# Reduced-bias estimation of models with ordinal responses

Ioannis Kosmidis

🐦 IKosmidis_

✉ ioannis.kosmidis@warwick.ac.uk

🌐 http://www.ikosmidis.com

Reader in Data Science

University of Warwick & The Alan Turing Institute

26 February 2020
NISOx group meetings
Big Data Institute, University of Oxford

# Outline

# Outline

# Wine tasting data[1]

| contact | temp | rating | | | | |
|---------|------|---|---|---|---|---|
| | | 1 | 2 | 3 | 4 | 5 |
| no | cold | 4 | 9 | 5 | 0 | 0 |
| | warm | 0 | 5 | 8 | 3 | 2 |
| yes | cold | 1 | 7 | 8 | 2 | 0 |
| | warm | 0 | 1 | 5 | 7 | 5 |



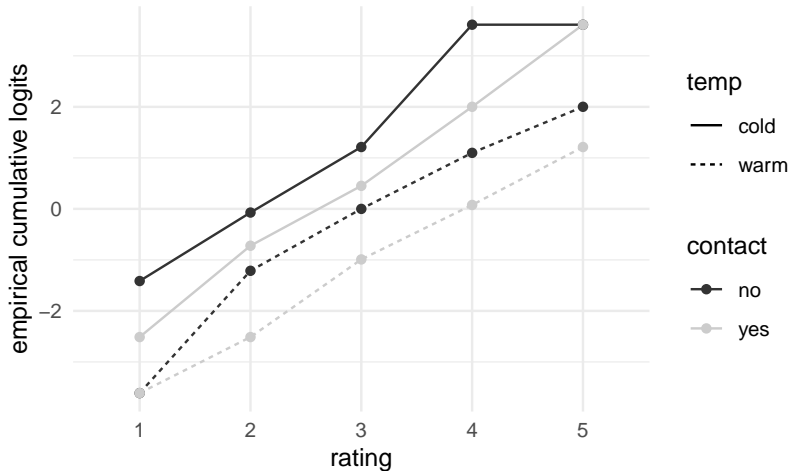Experiment on the effect of factors on the bitterness of white wine

contact of juice with skin and temperature when crushing the grapes

9 judges rated 2 bottles per combination of factors in terms of bitterness

---

[1]data from Randall (1989)

Empirical cumulative logits for factor combination $i$ and rating $j$

$$\log \frac{Y_{i1} + \ldots + Y_{ij} + 0.5}{Y_{ij+1} + \ldots + Y_{ik} + 0.5}$$

# Testing for proportional odds

Assume that counts for the $i$th factor combination are from independent

$$(Y_{i1}, \ldots, Y_{i5}) \sim \text{Mult}(18, (\pi_{i1}, \ldots, \pi_{i5}))$$

**Proportional odds model**[2]

$$\log \frac{\pi_{i1} + \ldots + \pi_{ij}}{\pi_{ij+1} + \ldots + \pi_{i5}} = \alpha_j - \beta w_i - \delta z_i$$

where $w_i$ is 0 (cold) or 1 (warm), $z_i$ is 0 (no) or 1 (yes),
$\beta, \delta \in \Re$, $\alpha_1 < \ldots < \alpha_4 < \alpha_5 = \infty$

---

[2]see, McCullagh (1980)
[3]see, Peterson and Harrell (1990)

# Testing for proportional odds

Assume that counts for the $i$th factor combination are from independent

$$(Y_{i1}, \ldots, Y_{i5}) \sim \text{Mult}(18, (\pi_{i1}, \ldots, \pi_{i5}))$$

**Proportional odds model**[2]

$$\log \frac{\pi_{i1} + \ldots + \pi_{ij}}{\pi_{ij+1} + \ldots + \pi_{i5}} = \alpha_j - \beta w_i - \delta z_i$$

where $w_i$ is 0 (cold) or 1 (warm), $z_i$ is 0 (no) or 1 (yes), $\beta, \delta \in \Re$, $\alpha_1 < \ldots < \alpha_4 < \alpha_5 = \infty$

**Partial proportional odds model**[3]

$$\log \frac{\pi_{i1} + \ldots + \pi_{ij}}{\pi_{ij+1} + \ldots + \pi_{i5}} = \alpha_j - \gamma_j w_i - \delta z_i$$

---

[2]see, McCullagh (1980)
[3]see, Peterson and Harrell (1990)

# Testing for proportional odds

Assume that counts for the $i$th factor combination are from independent

$$(Y_{i1}, \ldots, Y_{i5}) \sim \text{Mult}(18, (\pi_{i1}, \ldots, \pi_{i5}))$$

**Proportional odds model**[2]

$$\log \frac{\pi_{i1} + \ldots + \pi_{ij}}{\pi_{ij+1} + \ldots + \pi_{i5}} = \alpha_j - \beta w_i - \delta z_i$$

where $w_i$ is 0 (cold) or 1 (warm), $z_i$ is 0 (no) or 1 (yes),
$\beta, \delta \in \Re$, $\alpha_1 < \ldots < \alpha_4 < \alpha_5 = \infty$

**Partial proportional odds model**[3]

$$\log \frac{\pi_{i1} + \ldots + \pi_{ij}}{\pi_{ij+1} + \ldots + \pi_{i5}} = \alpha_j - \gamma_j w_i - \delta z_i$$

Proportional odds hypothesis $\gamma_1 = \gamma_2 = \gamma_3 = \gamma_4 = \beta$

---

[2]see, McCullagh (1980)
[3]see, Peterson and Harrell (1990)

# Testing for proportional odds

Use Wald statistic

$$(L\hat{\gamma})^\top \left\{ L F^{\gamma\gamma}(\hat{\theta}) L^\top \right\}^{-1} L\hat{\gamma}$$

with a $\chi_3^2$ limiting distribution under proportional odds

$F^{\gamma\gamma}(\theta)$ is $\gamma$-block of the inverse Fisher information matrix

$L$ is a matrix of $\gamma$-contrasts $\begin{bmatrix} 1 & . & . & -1 \\ . & 1 & . & -1 \\ . & . & 1 & -1 \end{bmatrix}$

---

[5]see, Pratt (1981) and Agresti (2010, §3.4.5) for sufficient conditions

# Testing for proportional odds

Use Wald statistic
$$(L\hat{\gamma})^{\top} \left\{ L F^{\gamma\gamma}(\hat{\theta}) L^{\top} \right\}^{-1} L\hat{\gamma}$$

with a $\chi_3^2$ limiting distribution under proportional odds

$F^{\gamma\gamma}(\theta)$ is $\gamma$-block of the inverse Fisher information matrix

$L$ is a matrix of $\gamma$-contrasts $\begin{bmatrix} 1 & . & . & -1 \\ . & 1 & . & -1 \\ . & . & 1 & -1 \end{bmatrix}$

Maximum likelihood[4] returns infinite estimates[5]

| $\alpha_1$ | $\alpha_2$ | $\alpha_3$ | $\alpha_4$ | $\gamma_1$ | $\gamma_2$ | $\gamma_3$ | $\gamma_4$ | $\delta$ |
|---|---|---|---|---|---|---|---|---|
| -1.27 | 1.10 | 3.77 | 24.90 | 21.10 | 2.15 | 2.87 | 22.55 | 1.47 |
| | | | Maximum absolute log-likelihood gradient: $10^{-6}$ | | | | | |
| -1.27 | 1.10 | 3.77 | 33.89 | 30.10 | 2.15 | 2.87 | 31.55 | 1.47 |
| | | | Maximum absolute log-likelihood gradient: $10^{-10}$ | | | | | |

---

[4] estimation here is done using the R package ordinal (Christensen, 2015)

[5] see, Pratt (1981) and Agresti (2010, §3.4.5) for sufficient conditions

# Requirements from a good estimator for PO models

Same or similar properties with the MLE (e.g. asymptotic efficiency)

Finite estimates and corresponding standard errors

Invariance to data (dis)aggregation

|         |      | Aggregated |   |   |   |   | Disaggregated |   |   |   |   |
|---------|------|---|---|---|---|---|---|---|---|---|---|
|         |      | rating | | | | | rating | | | | |
| contact | temp | 1 | 2 | 3 | 4 | 5 | 1 | 2 | 3 | 4 | 5 |
| no      | cold | 4 | 9 | 5 | 0 | 0 | 4 | 9 | 5 | 0 | 0 |
|         | warm | 0 | 5 | 8 | 3 | 2 | 0 | 4 | 6 | 1 | 2 |
|         | warm |   |   |   |   |   | 0 | 1 | 2 | 2 | 0 |
| yes     | cold | 1 | 7 | 8 | 2 | 0 | 1 | 7 | 8 | 2 | 0 |
|         | warm | 0 | 1 | 5 | 7 | 5 | 0 | 1 | 5 | 7 | 5 |

Optimal sampling properties which are preserved under linear parameter transformations (e.g. $L$ contrasts, reversal of categories and so on)

# Outline

# Cumulative link model[6]

Vectors of counts on $k$ ordered categories are from independent multinomial vectors $Y_1, \ldots, Y_n$ with

$$Y_i \,|\, x_i \sim \mathrm{Mult}(m_i, (\pi_{i1}, \ldots, \pi_{ik}))$$

$$g(\pi_{i1} + \ldots + \pi_{ij}) = \alpha_j + \beta^T x_i = \sum_{t=1}^{p+k-1} \theta_t z_{ijt}$$

$x_i$ is a $p$-vector of explanatory variables

$\alpha_1 < \ldots < \alpha_{k-1} < \alpha_k = \infty$ and $\beta \in \Re^p$

$\theta = (\alpha_1, \ldots, \alpha_{k-1}, \beta_1, \ldots, \beta_p)^T$

$g(.)$ is a monotone increasing, differentiable link function

Special cases

Proportional odds model: $g = \mathrm{logit}$

Proportional hazards model (grouped survival times): $g = \mathrm{cloglog}$

---

[6]see, McCullagh (1980) and Agresti (2010, §5.1)

# Bias reduction through adjusted score functions

**Maximum likelihood estimator**

$$\hat{\theta} \leftarrow \left\{ \sum_i \sum_{j=1}^{k-1} g'_{ij} \left( \frac{y_{ij}}{\pi_{ij}} - \frac{y_{ij+1}}{\pi_{ij+1}} \right) z_{ijt} = 0 \right\}$$

where $g'_{ij} = \mathrm{d}g^{-1}(\eta)/\mathrm{d}\eta$

**Bias-reduced estimator**[7]

An estimator with smaller asymptotic bias than $\hat{\theta}$ is

$$\theta^* \leftarrow \left\{ \sum_i \sum_{j=1}^{k-1} g'_{ij} \left( \overbrace{\frac{y_{ij} + c_{ij} - c_{ij-1}}{\pi_{ij}}}^{\text{adjusted response } y^*_{ij}} - \frac{y_{ij+1} + c_{ij+1} - c_{ij}}{\pi_{ij+1}} \right) z_{ijt} = 0 \right\}$$

where $c_{ij} = m_i g''_{ij} [Z_i F^{-1} Z_i^T]_{jj} / 2$ and $c_{i0} = c_{ik} = 0$

---

[7]see, K. (2014, RSSB) and K. and Firth (2009, B'ka) for method

# Iterative maximum likelihood fits

The kernel in the adjusted score (omitting $i$) is

$$\frac{y_j + d_j}{\pi_j} - \frac{y_{j+1} + d_{j+1}}{\pi_{j+1}}$$

where $d_j = c_j - c_{j-1}$

# Iterative maximum likelihood fits

The kernel in the adjusted score (omitting $i$) is

$$\frac{y_j + d_j}{\pi_j} - \frac{y_{j+1} + d_{j+1}}{\pi_{j+1}}$$

where $d_j = c_j - c_{j-1}$

**Empirical cumulative logits**

$$\log \frac{\pi_1 + \ldots + \pi_j}{\pi_{j+1} + \ldots + \pi_k} = \alpha_j$$

$d_1 = 0.5 - \pi_1$, $d_j = -\pi_j$ $(j = 2, \ldots, k-1)$, and $d_k = 0.5 - \pi_k$

1. add 0.5 to the counts of the first and last category only
2. use ML on the adjusted data

The bias-reduced estimators end up being the empirical cumulative logits

$$\alpha_j^* = \log \frac{Y_1 + \ldots + Y_j + 0.5}{Y_{j+1} + \ldots + Y_k + 0.5}$$

# Iterative maximum likelihood fits

The kernel in the adjusted score (omitting $i$) is

$$\frac{y_j + d_j}{\pi_j} - \frac{y_{j+1} + d_{j+1}}{\pi_{j+1}}$$

where $d_j = c_j - c_{j-1}$

**More general models**

The kernel can be re-expressed as

$$\frac{y_j + \overbrace{d_j l_j - \pi_j d_{j+1}(1 - l_{j+1})/\pi_{j+1}}^{\text{always} \geq 0}}{\pi_j} - \frac{y_{j+1} + d_{j+1}l_{j+1} - \pi_{j+1}d_j(1 - l_j)/\pi_j}{\pi_{j+1}}$$

where $l_j$ is 1 if $d_j > 0$ and 0 else

**Iterative maximum likelihood fits**

At the $u$th iteration

1 add $d_j^{(u)}l_j^{(u)} - \pi_j^{(u)}d_{j+1}^{(u)}(1 - l_{j+1}^{(u)})/\pi_{j+1}^{(u)}$ to $y_j$

2 fit the model on the adjusted counts with maximum likelihood

# Properties of bias-reduced estimator

$\theta^*$ is equivariant under linear transformations[8]

i.e. the bias-reduced estimator of $L\theta$ is $L\theta^*$

---

[8]see, K. (2014, RSSB, §6-7) for proofs

# Properties of bias-reduced estimator

$\theta^*$ is equivariant under linear transformations[8]

$\theta^*$ and $\hat{\theta}$ have the same asymptotic distribution, i.e. $N(\theta, F^{-1}(\theta))$[9]

First-order inference tools, like Wald tests, apply unaltered

Standard errors and estimated variance-covariance matrices, in general, can be computed using $F^{-1}(\theta^*)$

---

[8]see, K. (2014, RSSB, §6-7) for proofs
[9]see, Firth (1993) and K. and Firth (2009)

# Properties of bias-reduced estimator

$\theta^*$ is equivariant under linear transformations[8]

$\theta^*$ and $\hat{\theta}$ have the same asymptotic distribution, i.e. $N(\theta, F^{-1}(\theta))$[9]

$\theta^*$ has always finite components

|  | $\alpha_1$ | $\alpha_2$ | $\alpha_3$ | $\alpha_4$ | $\gamma_1$ | $\gamma_2$ | $\gamma_3$ | $\gamma_4$ | $\delta$ |
|---|---|---|---|---|---|---|---|---|---|
|  | Maximum likelihood | | | | | | | | |
| Estimates | -1.27 | 1.10 | 3.77 | $\infty$ | $\infty$ | 2.15 | 2.87 | $\infty$ | 1.47 |
| Std. errors | - | - | - | - | - | - | - | - | - |
|  | Bias reduction | | | | | | | | |
| Estimates | -1.19 | 1.05 | 3.50 | 5.20 | 2.62 | 2.05 | 2.65 | 2.96 | 1.40 |
| Std. errors | 0.50 | 0.44 | 0.74 | 1.47 | 1.52 | 0.58 | 0.75 | 1.50 | 0.46 |

Testing for proportional odds using $\hat{\theta}$

$W = 0.7502$ leading to a $p$-value of 0.861 (based on $\chi_3^2$)

---

[8]see, K. (2014, RSSB, §6-7) for proofs
[9]see, Firth (1993) and K. and Firth (2009)

# Properties of bias-reduced estimator

$\theta^*$ is equivariant under linear transformations[8]

$\theta^*$ and $\hat{\theta}$ have the same asymptotic distribution, i.e. $N(\theta, F^{-1}(\theta))$[9]

$\theta^*$ has always finite components

$\theta^*$ is invariant to data (dis)aggregation

| | | Aggregated | | | | | Disaggregated | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | rating | | | | | rating | | | | |
| contact | temp | 1 | 2 | 3 | 4 | 5 | 1 | 2 | 3 | 4 | 5 |
| no | cold | 4 | 9 | 5 | 0 | 0 | 4 | 9 | 5 | 0 | 0 |
| | warm | 0 | 5 | 8 | 3 | 2 | 0 | 4 | 6 | 1 | 2 |
| | warm | | | | | | 0 | 1 | 2 | 2 | 0 |
| yes | cold | 1 | 7 | 8 | 2 | 0 | 1 | 7 | 8 | 2 | 0 |
| | warm | 0 | 1 | 5 | 7 | 5 | 0 | 1 | 5 | 7 | 5 |

Adding constants + ML is dangerous for general models

[8]see, K. (2014, RSSB, §6-7) for proofs
[9]see, Firth (1993) and K. and Firth (2009)

## Aggregated representation

```
R> library("ordinal")
R> wine_agg <- xtabs(~ contact + temp + rating, data = wine)
R> ftable(wine_agg)

             rating 1 2 3 4 5
contact temp
no      cold        4 9 5 0 0
        warm        0 5 8 3 2
yes     cold        1 7 8 2 0
        warm        0 1 5 7 5

R> wine_agg <- data.frame(wine_agg)
```

# Disaggregated

```
R> inds <- with(wine_agg, temp == "warm" & contact == "no")
R> wine_sub1 <- rbind(wine_agg[inds, ], wine_agg[inds, ], wine_agg[inds, ])
R> freq1 <- c(0, 1, 2, 2, 0)
R> freq2 <- c(0, 1, 1, 1, 1)
R> wine_sub1$Freq <- c(wine_sub1$Freq[1:5] - freq1 - freq2, freq1, freq2)
R> wine_sub1$agg <- rep(1:3, each = 5)
R> wine_sub2 <- wine_agg[!inds,]
R> wine_sub2$agg <- 1
R> wine_dis <- rbind(wine_sub1, wine_sub2)
R> ftable(xtabs(Freq ~ agg + contact + temp + rating, data = wine_dis))

                 rating 1 2 3 4 5
agg contact temp
1   no      cold          4 9 5 0 0
            warm          0 3 5 0 1
    yes     cold          1 7 8 2 0
            warm          0 1 5 7 5
2   no      cold          0 0 0 0 0
            warm          0 1 2 2 0
    yes     cold          0 0 0 0 0
            warm          0 0 0 0 0
3   no      cold          0 0 0 0 0
            warm          0 1 1 1 1
    yes     cold          0 0 0 0 0
            warm          0 0 0 0 0
```

# Maximum likelihood over different data representations

```
R> m_agg <- clm(rating ~ contact, nominal = ~ temp, weights = Freq, data = wine_agg)
R> round(coef(summary(m_agg)), 3)

                Estimate Std. Error z value Pr(>|z|)
1|2.(Intercept)   -1.266         NA      NA       NA
2|3.(Intercept)    1.104         NA      NA       NA
3|4.(Intercept)    3.766         NA      NA       NA
4|5.(Intercept)   24.896         NA      NA       NA
1|2.tempwarm     -21.095         NA      NA       NA
2|3.tempwarm      -2.153         NA      NA       NA
3|4.tempwarm      -2.873         NA      NA       NA
4|5.tempwarm     -22.550         NA      NA       NA
contactyes         1.465         NA      NA       NA
```

```
R> m_dis <- clm(rating ~ contact, nominal = ~ temp, weights = Freq, data = wine_dis)
R> round(coef(summary(m_dis)), 3)

                Estimate Std. Error z value Pr(>|z|)
1|2.(Intercept)   -1.266         NA      NA       NA
2|3.(Intercept)    1.104         NA      NA       NA
3|4.(Intercept)    3.766         NA      NA       NA
4|5.(Intercept)   24.896         NA      NA       NA
1|2.tempwarm     -21.095         NA      NA       NA
2|3.tempwarm      -2.153         NA      NA       NA
3|4.tempwarm      -2.873         NA      NA       NA
4|5.tempwarm     -22.550         NA      NA       NA
contactyes         1.465         NA      NA       NA
```

# Maximum likelihood on adjusted data

```
R> m_agg_adj <- update(m_agg, weights = Freq + 0.5)
R> round(coef(summary(m_agg_adj)), 3)

                Estimate Std. Error z value Pr(>|z|)
1|2.(Intercept)   -1.280      0.474  -2.701    0.007
2|3.(Intercept)    0.865      0.397   2.179    0.029
3|4.(Intercept)    2.956      0.602   4.910    0.000
4|5.(Intercept)    4.442      1.056   4.206    0.000
1|2.tempwarm      -1.989      1.110  -1.792    0.073
2|3.tempwarm      -1.808      0.523  -3.460    0.001
3|4.tempwarm      -2.188      0.625  -3.500    0.000
4|5.tempwarm      -2.304      1.092  -2.110    0.035
contactyes         1.188      0.424   2.799    0.005
```

```
R> m_dis_adj <- update(m_dis, weights = Freq + 0.5)
R> round(coef(summary(m_dis_adj)), 3)

                Estimate Std. Error z value Pr(>|z|)
1|2.(Intercept)   -1.291      0.472  -2.733    0.006
2|3.(Intercept)    0.850      0.392   2.166    0.030
3|4.(Intercept)    2.936      0.597   4.918    0.000
4|5.(Intercept)    4.422      1.053   4.199    0.000
1|2.tempwarm      -1.431      0.854  -1.676    0.094
2|3.tempwarm      -1.707      0.495  -3.446    0.001
3|4.tempwarm      -2.210      0.617  -3.579    0.000
4|5.tempwarm      -2.367      1.084  -2.183    0.029
contactyes         1.158      0.410   2.825    0.005
```

# Bias reduction over different data representations

```
R> m_agg_br <- bpolr(rating ~ contact | 0 | temp, weights = Freq, data = wine_agg,
+                    method = "BR")
R> round(coef(summary(m_agg_br)), 3)

              Value Std. Error t value
contactyes    1.397      0.463   3.018
1|2.tempwarm  2.621      1.524   1.720
2|3.tempwarm  2.050      0.579   3.541
3|4.tempwarm  2.648      0.755   3.510
4|5.tempwarm  2.961      1.499   1.975
1|2          -1.195      0.499  -2.396
2|3           1.055      0.436   2.420
3|4           3.498      0.739   4.734
4|5           5.196      1.475   3.524
```

```
R> m_dis_br <- update(m_agg_br, data = wine_dis)
R> round(coef(summary(m_dis_br)), 3)

              Value Std. Error t value
contactyes    1.397      0.463   3.018
1|2.tempwarm  2.621      1.524   1.720
2|3.tempwarm  2.050      0.579   3.541
3|4.tempwarm  2.648      0.755   3.510
4|5.tempwarm  2.961      1.499   1.975
1|2          -1.195      0.499  -2.396
2|3           1.055      0.436   2.420
3|4           3.498      0.739   4.734
4|5           5.196      1.475   3.524
```

# Graduate admissions in Stanford U

**Data**

Admission scores and candidate characteristics
from 106 applications to the political science
PhD at Stanford University



rater's score ($1 < 2 < 3 < 4 < 5$)

interest in American politics and political theory ($z_{i1}$ and $z_{i2}$; 1:yes, 0:no)

standardized score on quantitative and verbal parts of GRE ($x_{i1}$ and $x_{i2}$)

gender ($g_i$; 0:male and 1:female)

**Proportional odds model**

$$\text{logit}(\pi_{i1} + \ldots + \pi_{ij}) = \alpha_j - \beta_1 x_{i1} - \beta_2 x_{i2} - \beta_3 z_{i1} - \beta_4 z_{i2} - \beta_5 g_i$$

ML estimates
$\hat{\beta}_1 = 1.993$, $\hat{\beta}_2 = 0.892$, $\hat{\beta}_3 = 2.816$, $\hat{\beta}_4 = 0.009$, $\hat{\beta}_5 = 1.215$

---

[10]rater F1 in the analysis in Jackman (2004); R package pscl (Jackman, 2015)

# Simulation results

|    |            | Bias  | MSE  | Bias$^2$/Variance (%) | Coverage (%) |
|----|------------|-------|------|-----------------------|--------------|
|    | $\beta_1$  | 0.13  | 0.14 | 13.90                 | 94.42        |
|    | $\beta_2$  | 0.05  | 0.06 | 5.02                  | 94.15        |
| ML | $\beta_3$  | 0.22  | 0.79 | 6.29                  | 94.68        |
|    | $\beta_4$  | 0.00  | 0.64 | 0.00                  | 94.50        |
|    | $\beta_5$  | 0.07  | 0.24 | 2.33                  | 94.21        |
|    | $\beta_1$  | 0.00  | 0.11 | 0.00                  | 95.05        |
|    | $\beta_2$  | 0.00  | 0.05 | 0.00                  | 95.09        |
| BR | $\beta_3$  | 0.01  | 0.59 | 0.01                  | 95.32        |
|    | $\beta_4$  | 0.00  | 0.56 | 0.00                  | 95.55        |
|    | $\beta_5$  | -0.00 | 0.21 | 0.00                  | 94.99        |

figures are based on 10000 samples under the maximum likelihood fit

# Outline

# Direction of shrinkage

Model is "shrunken" towards a binomial GLM for the boundary categories

**Demonstration**

Complete enumeration (3136) of tables of the form

|  | category | | | | | | total |
|---|---|---|---|---|---|---|---|
| $x$ | 1 | 2 | 3 | 4 | 5 | 6 | |
| -0.5 | | | | | | | 3 |
| 0.5 | | | | | | | 3 |

Model: $g(\pi_{i1} + \ldots + \pi_{ij}) = \alpha_j - \beta x_i$

Calculate fitted probabilities based on $\hat{\theta}$ and $\theta^*$ for each table and for $g = \mathrm{logit}$ and $g = \mathrm{cloglog}$.

```
Error in eval(expr, envir, enclos):  object 'ncat' not found
Error in par(mfrow = c(length(dat), ncat), mar = rep(c(0.3, 0.3/2), each = 2), :
object 'dat' not found
Error in make.link(link[[r]]):  object 'link' not found
Error in mtext(text = "fitted probability (BR)", line = 3, side = 2, outer =
TRUE): plot.new has not been called yet
Error in mtext(text = "fitted probability (ML)", line = 4, side = 1, outer =
TRUE): plot.new has not been called yet
```

BR probabilities for intermediate categories tend to shrink to 0

BR probabilities for 1st (6th) category tend to shrink to $g^{-1}(0)$ $(1 - g^{-1}(0))$

# Outline

# Discussion I

**Estimation properties**

$\theta^*$ has all the required properties when estimating cumulative link models and is always finite

First-order likelihood inference applies in a "plug-in" fashion

**Shrinkage**

Model is shrunken towards a binomial GLM for the boundary categories

Adjusted scores provide just enough regularization to correct for bias and improve inference. Different regularization schemes may be needed for other tasks (e.g. prediction)

**Other models**

Continuation ratio models with complementary log–log-link are equivalent to proportional hazards models in discrete time are equivalent[11]

# Discussion II

**Confidence intervals, hypothesis testing and model comparison**

can also be constructed using adjusted score statistics

$$s^*(\theta_-^*)^\top \left\{ F(\theta_-^*) \right\} s^*(\theta_-^*)$$

calibrated against $\chi^2$ distributions

Note: When testing for extreme effects, default tests (e.g. Wald or adjusted score) always reject due to the interplay of finiteness and discreteness

**High-dimensional nuisance specifications**

Bias reduction is particularly effective for inference about a low-dimensional parameter of interest in the presence of high-dimensional nuisance parameters[12]

e.g. panel-specific cutpoints on the latent scale with panel covariates

---

[11]see, Laara and Matthews (1985)

[12]see, Lunardon (2018)

# References I

Agresti (2010). *Analysis of Ordinal Categorical Data (2nd Edition)*. John Wiley & Sons.

Christensen, R. H. B. (2015). ordinal—regression models for ordinal data. R package version 2015.6-28. https://cran.r-project.org/package=ordinal.

Firth, D. (1993). Bias reduction of maximum likelihood estimates. *80*(1), 27–38.

Jackman, S. (2004). What do we learn from graduate admissions committees? a multiple rater, latent variable model, with incomplete discrete and continuous indicators. *Political Analysis 12*(4), 400–424.

Jackman, S. (2015). *pscl: Classes and Methods for R Developed in the Political Science Computational Laboratory, Stanford University*. Stanford, California: Department of Political Science, Stanford University. R package version 1.4.9.

K. (2014). Improved estimation in cumulative link models. *Journal of the Royal Statistical Society, Series B 76*(1), 169–196.

K. and D. Firth (2009). Bias reduction in exponential family nonlinear models. *Biometrika 96*(4), 793–804.

K. and D. Firth (2010). A generic algorithm for reducing bias in parametric estimation. *Electronic Journal of Statistics 4*, 1097–1112.

Kosmidis, I. (2018). *brglm2: Bias Reduction in Generalized Linear Models*. R package version 0.1.8.

Laara, E. and J. N. S. Matthews (1985). The equivalence of two models for ordinal data. *Biometrika 72*(1), 206–207.

Lunardon, N. (2018). On bias reduction and incidental parameters. *Biometrika 105*(1), 233–238.

McCullagh, P. (1980). Regression models for ordinal data. *42*, 109–142.

Palmgren, J. (1981). The Fisher information matrix for log linear models arguing conditionally on observed explanatory variables. *68*, 563–566.

# References II

Peterson, B. and J. Harrell, Frank E. (1990). Partial proportional odds models for ordinal response variables. *Applied Statistics 39*, 205–217.

Pratt, J. W. (1981). Concavity of the log likelihood (Corr: V77 p954). *Journal of the American Statistical Association 76*, 103–106.

Randall (1989). The analysis of sensory data by generalised linear model. *Biometrical Journal 7*, 781–793.

**Reduced-bias estimation of models with ordinal responses**  📷

$$\theta^* \leftarrow \left\{ \sum_i \sum_{j=1}^{k-1} g'_{ij} \left( \frac{y_{ij} + c_{ij} - c_{ij-1}}{\pi_{ij}} - \frac{y_{ij+1} + c_{ij+1} - c_{ij}}{\pi_{ij+1}} \right) z_{ijt} = 0 \right\}$$

**Paper**

Kosmidis (2014). Improved estimation in cumulative link models. *Journal of the Royal Statistical Society: Series B*, **76**. DOI:10.1111/rssb.12025

**Software**

**bpolr** R function in the supplementary material of the paper, soon be part of the **brglm2** R package[13]

🐦 IKosmidis_

✉ ioannis.kosmidis@warwick.ac.uk

🌐 http://www.ikosmidis.com

[13]Kosmidis (2018); https://github.com/ikosmidis/brglm2

# Outline

# Poisson trick and likelihood inference

Multinomial logistic regression: Multinomial counts $Y_{i1}, \ldots, Y_{ik}$ with $\sum_{j=1}^{k} Y_{ij} = t_i$ and probabilities $\pi_{i1}, \ldots, \pi_{ik}$ with $\sum_{j=1}^{k} \pi_{ij} = 1$ and

$$\log \frac{\pi_{ij}}{\pi_{ik}} = \beta_j^\top x_i \tag{1}$$

Poisson trick: View $Y_{ij}$ as a Poisson count with rate $\mu_{ij}$

$$\begin{aligned}
\log \mu_{ij} &= \lambda_i + \beta_j^\top x_i \quad (j = 1, \ldots, k) \\
\log \mu_{ik} &= \lambda_i
\end{aligned} \tag{2}$$

Then,

- the score equations for $\lambda_i$ cause $\sum_{j=1}^{k} \mu_{ij} = t_i$
- the score equations for $\beta_j$ are the same to those from model (1)
- Under $\sum_{j=1}^{k} \mu_{ij} = t_i$, model (2) gives the same inverse expected information for $\beta_1, \ldots, \beta_k$ as does model (1)[14].
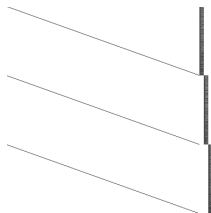
---

[14]see, Palmgren (1981)

# Poisson trick and bias reduction

The Poisson trick can be used for bias reduction if both $i(\theta)$ and $b(\theta)$ are evaluated under the restriction $\sum_{j=1}^{k} \mu_{ij} = t_i$ [15]

`brglm2::brmultinom` is just a wrapper that
- constructs the model matrix for model (2): [16]



- asks `brglmFit` to fit a Poisson log-linear model by re-calibrating $\mu_{ij}$ to satisfy $\sum_{j=1}^{k} \mu_{ij} = t_i$ at each iteration when calculating $i(\theta)$ and $b(\theta)$

---

[15]see, K. and Firth (2010) for theoretical justification

[16]Matrix is used to exploit sparsity; *eliminate* device can also be used

# brmultinom

```
R> library("brglm2")
R> library("nnet")
R> data("hepatitis", package = "pmlr")
R> ## Construct a variable with the multinomial categories according to
R> ## the HCV and nonABC columns
R> hepat <- hepatitis
R> hepat$type <- with(hepat, factor(1 - HCV*nonABC + HCV + 2 * nonABC))
R> hepat$type <- factor(hepat$type, labels = c("noDisease", "C", "nonABC"))
R> contrasts(hepat$type) <- contr.treatment(3, base = 1)
```

## multinom

```
R>    summary(multinom(type ~ group + time + group:time, weights = counts,
+                      data = hepat, trace = 0))

Call:
multinom(formula = type ~ group + time + group:time, data = hepat,
    weights = counts, trace = 0)

Coefficients:
       (Intercept) groupwithhold   timepost groupwithhold:timepost
C        -4.354779   -10.3461615 -1.5671468              9.8182339
nonABC   -4.866260    -0.4323412 -0.2655069              0.3192401

Std. Errors:
       (Intercept) groupwithhold   timepost groupwithhold:timepost
C        0.4502163    77.8477205  0.6351445             77.851156
nonABC   0.5799391     0.9159536  0.6539885              1.015315

Residual Deviance: 488.0431
AIC: 504.0431
```

# brmultinom

```
R> summary(brmultinom(type ~ group + time + group:time, weights = counts,
+                     data = hepat))

Call:
brmultinom(formula = type ~ group + time + group:time, data = hepat,
    weights = counts)

Coefficients:
       (Intercept) groupwithhold   timepost groupwithhold:timepost
C        -4.260116    -2.4257452 -1.5658843              1.9567429
nonABC   -4.712101    -0.3643221 -0.3763003              0.2563332

Std. Errors:
       (Intercept) groupwithhold   timepost groupwithhold:timepost
C         0.4302119    1.4815705  0.6062856              1.6321819
nonABC    0.5379321    0.8326882  0.6139492              0.9363615

Residual Deviance: 489.4283
Log-likelihood: -244.7142
AIC: 505.4283
```