

This is a repository copy of *Parametric models for biomarkers based on flexible size distributions*.

White Rose Research Online URL for this paper:

<https://eprints.whiterose.ac.uk/131018/>

Version: Accepted Version

---

**Article:**

Davillas, Apostolos and Jones, Andrew Michael [orcid.org/0000-0003-4114-1785](https://orcid.org/0000-0003-4114-1785) (2018)  
Parametric models for biomarkers based on flexible size distributions. *Health Economics*.  
pp. 1617-1624. ISSN 1057-9230

<https://doi.org/10.1002/hec.3787>

---

**Reuse**

This article is distributed under the terms of the Creative Commons Attribution (CC BY) licence. This licence allows you to distribute, remix, tweak, and build upon the work, even commercially, as long as you credit the authors for the original work. More information and the full terms of the licence here:

<https://creativecommons.org/licenses/>

**Takedown**

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing [eprints@whiterose.ac.uk](mailto:eprints@whiterose.ac.uk) including the URL of the record and the reason for the withdrawal request.

# Parametric models for biomarkers based on flexible size distributions

## Abstract

Recent advances in social science surveys include collection of biological samples. Although biomarkers offer a large potential for social science and economic research, they impose a number of statistical challenges, often being distributed asymmetrically with heavy tails. Using data from the UK Household Panel Survey (UKHLS), we illustrate the comparative performance of a set of flexible parametric distributions, which allow for a wide range of skewness and kurtosis: the four-parameter generalized beta of the second kind (GB2), the three-parameter generalized gamma (GG) and their three-, two- or one-parameter nested and limiting cases. Commonly used blood-based biomarkers for inflammation, diabetes, cholesterol and stress-related hormones are modelled. Although some of the three-parameter distributions nested within the GB2 outperform the latter for most of the biomarkers considered, the GB2 can be used as a guide for choosing among competing parametric distributions for biomarkers. Going “beyond the mean” to estimate tail probabilities, we find that GB2 performs fairly well with some disparities at the very high levels of HbA1c and Fibrinogen. Commonly used linear models are shown to perform worse than almost all the flexible distributions.

**Keywords:** biomarkers; generalised beta of second kind, heavy tails, tail probabilities

**JEL codes:** C18, C52, I14.

## 1. Introduction

Recent developments in social surveys include the integration of biomarkers and self-reported health measures. Biomarkers are objectively measured indicators of normal biological or pathogenic processes and, as such, offer at least two key advances over self-report health. First, biomarkers are not subject to reporting bias; given evidence for socioeconomic-related reporting bias in health, biomarkers offer a significant advantage in socioeconomic inequalities research (Bago d’Uva et al., 2008; Carrieri and Jones, 2017). Second, biomarkers can contribute to our understanding of the underlying biological factors through which socioeconomic conditions get “under the skin” (for example, thought stress-related physiological responses) and the role of socioeconomic exposures at earlier pre-symptomatic health states (Davillas et al., 2016; Jürges et al., 2013).

A growing literature analyses the effect of socioeconomic position on the conditional mean of biomarkers (e.g., Davillas et al., 2016, Jürges et al., 2013). However, biomarkers create several statistical modelling challenges as they often have skewed distributions with heavy tails (Jones, 2017). Furthermore, errors are likely to be heteroskedastic and responses to covariates may be nonlinear. Existing studies have estimated linear regression models using ordinary least squares (OLS) on raw or log transformed biomarkers (Jürges et al., 2013) and alternative inherently nonlinear specifications, such as the generalized linear models (GLM) (Davillas et al., 2016). While OLS on log rather than on levels might improve performance by reducing skewness, re-transformation to the raw scale –as health policymakers require– is highly challenging, requiring knowledge of the degree and form of heteroscedasticity (Jones et al., 2014). Although the GLM family deals with heteroskedasticity, it fails to explicitly account for skewness and kurtosis, imposing potential bias and efficiency losses (Jones et al., 2014).

Our paper contributes to the literature on modelling biomarkers by comparing the performance of a set of more flexible parametric distributions, the generalized beta of the second kind (GB2), the generalized gamma (GG) and their nine nested and limiting cases; we use nationally representative UK data on commonly used blood-based biomarkers for inflammation, diabetes, cholesterol and stress-related hormones (Carrieri and Jones, 2017). The GG and GB2 allow for a wide range of skewness and kurtosis to better

accommodate the biomarker data generation processes; these distributions have been proposed for fitting heavily skewed outcomes (for example, health care costs; Jones et al., 2014), to which biomarkers share similar distributional features. For comparison purposes, linear regression models using OLS are also estimated. Given that different biomarkers exhibit different distributions, identifying GB2 as a discriminatory tool amongst competing distributions might be useful for health researchers. Going “beyond the mean”, we also explore to what extent the GB2 and its nested cases that exerted the best goodness of fit (for each biomarker) regarding the whole distribution also perform well to predict tail probabilities.

## 2. Methods

The three-parameter GG distribution has been introduced as robust alternative to common estimation techniques for asymmetric data (Manning et al., 2005). More recently, Jones et al. (2014) have suggested adding further flexibility based on the four-parameter GB2 distribution. GB2 allows for a wider range of skewness and kurtosis, choosing among its several special or nested cases, while GB2’s extra flexibility may also enhance performance (Jones et al., 2014).

The GG distribution has a density function and conditional expectation that take the form:

$$f(y; \kappa, \mu, \sigma) = \frac{\gamma^\gamma}{\sigma y \sqrt{\gamma} \Gamma(\gamma)} \exp(z\sqrt{\gamma} - u) \quad (1)$$

and

$$E(y|x) = \exp(x'\beta) \left[ \kappa^{2\sigma/\kappa} \frac{\Gamma\left(\frac{1}{\kappa^2} + \frac{\sigma}{\kappa}\right)}{\Gamma\left(\frac{1}{\kappa^2}\right)} \right] \quad (2)$$

where,  $\gamma = |\kappa|^{-2}$ ,  $z = \text{sign}(\kappa)\{\ln(y) - \mu\}$ ,  $u = \gamma \exp(|\kappa|z)$ ,  $\mu = x'\beta$  and  $\Gamma(\cdot)$  is the gamma function. Parameters  $\kappa$  and  $\sigma$  are the shape parameters (Manning et al., 2005). The GG nests the gamma ( $\kappa = \sigma$ ), Weibull ( $\kappa = 1$ ), exponential ( $\kappa = 1, \sigma = 1$ ), and lognormal ( $\kappa = 0$ ) distributions.

The 4-parameter GB2 distribution adds further flexibility and has a probability density function and conditional mean of:

$$f(y; a, b, p, q) = \frac{ay^{ap-1}}{b^{ap} B(p, q) \left[ 1 + \left(\frac{y}{b}\right)^a \right]^{(p+q)}} \quad (3)$$

and

$$E(y|x) = b \left[ \frac{\Gamma(p+\frac{1}{a})\Gamma(q-\frac{1}{a})}{\Gamma(p)\Gamma(q)} \right] \quad (4)$$

where,  $b = \exp(x'\beta)$ ,  $B(\cdot)$  and  $\Gamma(\cdot)$  are the beta and gamma functions (Jones et al., 2014). Parameter  $a$  influences kurtosis and  $p$  and  $q$  the skewness of the distribution. We also estimate the nested and limiting cases of GB2; the three-parameter Beta of the second kind (B2) [ $a = 1$ ], Singh-Maddala (SM) [ $p = 1$ ] and Dagum [ $q = 1$ ]; the two-parameter Fisk [ $p = q = 1$ ], and Lomax [ $p = a = 1$ ]. GG itself is also a limiting case of the GB2, where  $b = q^{1/a}\beta$  and  $q \rightarrow \infty$  (Jones et al., 2014). We also estimate linear regression models using OLS.

The restrictions imposed by each of the special and limiting cases within the GG and GB2 are evaluated using Wald and likelihood-ratio (LR) tests. To assess the comparative performance of beta- with gamma-family distributions (being limited cases and not a linear restriction of a parameter), we compare Akaike (AIC) and Bayesian (BIC) information criteria across all models (Jones et al., 2014).

### 3. Data

The UK Household Panel Study (UKHLS) is a large, nationally representative UK study. At UKHLS wave 2, participants from its predecessor, the British Household Panel Survey (BHPS), were also incorporated. Non-fasted blood samples were collected, after the UKHLS wave 2 interview for the original UKHLS respondents and, at wave 3, for the BHPS sample. Pooling biomarker data from UKHLS waves 2 and 3 (2010-2013), resulted in a potential sample of 13,107 respondents.

Four biomarkers are used. Fibrinogen is an inflammatory biomarker, with higher values linked to cardiovascular morbidity and all-cause mortality risks (Davillas et al., 2017). Glycated haemoglobin (HbA1c) is a diagnostic biomarker for diabetes. The ratio of total cholesterol to high-density lipoprotein cholesterol is used as a marker for fatty substances in the blood. Dehydroepiandrosterone sulfate (DHEAS) is a steroid hormone and one of the mechanisms through which psychosocial stressors might affect health (Vie et al., 2014). Given our focus on the comparative performance of parametric distributions regarding goodness of fit, rather than explore potential effects from covariates, a

parsimonious set of covariates is used; polynomials of age (cubic or quartic depending on the biomarker used), gender, and their interactions to allow for flexible gender effects (Figure A1, appendix).<sup>1</sup>

#### 4. Results

Figure 1 presents the distribution of biomarkers (descriptive statistics in Table A1, appendix). Fibrinogen has a symmetric distribution but with fat tails (Figure 1). HbA1c is much more skewed (skewness statistic of 4.2 compared to zero for normal data) with long right-hand tails and excess kurtosis (31.15 versus 3 for normal data; Table A1). The cholesterol ratio and DHEAS also exhibits long right-hand tails and high kurtosis.

Table 1 contains restriction tests for the nested and limiting cases within the GG and GB2. Across all biomarkers, we find no evidence in support of any of the special cases within the GG distribution. For fibrinogen, we are unable to reject the null hypothesis of the restriction being valid for the SM model. Our results for HbA1c do not support any of the nested distributions. For the cholesterol ratio, both the LR and Wald tests favour the B2 distribution. Although the Wald test also fails to reject the null hypothesis for SM, this is not confirmed by the LR test; this disparity reflects the wide confidence intervals for GB2's  $p$  parameter (which include both one, satisfying the SM restriction, but also zero; Table A2, appendix). Our results for DHEAS favour the SM distribution.

**Table 1. LR and Wald tests (p-values) for special cases of the GB2 and GG.**

	<i>Fibrinogen</i>		<i>HbA1c</i>		<i>Cholesterol ratio</i>		<i>DHEAS</i>	
	LR	Wald	LR	Wald	LR	Wald	LR	Wald
<b><i>GB2 vs...</i></b>								
B2	0.000	0.000	0.000	0.000	<b>0.247</b>	<b>0.193</b>	0.000	0.000
SM	<b>0.208</b>	<b>0.236</b>	0.000	0.000	0.000	0.188	<b>0.703</b>	<b>0.710</b>
Dagum	0.004	0.013	0.000	0.000	0.000	0.020	0.000	0.000
Fisk	0.002	0.000	0.000	0.000	0.000	0.000	0.000	0.000
Lomax	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
<b><i>GG vs...</i></b>								
Gamma	0.000	0.024	0.000	0.000	0.000	0.000	0.000	0.000
Log Normal	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
Weibull	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
Exponential	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000

<sup>1</sup> The limited number of covariates may also alleviate concerns that, for less parsimonious specifications, the best specification for each model need to be compared rather than using the same covariates (Jones et al., 2014). However, the relative performance of our models (Table 2) remained the same in the case of no covariates.

Table 2 shows that AIC and BIC results are in accordance with the tests of Table 1. For all biomarkers, linear regressions estimated by OLS perform worse than each of the four- and three-parameter and most of the more parsimonious distributions. For fibrinogen, GB2 and SM perform best according to AIC and BIC criteria, with the latter showing the best performance. GB2 outperforms all the competing distributions regarding HbA1c. While the B2 and the SM distribution exhibit the best performance for the Cholesterol ratio and DHEAS, GB2 is ranked the second best.

Figure 2 presents the conditional tail probabilities (at  $k$  equal to 10<sup>th</sup>, 25<sup>th</sup>, 50<sup>th</sup>, 75<sup>th</sup> and 90<sup>th</sup> quantile) and spike plots of the actual-fitted difference (bias) for the GB2 distribution and its nested cases exerted the best performance for each biomarker (Table 2).

Specifically, 20-quantiles of the fitted values from these models are used to split the sample to calculate within-quantiles means of actual  $[P(y > k)]$  and predicted  $[P(y > k|X)]$  probabilities.

There are limited differences in the predictive ability of the more parsimonious distributions compared to GB2, confirming previous evidence that a flexible distribution is not a substitute for finding the correct distribution (Jones et al., 2014). GB2 performs reasonably well at predicting tail probabilities, although there are some disparities at the very high fibrinogen levels (90<sup>th</sup> quantile) and HbA1c above the pre-diabetes threshold ( $HbA1c \geq 42$ ).

**Table 2. AIC and BIC for each model.**

	Fibrinogen		Hba1c		Cholesterol ratio		DHEAS	
	AIC	BIC	AIC	BIC	AIC	BIC	AIC	BIC
GB2	<b>20866</b>	20948	<b>72138</b>	<b>72219</b>	<b>39175</b>	39257	<b>53800</b>	53889
B2	21221	21296	76134	76371	<b>39173</b>	<b>39249</b>	53897	53979
SM	<b>20865</b>	<b>20939</b>	72329	72404	39432	39506	<b>53798</b>	<b>53880</b>
Dagum	20872	20947	72927	73001	39315	39390	53855	53937
Fisk	20883	20950	73563	73629	39482	39549	54149	54223
Lomax	51843	51910	112182	112249	59542	59624	61959	62040
GG	21204	21278	74986	75060	39180	39270	53927	54016
Log-normal	21502	21569	77305	77372	39306	39373	54407	54482
Gamma	21219	21287	79049	79116	39867	39934	53942	54016
Weibull	22804	22871	88676	88743	42443	42518	54640	54715
Exponential	51841	51900	112180	112239	59540	59615	61957	62031
OLS	21500	21558	84119	84178	42875	42950	58371	58446

## 5. Conclusion

We illustrate the comparative performance of a set of more flexible parametric distributions, the GB2, GG, and their nested and limiting cases for a set of biomarkers. Although some of the three-parameter distributions nested within the GB2 (mainly the B2 and SM) outperform the latter in most of the biomarkers considered, GB2 can be used as a guide for choosing among competing distributions; a potentially useful message for applied researchers given that different biomarkers follow different distributions. The linear models estimated by OLS are dominated by almost all the competitive models. GB2 performs well at predicting biomarkers' tail probabilities, although with some disparities at the very high levels of fibrinogen and HbA1c.



## References

Bago d'Uva, T., O'Donnell, O., van Doorslaer, E., 2008. Differential health reporting by education level and its impact on the measurement of health inequalities among older Europeans. *International Journal of Epidemiology*, 37(6), 1375–1383.

Carrieri, V., Jones, A.M. (2017). The Income–Health Relationship ‘Beyond the Mean’: New Evidence from Biomarkers. *Health Economics*, 26(7), 937-956.

Davillas, A., Benzeval, M., Kumari, M. (2017). Socio-economic inequalities in CRP and fibrinogen across the adult age span. *Scientific Reports*, 7(1), 2641.

Jones, A.M. (2017). Data visualization and health econometrics, *Foundations and Trends in Econometrics*, 9, 1-78.

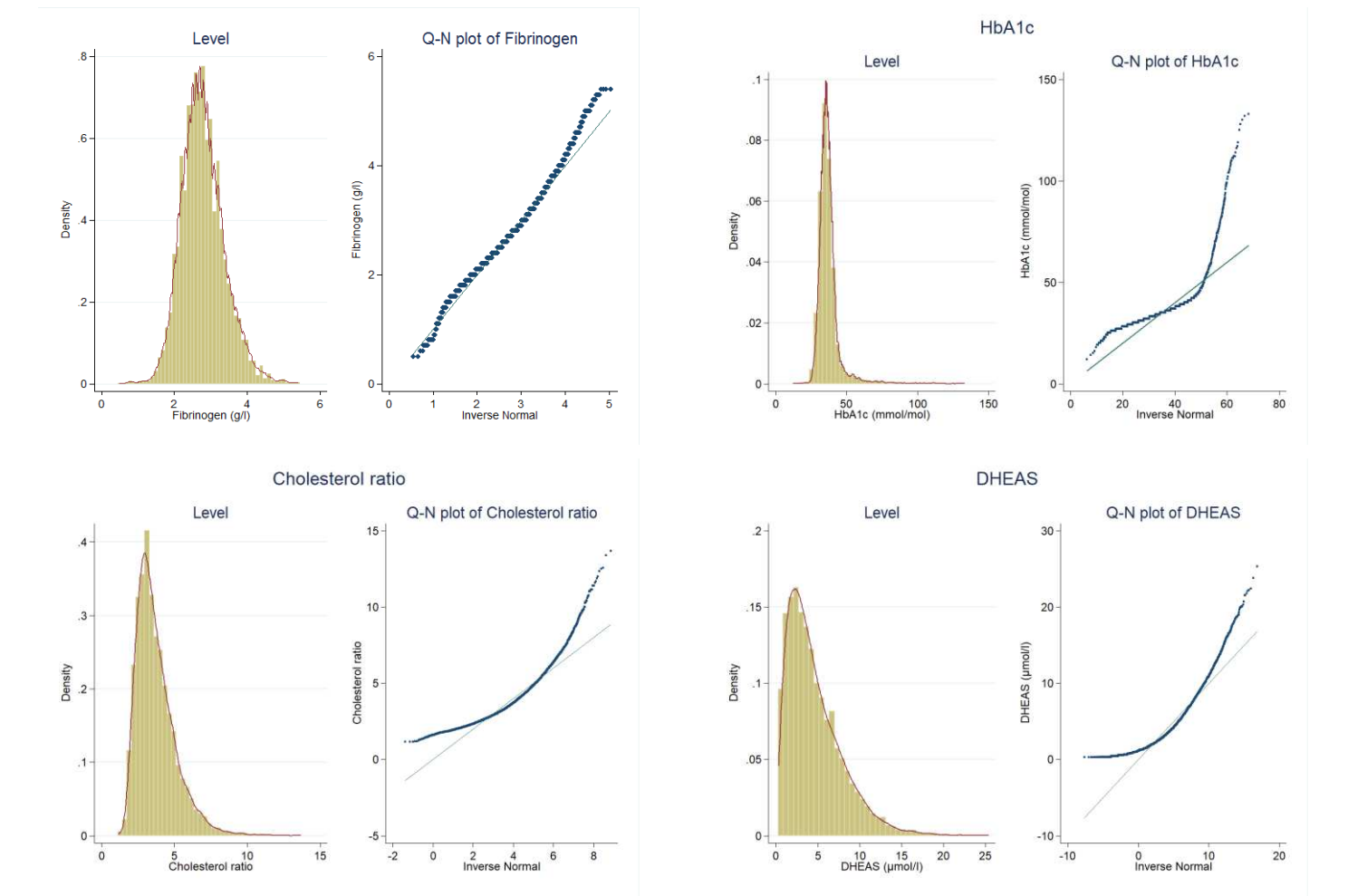
Jones, A.M., Lomas, J., Rice, N. (2014). Applying beta-type size distributions to healthcare cost regressions. *Journal of Applied Econometrics* **29**, 649-670.

Jürges, H., Kruk, E., Reinhold, S. (2013). The effect of compulsory schooling on health-evidence from biomarkers. *Journal of Population Economics*, 26(2), 645-672.

Manning, W. G., Basu, A., Mullahy, J. (2005). Generalized modeling approaches to risk adjustment of skewed outcomes data. *Journal of Health Economics*, 24(3), 465-488.

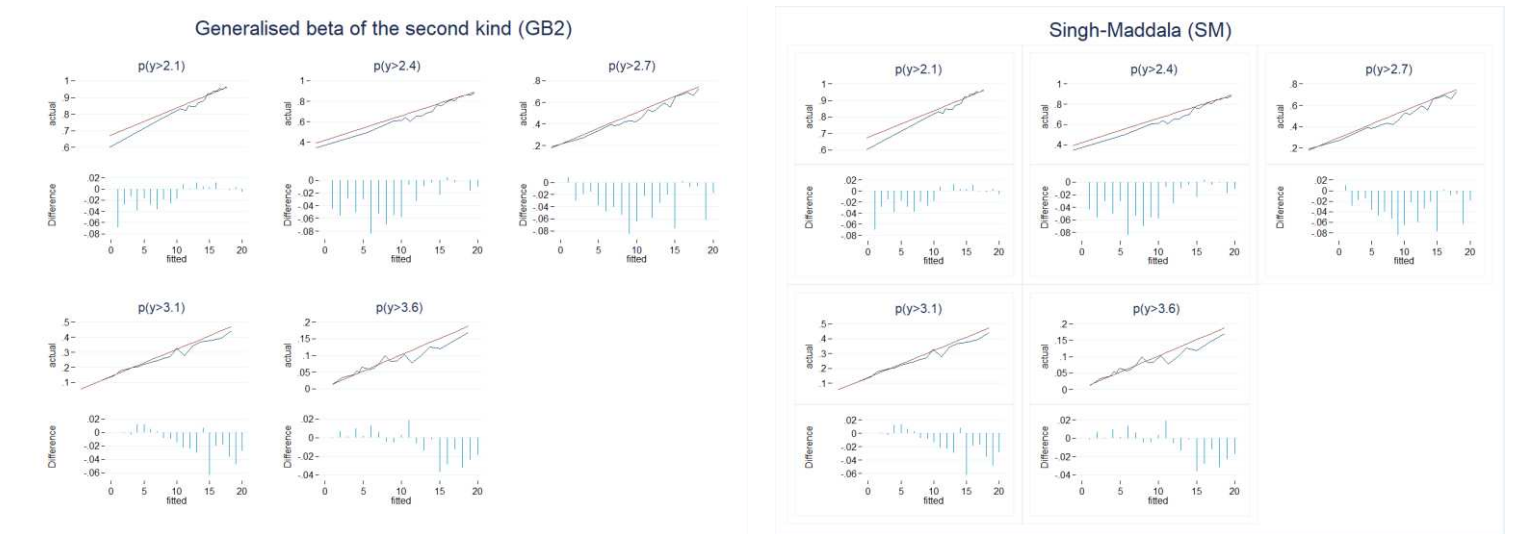
Vie, T., Hufthammer, K. O., Holmen, T. L., Meland, E., Breidablik, H. J. (2014). Is self-rated health a stable and predictive factor for allostatic load in early adulthood? Findings from the Nord Trøndelag Health Study. *Social Science & Medicine*, 117, 1-9.

**Figure 1. Distribution of biomarkers and quantile-normal (Q-N) plots**

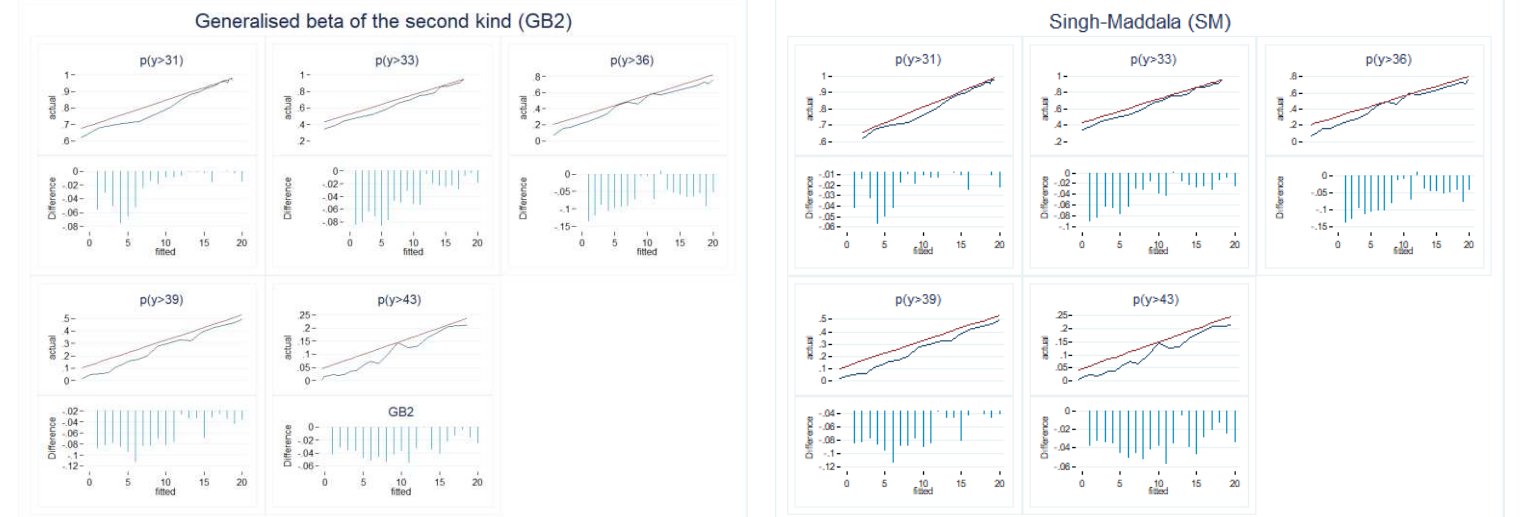


**Figure 2. Actual versus fitted tail probabilities.**

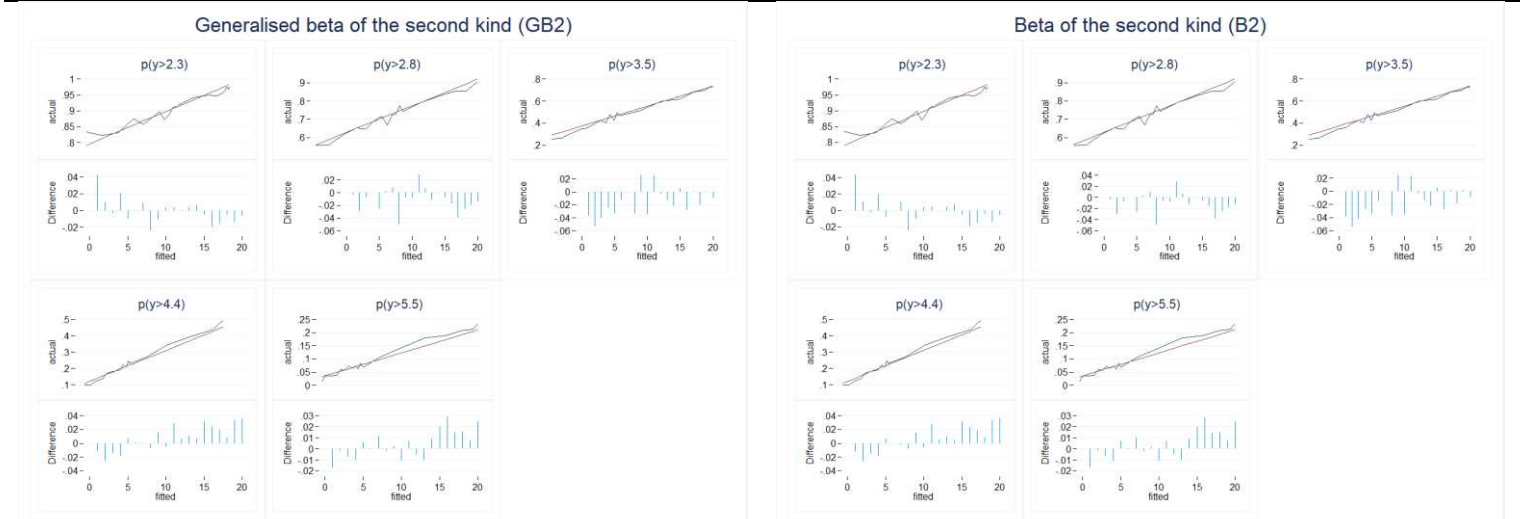
*Fibrinogen*



*HbA1c*

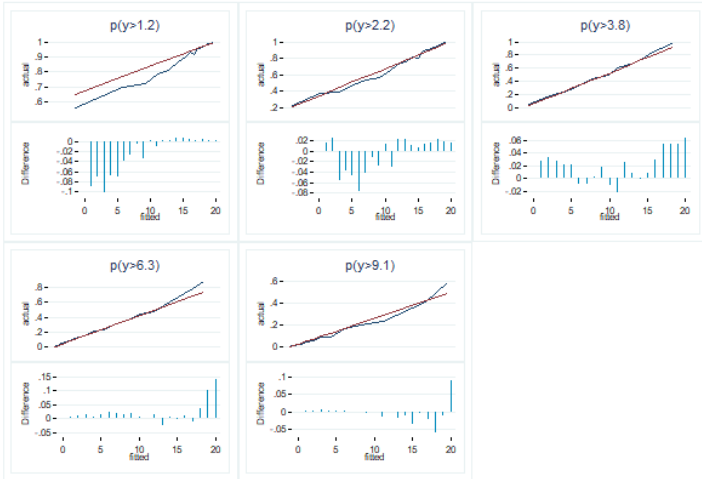


*Cholesterol ratio*



*DHEAS*

Generalised beta of the second kind (GB2)



Singh-Maddala (SM)

