

gkwreg: An R Package for Generalized Kumaraswamy Regression Models for Bounded Data

José Evandelton Lopes¹ and Wagner Hugo Bonat¹

DOI:

1 Paraná Federal University, Brazil

Software

- [Review ↗](#)
- [Repository ↗](#)
- [Archive ↗](#)

Submitted:

Published:

License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC-BY](#)).

Summary

`gkwreg` is an R package for fitting regression models to data restricted to the unit interval $(0, 1)$, such as proportions, rates, and indices. The package implements the flexible five-parameter Generalized Kumaraswamy (GKw) distribution and its seven main subfamilies, including the widely used Beta and Kumaraswamy distributions. A key feature of `gkwreg` is its use of the Template Model Builder (TMB) framework, which leverages automatic differentiation and C++ templates to provide fast, stable, and accurate maximum likelihood estimation. This overcomes the significant computational challenges typically associated with such complex multiparametric models, making them accessible for practical application. The package provides a user-friendly interface with standard R methods for model specification, inference, and diagnostics.

Statement of need

Statistical modeling of data bounded in the interval $(0, 1)$ is frequent across fields such as economics, epidemiology, and social sciences. Traditional methods like variable transformations followed by linear regression often present interpretability issues and fail near boundary points.

Direct modeling using distributions defined on $(0, 1)$ is preferable. While the Beta distribution is commonly used, it can be insufficient for complex patterns and lacks a closed-form cumulative distribution function (CDF). The Kumaraswamy (Kw) distribution (Kumaraswamy, 1980) offers an analytically simple CDF, yet its two-parameter form may be overly restrictive. To overcome these limitations, the Generalized Kumaraswamy (GKw) distribution, a flexible, five-parameter family incorporating the Beta and Kw distributions introduced by (Carrasco, Ferrari, & Cordeiro, 2010) was developed. However, practical applications of GKw in regression contexts have faced computational challenges. Its complex likelihood function makes Maximum Likelihood Estimation (MLE) computationally demanding and unstable, necessitating efficient and user-friendly computational tools.

This work introduces the R package `gkwreg`, which addresses this need. Built on the Template Model Builder (TMB) package (Kristensen, Nielsen, Berg, Skaug, & Bell, 2016), it leverages automatic differentiation (AD) in C++ to efficiently compute gradients and Hessians, enhancing speed, accuracy, and stability of MLE, especially when distribution parameters vary with covariates. `gkwreg` offers an intuitive interface aligned with standard R modeling conventions. Its integration with the multi-part formula syntax of the `Formula` package (Zeileis & Croissant, 2010) allows flexible specification of regression structures. Additionally, it provides comprehensive S3 methods (`summary()`, `predict()`, `plot()`, `residuals()`) and randomized quantile residuals (Dunn & Smyth, 1996) for model diagnostics, facilitating robust goodness-of-fit assessments.

Mathematics

The Probability Density Function (PDF) of the five-parameter Generalized Kumaraswamy (GKw) distribution is given by:

$$f(y; \boldsymbol{\theta}) = \frac{\lambda \alpha \beta y^{\alpha-1}}{B(\gamma, \delta+1)} (1-y^\alpha)^{\beta-1} [1 - (1-y^\alpha)^\beta]^{\gamma\lambda-1} \left\{ 1 - [1 - (1-y^\alpha)^\beta]^\lambda \right\}^\delta$$

where $\boldsymbol{\theta} = (\alpha, \beta, \gamma, \delta, \lambda)^\top$ is the vector of positive shape parameters and $B(\cdot, \cdot)$ is the beta function.

Distributional Regression Framework

`gkwreg` implements a comprehensive distributional regression framework where all relevant distribution parameters can be modeled as functions of covariates through flexible link functions. For a response variable $y_i \in (0, 1)$ following a GKw family distribution, each parameter $\theta_{ip} \in \{\alpha_i, \beta_i, \gamma_i, \delta_i, \lambda_i\}$ is related to a linear predictor via:

$$g_p(\theta_{ip}) = \eta_{ip} = \mathbf{x}_{ip}^\top \boldsymbol{\beta}_p$$

where $g_p(\cdot)$ is a suitable link function, \mathbf{x}_{ip} is the covariate vector for the i -th observation and p -th parameter, and $\boldsymbol{\beta}_p$ is the corresponding coefficient vector. The package employs an extended formula syntax allowing users to specify parameter-specific linear predictors through the notation `y ~ alpha_predictors | beta_predictors | gamma_predictors | delta_predictors | lambda_predictors`. Multiple link functions are supported: for positive parameters ($\alpha, \beta, \gamma, \lambda$), options include logarithmic (default), square root, inverse, and identity links; for probability parameters ($\delta \in (0, 1)$), logit (default), probit, complementary log-log, and Cauchy links are available. Additionally, link scaling functionality allows control over transformation intensity, enhancing numerical stability in challenging optimization scenarios. Maximum likelihood estimation maximizes the log-likelihood function:

$$\ell(\boldsymbol{\Theta}; \mathbf{y}, \mathbf{X}) = \sum_{i=1}^n \log f(y_i; \boldsymbol{\theta}_i(\boldsymbol{\Theta}))$$

where TMB computes exact gradients $\nabla \ell$ and Hessian matrix \mathbf{H} via automatic differentiation, enabling fast and stable optimization through efficient algorithms such as `nlminb` (default) or alternative methods like BFGS, Nelder-Mead, and L-BFGS-B, all configurable through the `gkw_control()` function.

Model Diagnostics

Model diagnostics in `gkwreg` are primarily based on randomized quantile residuals, defined as:

$$r_i^Q = \Phi^{-1}(F(y_i; \hat{\boldsymbol{\theta}}_i))$$

where $F(y_i; \hat{\boldsymbol{\theta}}_i)$ is the fitted CDF evaluated at observation y_i with estimated parameters $\hat{\boldsymbol{\theta}}_i$, and $\Phi^{-1}(\cdot)$ is the quantile function of the standard normal distribution. If the model is correctly specified, these residuals should follow a standard normal distribution. The package provides six diagnostic plot types to assess model adequacy: residuals versus indices for detecting autocorrelation, Cook's distance for identifying influential observations, leverage versus fitted values for flagging high-leverage points, residuals versus linear predictors for checking linearity and heteroscedasticity, half-normal plots with simulated envelopes for distributional assessment, and predicted versus observed plots for overall goodness-of-fit evaluation. These diagnostics are accessible through a unified `plot()` method supporting both base R graphics and `ggplot2` visualization systems.

References

- Carrasco, J. M. F., Ferrari, S. L. P., & Cordeiro, G. M. (2010). A new generalized Kumaraswamy distribution. *arXiv preprint arXiv:1004.0911*.
- Dunn, P. K., & Smyth, G. K. (1996). Randomized quantile residuals. *Journal of Computational and Graphical Statistics*, 5(3), 236–244. doi:[10.1080/10618600.1996.10474708](https://doi.org/10.1080/10618600.1996.10474708)
- Kristensen, K., Nielsen, A., Berg, C. W., Skaug, H., & Bell, B. M. (2016). TMB: Automatic differentiation and Laplace approximation. *Journal of Statistical Software*, 70(5), 1–21. doi:[10.18637/jss.v070.i05](https://doi.org/10.18637/jss.v070.i05)
- Kumaraswamy, P. (1980). A generalized probability density function for double-bounded random processes. *Journal of Hydrology*, 46(1-2), 79–88. doi:[10.1016/0022-1694\(80\)90036-0](https://doi.org/10.1016/0022-1694(80)90036-0)
- Zeileis, A., & Croissant, Y. (2010). Extended model formulas in R: Multiple parts and multiple responses. *Journal of Statistical Software*, 34(1), 1–13. doi:[10.18637/jss.v034.i01](https://doi.org/10.18637/jss.v034.i01)