

Generalized Linear Models (ST425A)

(Advanced level course, 7.5 hec, Aut. 2019)

Examnation (Part 1)

Gebrenewus Ghilagaber (Professor & Head)
Department of Statistics, Stockholm University

- **Date and time:** Monday 30 September 2019, 16:00 - 19:00
- **Permitted facilities:** Pocket calculator and sample of formula attached at the end of this exam.
- **Return of exam:** Not yet decided (information will be sent via E-mail or Athena).
- **Instructions:**
 - The exam consists of 4 questions and a sample of formula is given at the end
 - The total amount of points for this part of the exam is 30.
 - The minimum requirement to pass this part of examination is 20 points.
 - Solutions to each question should be detailed enough and well-motivated in order to get full points.

1 Question 1

Define a Generalized Linear Model and describe its "components".

2 Question 2

Consider the density function

$$f(y; \theta) = \theta \exp(-\theta y), \quad 0 < y < \infty, \theta > 0$$

- a) Show that the function belongs to the exponential family and indicate the canonical parameter
- b) Derive the maximum likelihood estimator of θ
- c) Suppose now we have 4 observations from the above density function with values $y_1 = 0.06$, $y_2 = 0.09$, $y_3 = 0.1$, and $y_4 = 0.2$. Use these values to compute a numerical value of the maximum likelihood estimator of θ
- d) Use your result in (c) to compute the value of the corresponding likelihood function
- e) Does your result in (d) look strange (somehow)? If yes, in what way and what lesson can one learn from the result?

3 Question 3

The Tables below contains results from analysis of effects of *Education* (with secondary-level education as a baseline) and *Age* on some demographic outcome among a sample of women.

Parameter	Levels	Symbol	Estimate	Std. Error	t-ratio
Constant	-	μ	-5.954	7.166	-0.83
Education	None	α_1	19.448	3.729	5.21
	Primary	α_2	4.144	3.191	1.30
Age	(linear)	β	0.1693	0.1056	1.60

Source of variation	Degrees of freedom	Sum of squares	Mean Square	F-Ratio
Age (linear)	1	1201.1	1201.1	36.5
Education Age	2	923.4	461.7	14.1
Residual	16	525.7	32.9	
Total	19	2650.2		

- a) What type of analysis is likely to have given rise to the above tables.
- b) Interpret the effect of Age on the outcome of interest.
- c) How large of the total variation in the outcome variable is due to variability in Age?
- d) Do the results in the table support the hypothesis of no education-effect (i.e. $H_0 : \alpha_2 = \alpha_3 = 0$)?

4 Question 4

The Table below presents results of analysis of count data across some models:

Estimates of interest	Fitted models			
	Poisson	Overdispersed Poisson	Negative Binomial	Zero-inflated
Coefficient (β)	-0.187	-0.187	-0.306	-0.187
Z-value	-5.12	-1.44	-3.10	-5.00
# param. estimated	p	$p + 1$	$p + 1$	$2p$
Log-likelihood	-18373	-18373	-9829	-15957

- a) Based on the results above, which of the fitted models provides the best description of the data? Don't forget to justify your answer
- b) Does the table indicate the original data is overdispersed?
- c) Does the table indicate the original data contains excess zeroes than would be expected?
- d) What is your conclusion concerning the explanatory variable of interest (whose coefficient is denoted by β)

A Formulas for Generalized Linear Models (A sample)

• Exponential family

$$f(y; \theta) = \exp[a(y)b(\theta) + c(\theta) + d(y)]$$

$$E[a(Y)] = \frac{c'(\theta)}{b'(\theta)}$$

$$var[a(Y)] = \frac{[b''(\theta)c'(\theta) - c''(\theta)b'(\theta)]}{[b'(\theta)]^3}$$

• Score statistic

$$U = a(y)b'(\theta) + c'(\theta)$$

$$U_j = \sum_{i=1}^N \left[\frac{var(Y_i)}{n_i} \right] x_{ij} \left(\frac{\partial \eta_i}{\partial \theta_j} \right)$$

• GLM weights

$$w_{ii} = \frac{1}{var(Y_i)} \left(\frac{\partial \eta_i}{\partial \theta_j} \right)^2$$

• Deviance

$$2 \left[\log L(\hat{\beta}_{max}; y) - \log L(\hat{\beta}; y) \right]$$

• Binomial distribution

– Probability function (using π in place of $\theta = p$):

$$f(y; n, \pi) = \binom{n}{y} \pi^y (1 - \pi)^{n-y}, \quad y = 0, 1, \dots, n$$

– Probit model

$$\pi_i = \Phi(\beta' x_i) \iff \pi_i = \Phi(\beta' x_i)$$

– Logit model

$$\beta' x_i = \ln \left(\frac{\exp(\beta' x_i) + 1}{\exp(\beta' x_i)} \right) = \ln \left(\frac{1 + \exp(-\beta' x_i)}{\exp(\beta' x_i)} \right)$$

$$\ln \left(\frac{\pi_{j+1} + \dots + \pi_j}{\pi_j} \right) = \mathbf{x}_T^j \beta_j, \quad j = 1, \dots, J - 1.$$

* Continuation ratio logit model

$$\ln \left(\frac{\pi_{j+1}}{\pi_j} \right) = \mathbf{x}_T^j \beta_j, \quad j = 1, \dots, J - 1.$$

* Adjacent categories model

$$\ln \left(\frac{\pi_{j+1} + \dots + \pi_j}{\pi_1 + \dots + \pi_j} \right) = \beta_{0j} + \beta_1 x_1 + \dots + \beta_{p-1} x_{p-1}, \quad j = 1, \dots, J.$$

* Proportional odds model

$$\ln \left(\frac{\pi_1 + \dots + \pi_j}{\pi_{j+1} + \dots + \pi_j} \right) = \mathbf{x}_T^j \beta_j, \quad j = 1, \dots, J - 1.$$

* Cumulative logit model

– Ordinal response:

$$\ln \left(\frac{\pi_j}{\pi_1} \right) = \mathbf{x}_T^j \beta_j, \quad j = 2, \dots, J.$$

– Normal logistic model

for $y_j = 0, \dots, n; j = 1, \dots, J$; and $\sum_j y_j = n$.

$$f(y_1, \dots, y_J; n_1, \dots, n_J; \pi_1, \dots, \pi_J) = \binom{n}{y_1, \dots, y_J} \pi_1^{y_1} \pi_2^{y_2} \dots \pi_J^{y_J},$$

– Probability function (again using π_i in place of $\theta_i = p_i$):

• Multinomial distribution

$$\pi_i = 1 - \exp[-\exp(\mathbf{x}_T^i \beta)] \iff \ln[-\ln(1 - \pi_i)] = \mathbf{x}_T^i \beta$$

– Complementary log-log link

$$f(y; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma}} \exp\left[-\frac{1}{2\sigma^2}(y - \mu)^2\right], \quad -\infty < y < \infty.$$

– Probability function (using μ and σ^2 in place of θ_1 and θ_2):

• Normal distribution

$$\ln(y_i) = \ln(n_i) + \mathbf{x}_i^T \boldsymbol{\theta}_j$$

– log link function

$$f(y; \lambda) = \frac{\lambda^y}{y!} e^{-\lambda}, \quad y = 0, 1, \dots$$

– Probability function (using λ in place of θ):

• Poisson distribution