multivariate normal distribution with unbiased mean and covariance matrix given by the inverse of the expected Fisher information matrix.

Summing up

## R packages

```
install.packages(c("VGAM", "ggplot2", "statmod", "knitr"))
```

## Further reading

- ▶ A. Agresti (2015): “Foundations of Linear and Generalized Linear Models.” Chapter 6. Wiley.