

Prediction can be safely used as a proxy for explanation in causally consistent Bayesian generalized linear models

Maximilian Scholz^{1,*}

Paul-Christian Bürkner¹

¹ Cluster of Excellence SimTech, University of Stuttgart, Germany

* Corresponding author; email: maximilian.scholz@simtech.uni-stuttgart.de

Abstract

Bayesian modeling provides a principled approach to quantifying uncertainty in model parameters and structure and has seen a surge of applications in recent years. Despite a lot of existing work on an overarching Bayesian workflow, many individual steps still require more research to optimize the related decision processes. In this paper, we present results from a large simulation study of Bayesian generalized linear models for double- and lower-bounded data, where we analyze metrics on convergence, parameter recoverability, and predictive performance. We specifically investigate the validity of using predictive performance as a proxy for parameter recoverability in Bayesian model selection. Results indicate that – for a given, causally consistent predictor term – better out-of-sample predictions imply lower parameter RMSE, lower false positive rate, and higher true positive rate. In terms of initial model choice, we make recommendations for default likelihoods and link functions. We also find that, despite their lacking structural faithfulness for bounded data, Gaussian linear models show error calibration that is on par with structural faithful alternatives.

Keywords: Bayesian workflow, causal modeling, predictive performance, parameter recoverability, generalized linear models, simulation study

1. Introduction

Probabilistic (Bayesian) modeling provides a principled approach to quantifying uncertainty in model parameters and model structure. In recent years, it has seen a surge of applications in almost all quantitative sciences and industrial areas [43, 82, 46]. To support the principled application of Bayesian methods, [46] propose an overarching workflow to conduct Bayesian data analysis. In short, the workflow asks users to pick an initial model and iteratively refine it, performing various checks on the way to ensure that probabilistic assumptions are sensible, computations are valid, and model results are trustworthy for the intended purpose. This basic model building loop is repeated until the user’s requirements are satisfied or no satisfactory model can be found with the available resources. Within the overarching workflow, there remain many open questions with regards to individual steps, sub-workflows and how to make optimal decisions to move forward through the iterative steps of the workflow. In this paper, we want to investigate aspects of both the initial model choice and its subsequent modifications during the model building iterations, in the context of *latent inferential goals* as introduced below.

1.1. *Latent Inferential Goals*

Modeling choices depend on the context of the analysis and users’ needs. In [23], we proposed a systematic way to describe desirable properties of Bayesian models in terms of utilities, which both provides a framework for principled modeling decisions and helps making user requirements and chosen trade-offs explicit. Here, we focus on latent inferential goals, that is, goals that aim to make inference about unobservable variables (parameters) by means of data-informed models [23]. Below, we give a short introduction to some of the important modeling utilities relevant for the latent inferential goals studied in this paper. The presented utilities (as well as others discussed in [23]) are applicable and useful more generally beyond the Bayesian perspective on modeling.

Causal consistency is the utility about the degree to which a model can support claims about cause and effect, for example, about the effect of a treatment on a medical condition. While standard statistical inference can handle the static nature of associations, causality requires additional assumptions to build upon [93]. In terms of a utility for statistical models, causal consistency means that the model is consistent with a theoretically justified causal graph of its contributing variables. Causal consistency can be seen as a prerequisite for statistical models to provide valid answers for latent inferential goals [23]. However, as we also learn from the results of the present study, certain kinds of causal inconsistencies can be much more detrimental than others [29].

Convergence is the degree to which a posterior approximator approaches its closest possible approximation to the analytic posterior. Under certain conditions, Markov-Chain Monte Carlo (MCMC), for example, provides an asymptotically unbiased posterior approximation given infinite samples [43]. As we rarely have time to wait an infinity, convergence is a measure of whether we are sufficiently close to the ideal target. This is relevant in practice as model parameter estimates only have the chance to be trustworthy after convergence has been reached [74]. Accordingly, convergence is another prerequisite for making valid latent inference on real data.

Parameter recoverability (PR) is a model’s ability to recover the (latent) parameters of an assumed true data generating process (DGP; [23]). It is also known as the “explanation” perspective on statistical modeling [107, 133] and lies at the core of latent inferential goals. The practical problem with PR is that it requires concrete knowledge of the true DGP. Accordingly, it has to be studied in artificial settings where the true DGP is known, with the hope that the real data of interest fulfills the assumptions of these artificial settings sufficiently well [23]. The latter hope can be supported by evidence obtained from model utilities that are available at real data inference time. In this paper, we will specifically focus on predictive performance as one such supporting utilities, as elaborated further in Section 1.2.

Predictive performance (PP) generally describes the ability of a model to accurately predict new outcome data from existing observations. It is one of the most prominent utilities of model performance in Bayesian data analysis and a central tool for model comparisons [120, 43, 82]. In most cases, one is interested in out-of-sample (OOS) PP, as

predicting unseen data is almost always the actual real-world goal [120]. PP is conceptually similar to PR in that both target the accurate estimation of model-implied quantities [23], with the main difference that the former targets quantities that are observable at real data inference time (i.e., outcome variables), which allows to derive estimates that are independent of knowledge about the true DGP [120].

1.2. Prediction as Proxy for Explanation

From a philosophical perspective, prediction and explanation may be viewed as fully compatible [133]. However, in reality, the available finite data, as well as unknown misspecifications of causal or statistical assumptions in the fitted models often renders explanation and prediction at least partially opposing [17, 107, 133, 91]. In the context of the here-considered latent inferential goals, explanation is the primary target. Accordingly, in this case, PP reduces to a conveniently available supporting utility that *ideally* helps to select models with better PR at real data inference time [23]. In practice at least, and despite the theoretical arguments for caution, this assumption is very commonly (and often implicitly) made whenever explanatory model choices are based on OOS posterior predictive metrics or their approximations, such as AIC (e.g., [125, 82]), DIC [111], WAIC [126, 122] or ELPD-LOO [120, 122].

But under which circumstances is it actually valid to select explanatory models according to such predictive estimates (or average over them according to their predictive weights; [132])? It is one of the goals of this paper to shed some light on this question. Based on a number of counter-examples [82, 107, 13, 51], the available research demonstrates that this relationship cannot hold in general. Yet, despite a lot of literature on the theoretical distinction between explanation and prediction (e.g., see [107, 133] for an overview), the practical side of this topic has received comparably less attention.

For the remainder of this paper, we will refer to the question of using predictive performance as proxy for parameter recoverability as the *PP4PR question*. Among the factors that are likely influencing the (local) answer to this question, are the (in-)consistencies between the true model and the fitted models, both causally and statistically, employed estimation algorithms, and even the particular realization of observed data. To study these influencing factors systematically, we focus on the model class of (Bayesian) generalized linear models, as they constitute a single framework of highly common, well structured yet flexible models.

1.3. Generalized Linear Models

Despite or perhaps because of their simplicity, regression models make up a big part of all statistical data analyses. Their success can be explained by several factors, involving the ease of interpretation of their additive structure, rich mathematical theory, and (relative) simplicity of their estimation [45, 53, 92, 47]. With all of these advantages, the vast amount of modeling options and required analyst choices still render building trustworthy, well-predicting, and well-explaining regression models a difficult task.

The generalized linear models (GLMs) we focus on here all consider a univariate response variable y that is assumed to follow a parametric *likelihood* distribution, often called likelihood family [10, 19], with one main centrality parameter μ that is predicted

as well as zero or more auxiliary distributional parameters ψ_1, \dots, ψ_P that are assumed constant over observations. For the n th of a total of N observations in the training data, we write

$$y_n \sim \text{likelihood}(\mu_n, \psi_1, \dots, \psi_P). \quad (1)$$

The domain of all distributional parameters is specific to the given likelihood family. Regardless of its domain, μ may be predicted by a vector of predictor variables (also called features or covariates) $X = (x_1, \dots, x_K)$ where each variable x_k is itself a vector of length N . However, if the domain of μ does not span to whole real line, a *link function* has to be introduced such that the transformed domain becomes fully unbounded. For the n th observation we write

$$\text{link}(\mu_n) = \sum_{j=0}^J b_j f_j(X_n). \quad (2)$$

or equivalently

$$\mu_n = \text{inv_link} \left(\sum_{j=0}^J b_j f_j(X_n) \right). \quad (3)$$

In Equations (2) and (3), X_n denotes the vector (x_{1n}, \dots, x_{Kn}) of predictor values of the n^{th} observation, f_j are deterministic (possibly non-linear) transformations of the predictor variables and b_j are the regression coefficients. Typically, $f_0 = 1$ is a constant function and used as an intercept. The inverse link function inv_link is also known as *response function* as it transforms the unbounded linear predictor onto the possibly restricted domain of μ . In the context of artificial neural networks, the inverse link is also known under the term *activation function*.

When setting up GLMs, the four important choices the analyst has to make are (a) the likelihood family, (b) the link function, (c) the linear predictor term, and (d) whether and how to regularize, that is, for Bayesian GLMs, what prior distributions to use. All of them are mutually related [44], but specifically (a) and (b) are closely intertwined as the choice of link function depends on the support of μ and thus on the chosen likelihood. It is these two choices that we concentrate most of our efforts on here.

1.4. Choice of Likelihood and Link

In a nutshell, there are two kinds of approaches to choosing a likelihood. The first is to search within the space of *structurally faithful* likelihoods that respect the variable type of y [23, 43, 82], for example, an exponential or Gamma likelihood for positive continuous data that has no or no known upper bound. The second is just using a normal likelihood with identity link (i.e., linear regression) regardless of response type. The latter approach is openly advocated for comparably rarely [54] but de-facto applied across many disciplines because of its convenience and interpretability of the obtained regression coefficients. Still, there are obvious drawbacks of the "linear regression for all" approach, for example, that it can produce predictions that are impossible in reality (e.g., negative counts). What is more, it may seriously distort effect size estimates, their uncertainty, and sometimes even their sign [113, 78, 130, 76]. When going the structural faithful route, a lot of different options become available to the analyst, which adds the burden of choosing among them.

In theory, there are of course an infinite number of possible distributions to consider. In practice, the number of useful and practical available distributions is much smaller, but still substantial. For example, in general-purpose regression software such as the R packages `gamlss` [112], `VGAM` [134], or `brms` [19], analysts can choose among dozens likelihoods in total with at least several being available for each response type.

Within the space of structural faithful likelihoods, distributions can still differ substantially. One important aspect is how they deal with aleatoric (irreducible) variability in the responses. Location-scale families such as normal, logistic, or Student-t have a dispersion parameter that takes full control over the variability independently of the distribution’s location. In families with bounded support, variability and location are usually related in the sense that variability increases the further away the location is from the boundaries. But even then families differ in how they account for variability. For example, Poisson would assume equi-dispersion (equal variance and mean), negative binomial would offer equi- and over-dispersion (higher variance than mean), while Conway-Maxwell Poisson would offer equi-, over-, and under-dispersion [108, 131]. Other important aspects that differ among structurally faithful likelihoods is how to deal with skewness or tail-heaviness. They could either be fully dependent on location and variability or can have their own specifically dedicated parameters as in student-t (for heavy tails), skew-normal (for skewness), or skew-Student-t (for both), which are all generalization of the symmetric, thin-tailed normal distribution [5, 2]. Especially in the context of regression models, the kind of location parameter that μ represents is another crucial aspect in which likelihoods differ. Most canonically, μ is the mean but for some likelihood parameterizations it rather represents the median or mode, specifically if the mean has no analytic form [108, 52]. Of course, the meaning of the regression coefficients change according to the meaning of μ but this is not typically considered or acknowledged in practice.

The choice of link function can further complicate the situation. In common text books, links are often just presented as an immediate implication of the likelihood choice, such that the focus is on default links [82, 43, 47, 79], that is the identity link if μ is unbounded, the log-link if μ is lower-bounded and the logit-link if μ is double-bounded¹. While we are not aware of any relevant competitor to the identity link in the unbounded case, there exist multiple alternatives for both the lower- and double-bounded cases (see Appendix A.2 and A.4).

The focus on the likelihood rather than the link can be understood from multiple perspectives: First, the distributional assumptions for probabilistic quantities is at the heart of statistical modeling and we have a better basis to argue about it (see above) then about the choice of link function. Second, specifically in the context of binary regression (requiring links for double-bounded parameters), it has been repeatedly argued [36, 47, 98] that, as [47] puts it, different links ”provide essentially identical substantive conclusions”. Overall, research on dedicated link function choice remains comparable thin and mostly focuses on specific application in an applied field [e.g., 31, 57, 75]. In some contrast to the apparently common understanding, the results we present in this paper suggest that

¹Any continuous lower-bounded variable can be linearly shifted such that its support is the positive numbers. Any continuous upper-bounded variable (without a lower bound) can be sign reversed and shifted such that its support is again the positive numbers. Any continuous double-bounded variable can be shifted and scaled such that its support is the unit interval. Accordingly, it is sufficient to consider link functions with either positive or unit interval support.

the choice of link function may, in certain scenarios, very well be important for latent inferential goals.

1.5. *Aims and Structure of This Article*

The range of practically relevant likelihood classes is quite big and encompasses, among others, likelihoods for unbounded, lower-bounded, and double-bounded continuous data, as well as binary, categorical, ordinal, count, and proportional (sum-to-one) data [58, 59, 112, 134, 21]. Studying all of them at once would be too large of a scope for a single paper. Here, we focus our efforts on GLMs with lower-bounded or double-bounded continuous likelihoods. Within these classes, we not only have several qualitatively different (non-nested) likelihood options, but can also study both main classes of non-identity link functions.

The first aim is to study the utility of prediction as proxy for explanation combined with several causal (mis-)specification mechanisms. As we will do so by means of extensive simulations, we can use the large amount of investigated modeling cases to make practical recommendations for the choice of both likelihoods and link functions in the context of Bayesian GLMs for lower-bounded and double-bounded continuous data, which constitutes our second aim. As an additional contribution, we have developed the R package *bayesim* [105] in the process of working on this paper, which facilitates large-scale simulation studies of Bayesian models under both causal and statistical misspecification.

2. Methods

In this section, we will explain the simulation study’s design, including the considered likelihoods and link, the data generation process, the fitting of Bayesian GLMs to the simulated data, and the metrics we calculate for each fitted model. The basic architecture of the study generates datasets (D) from a list of data generation configurations (G), fits models using model fitting configurations (F) from all relevant combinations of likelihoods, links, and linear predictor terms on each dataset, and calculate metrics of interest for each model. The flow of the simulation is illustrated in Figure 1. Many aspects of this simulation study are not only dependent on the methods and algorithms used, but also on their current implementations. The simulation was written in R [99] using Stan [1, 25] and brms [19]. Our simulation framework is available as an R package [105] and the code for the simulation configurations can be found in our online appendix [104].

2.1. *Considered Likelihoods and Link Functions*

Below, for brevity, we only shortly list the likelihoods and links included in the simulation study. A detailed review of the considered options and our inclusion criteria are available in Appendix A.

Models for double-bounded responses. We included the beta, Kumaraswamy, simplex, and transformed-normal likelihoods. The latter arise from applying the double-bounded links to the response variable y , instead of to the location parameter μ , as detailed in Appendix A.1. All of these likelihoods have two distributional parameters.

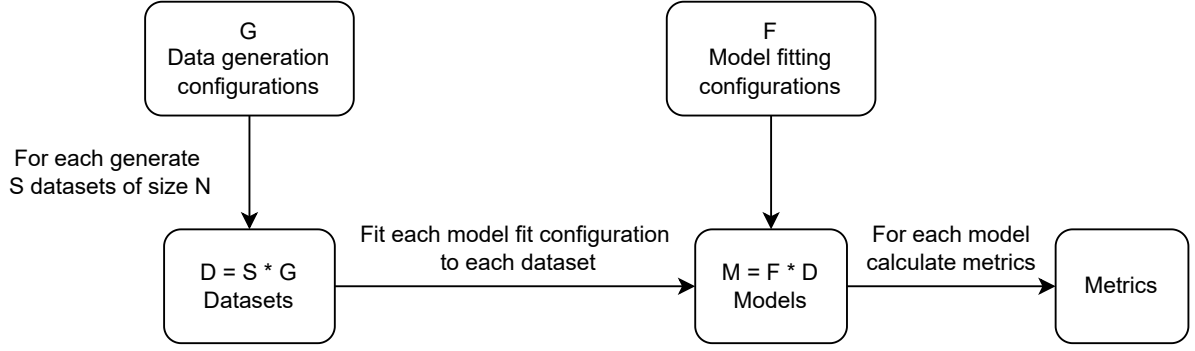


Figure 1: Conceptual simulation architecture. Multiple datasets are generated for each entry from a list of data generation configurations. Each dataset is then fitted with every entry from a list of model configurations and metrics are calculated for each fitted model.

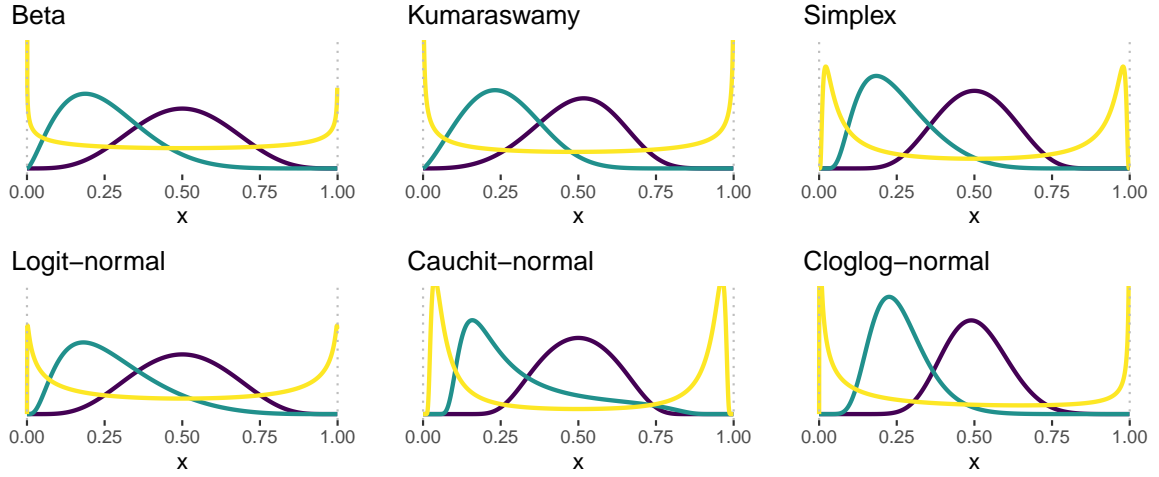


Figure 2: Exemplary illustrations of all included double-bounded densities each with three different shapes. The y-axis is truncated at 5 from above for better visibility of different shapes. For details, see Appendix A.1.

Figure 2 shows some exemplary densities for each of them, illustrating qualitatively different kinds of shapes they can accommodate. The three distinct shapes are uni-modal symmetric and asymmetric shapes as well as a bi-modal bathtub shape. For the remainder of the paper, we will refer to these shapes as symmetric, asymmetric, and bathtub, respectively. As link function, we included the logit, cauchit, and cloglog links, each of them having qualitatively different properties (see Appendix A.2). Figure 3 illustrates the included link functions and their corresponding response functions.

Models for lower-bounded responses. We included the gamma, Weibull, Fréchet, inverse Gaussian, beta prime, Gompertz and transformed-normal likelihoods as detailed in Appendix A.3. All of these likelihoods have two distributional parameters. Figure 4 shows some exemplary densities for each of them, illustrating qualitatively different kinds of shapes they can accommodate. The three distinct shapes are uni-modal thin tailed and heavy tailed shapes as well as a ramp shape. For the remainder of the paper we will

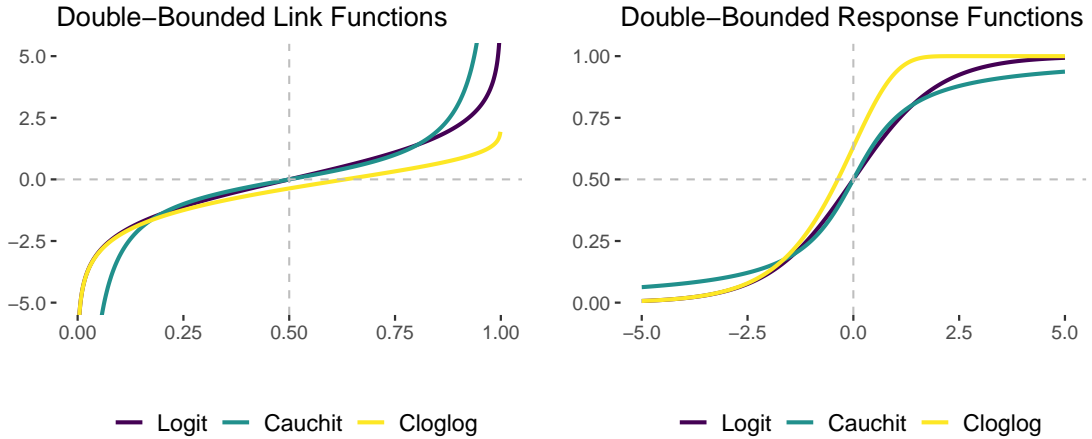


Figure 3: Double-bounded link and corresponding response functions. For details, see Appendix A.2.

refer to these shapes as thin tail, heavy tail, and ramp respectively. As link functions, we included the log and softplus links (see Appendix A.4). Figure 5 illustrates the included link functions and their corresponding response functions.

2.2. Data Generation

The foundation of the simulation study are the individual simulated datasets D . We generate them from the combinations of likelihoods and link functions as presented in Section 2.1. In practice, we usually cannot know if a model actually represents the true DGP², not only in terms of likelihood and link but also in terms of included predictors and their (non-)linear combination on the latent space. Here, we use a causal directed acyclic graph (DAG) [93] as the foundation for our data generating scenario of the linear predictor term. Using DAG-based data generation allows us to add causally misspecified models to our simulation, which use a set of predictors that differs from the true DGP. Using causal DAGs has the added benefit of good theoretical understanding of how the misspecification should influence parameter estimation. Figure 6 shows the DAG we use for data generation. It is based on the work of [29] and combines four types of controls into one structure. The DAG consists of an outcome y , a treatment x , and four additional variables z_1, z_2, z_3, z_4 that correspond to four qualitatively different types of controls. Our aim in this study is to estimate the causal effect of x on y .

The *do*-calculus framework [95] utilizes causal graphs representing assumptions about cause and effect of the data at hand to state if a specific model can make unbiased estimation about a causal effect. See Section 2.3 for how we use this when fitting models.

²Below a certain level of abstraction, the assumed model is probably never equal to the true DGP in reality.

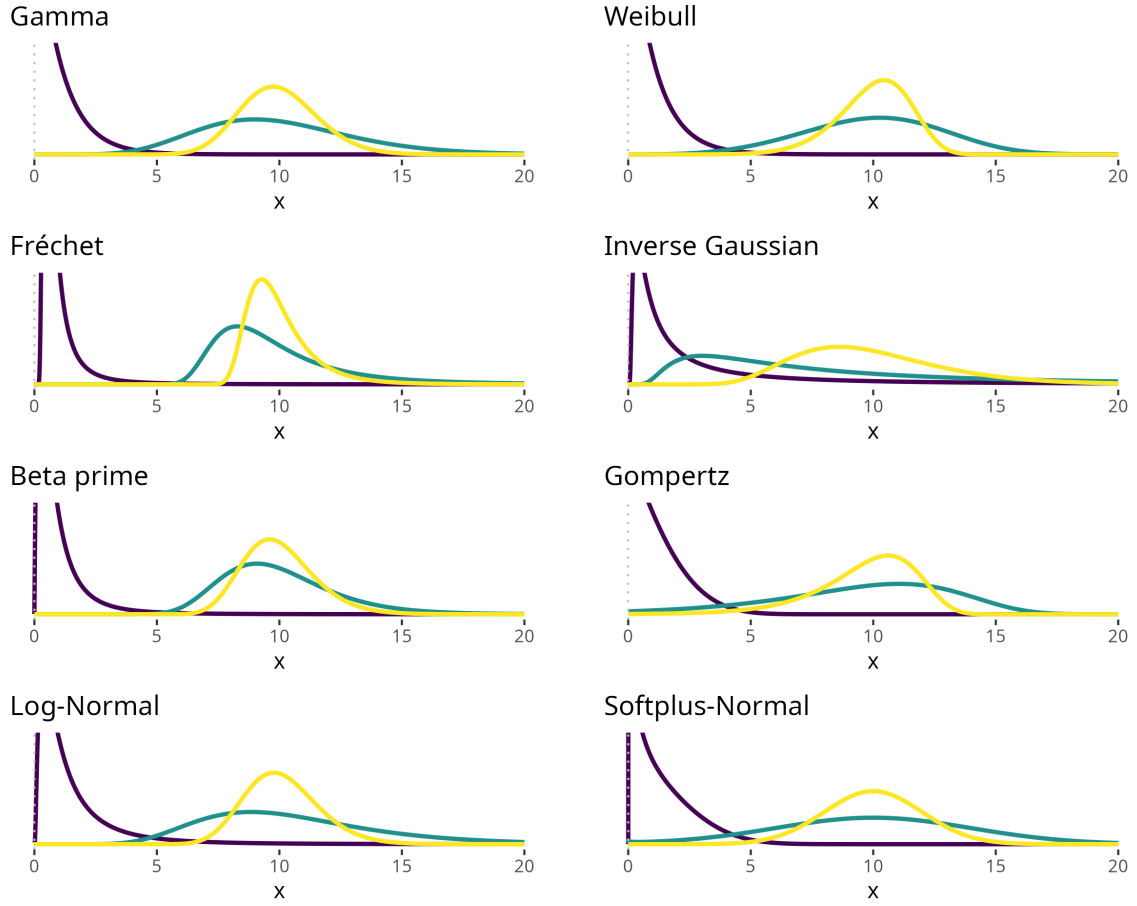


Figure 4: Exemplary illustrations of all included lower-bounded densities each with three different shapes. For details, see Appendix A.3. The y-axis is truncated at 0.4 from above for better visibility of different shapes.

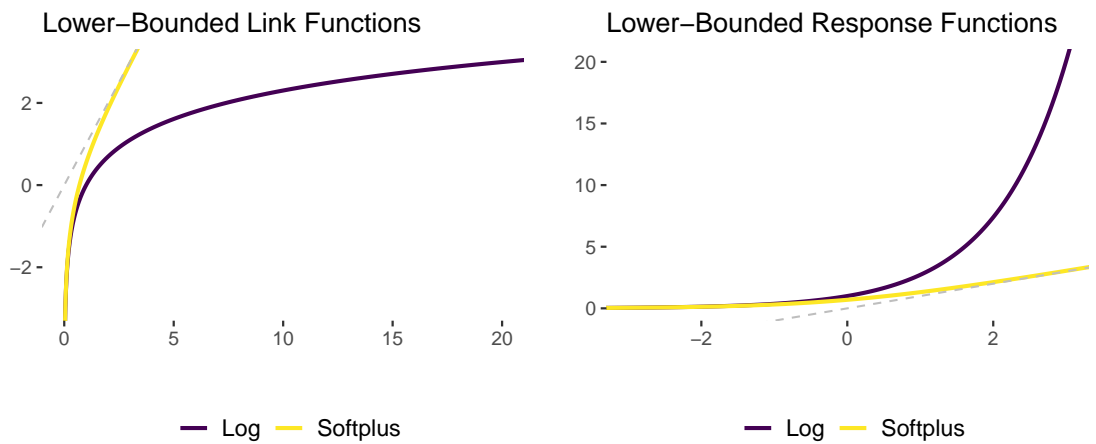


Figure 5: Lower-bounded link and corresponding response functions. For details, see Appendix A.2. The gray dashed line indicates the identity function.

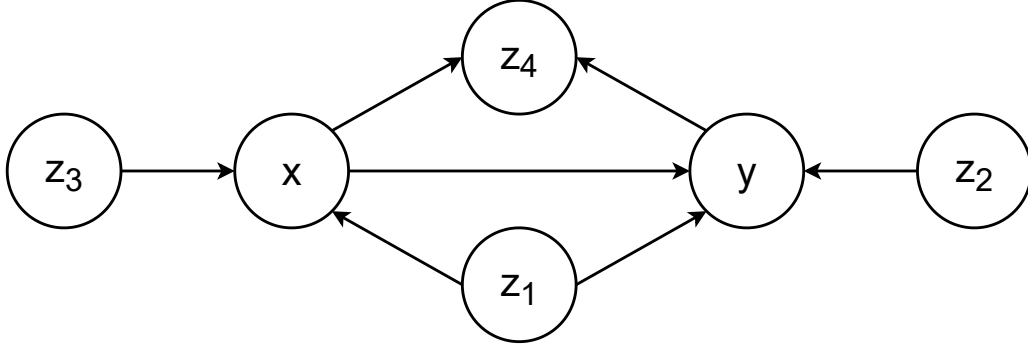


Figure 6: Full simulation study DAG. The aim is to estimate the causal effect of x on y . The four z_i allow for different causal misspecifications in the assumed models.

From there, we generated individual datasets of $N = 100$ observations as follows:

$$\begin{aligned}
z_1 &\sim \text{normal}(0, \sigma_{z_1}) \\
z_2 &\sim \text{normal}(0, \sigma_{z_2}) \\
z_3 &\sim \text{normal}(0, \sigma_{z_3}) \\
x &\sim \text{normal}(\beta_{z_1x}z_1 + \beta_{z_3x}z_3, \sigma_x) \\
y &\sim \text{likelihood}(\text{inv_link}(\alpha_y + \beta_{xy}x + \beta_{z_1y}z_1 + \beta_{z_2y}z_2), \phi) \\
z_4 &\sim \text{normal}(\beta_{xz_4}x + \beta_{yz_4}y, \sigma_{z_4})
\end{aligned}$$

where ϕ denotes the second distributional parameter of the specific likelihood family (see Appendix A). We chose normal distributions to generate all data variables besides y to limit the scope of our simulations. Parameters of the true DGP were fixed to sensible values as detailed below. A fully Bayesian approach to simulations would have meant to draw parameters from a prior distribution, instead of setting them to some fixed values [115]. However, in the present simulations, we choose against the former, because (a) the variation in the parameter caused by the priors would have blurred the effect of our carefully crafted causal misspecifications as well as different likelihood shapes, and (b) for exponential-based link functions, the variance even from reasonably varying prior distributions would lead to extreme datasets that are impossible in practice, a problem that bugs Bayesian simulations more generally [42, 85].

Each likelihood’s second distributional parameter ϕ as well as the intercept α_y were chosen to represent the qualitative shapes we presented in Section 2.1. The individual coefficients for x and z_i were then calibrated so that the parameter recovery was non-trivial for the model representing the true DGP while also preventing the causal misspecified models from obviously failing every time. The calibration procedure and simulation parameters chosen on that basis are documented in our online appendix [104]. An overview of the data generation configurations is given in Table 1. Despite extensive precautions in choosing the true DGP’s parameters, in few cases, response values under- or overflowed numerically to the lower and upper boundaries, respectively. To prevent this, we introduced a lower truncation bound of 0.000001 and (for the double-bounded data only) an upper truncation bound of 0.999999, such that all floating point operations were stable for all fitted models.

While we initially included the inverse Gaussian (IG) likelihood in our study, we

observed severe sampling problems during model fitting when using an IG likelihood, with more than half of the models failing to converge. This problem occurred consistently in multiple independent IG implementations. In addition, both the inclusion of the IG as true and assumed likelihood greatly increased the duration of our simulation study with some models taking more than 100 times longer to fit than the average model. For these reasons we decided to exclude the IG likelihood from our study.

We generated $S = 200$ datasets for each data generation configuration, as shown in Table 1 in a fully crossed design, which lead to a total of $D = 14,400$ datasets for both double- and single-bounded data.

Table 1: Data generation configurations

Factor	Levels
Double-bounded likelihoods	beta, Kumaraswamy, simplex, transformed-normal
Double-bounded links	logit, cauchit, cloglog, (identity)
Double-bounded shapes	symmetric, asymmetric, bathtub
Lower-bounded likelihoods	gamma, Weibull, transformed-normal, Fréchet, beta-prime, Gompertz
Lower-bounded links	log, softplus, (identity)
Lower-bounded shapes	ramp, heavy tail, thin tail
True β_{xy}	zero, positive

Note: The identity link is listed in parenthesis as the transformed-normal datasets technically are generated using an identity link, as the link transformation is part of the likelihood.

2.3. Model Fitting

After a dataset was generated, all the relevant models were fitted on the same dataset. Each model is based on a fit configuration that consists of a likelihood, a link, and an linear predictor term. The latter we will also refer to as *formula* from here on in reference to R formula syntax. The likelihoods and links are the ones presented in Section 2.1. So given a dataset generated from a double-bounded likelihood, we fitted several models on that dataset using all (fully crossed) combinations of double-bounded likelihoods, links, and DAG-based formulas (see below). The same approach was followed in the lower-bounded scenario. Additionally, we fitted models using a normal likelihood with the respective double- or lower-bounded links. While the normal likelihood does not respect the bounds of the data, it is a convenient default choice that many people use in practice (see Section 1.3). Finally, we also fit a model with a normal likelihood and identity link on each dataset to serve as a baseline, due to how commonly used linear regression is (see Section 1.4).

Using the graph in Figure 6, we can pose the query $P(y \mid \text{do}(x))$ to assess if it is possible to estimate the unbiased causal effect β_{xy} of x on y . If the do operator can be eliminated, using the rules of do-calculus, the query is valid and one can use the proposed model to make valid causal statements [93]. In this case we can use the (simplified) second rule of do-calculus: $P(y \mid \text{do}(x), z) = P(y \mid x, z)$ if z satisfies the back-door criterion. The back-door criterion is satisfied if a set of variables z blocks all back-door paths from x to y . Back-door paths are paths that start with an incoming arrow into x . The only backdoor path from x to y goes through z_1 , a fork [93]. There is another possible causal path through z_4 , a collider. However, collider paths are closed unless conditioning on

the collider. In order to build a valid causal query for the effect of x on y , we need to additionally condition on z_1 to close the backdoor path. While this model would result in an unbiased estimation of β_{xy} , adding z_2 will increase the precision of the estimation [29]. This results in the modified query $P(y \mid \text{do}(x), z_1, z_2)$ which, using the rule above turns into $P(y \mid x, z_1, z_2)$. While this model does not reflect the entire data generating process, we will refer to it as the *true (causal) model* for the remainder of the paper, as it is the subset of the true DGP that is sufficient for estimating β_{xy} optimally. By additionally in- or excluding every z_i , we get five possible formulas implying different (assumed) linear predictor terms for y each of which are causally misspecified:

$y \sim x + z_1 + z_2$	(true model of y)
$y \sim x + z_1$	(leaving out z_2 reduces the precision of the estimation)
$y \sim x + z_1 + z_2 + z_3$	(including z_3 reduces the precision of the estimation)
$y \sim x + z_2$	(leaving out z_1 biases the estimation)
$y \sim x + z_1 + z_2 + z_4$	(including z_4 biases the estimation)

One central question in a Bayesian modelling workflow is the choice of priors. Contrary to what we would recommend in practice, we used flat priors for all model parameters, as it is not clear to us how one would specify equivalent priors for the different auxiliary parameters ϕ (e.g., the precision parameter of the beta distribution) across all likelihoods. What is more, different links imply different latent scales, which renders the regression coefficients' scales incomparable across (assumed) links and thus further complicates equivalent prior specification. In a real-world analysis, we would prefer to use more informative priors, at least some that are weakly-informative [1, 43, 82]. Accordingly, the here-made choice is to be understood only in the context of the present simulations to achieve better comparability across likelihoods and links. Since the fitted models only have between 4 and 6 parameters, they are simple enough to be well identified on the basis of $N = 100$ observations alone. In initial adhoc experiments (not shown here), we have confirmed that the differences in posteriors, their estimation efficiency as well as the implied prediction metrics between models with flat vs. weakly informative was been minimal. In our opinion such minimal differences do not justify extensive evaluation of different prior choices, since prior specification is not in focus of the present paper (but see [1, 42] for current recommendations).

All models were fit using Stan [25, 1] via brms [19] with two chains, 500 warmup- and 2000 post-warmup samples, which leaves us with $S = 4000$ total post-warmup posterior samples per model that can be used for inference [1]. We used an initialization range of $\text{init} = 0.1$ around the origin on the unconstrained space as some models, especially the ones using a cloglog link, were struggling with the default initial value settings. For all other MCMC hyperparameters, we applied the brms defaults [19]. The exact R [99] code used for the model fitting is available in bayesim [105] and the online appendix [104]. An overview of the model fit configurations is given in Table 2.

The full factorial design results in $F = 80$ fit configurations for the double-bounded models and $F = 75$ fit configurations for the single-bounded models. Combined with the $D = 14,400$ datasets, this leads to a total of $M = 1,152,000$ double-bounded and

$M = 1,080,000$ single-bounded models fitted in our simulations.

Table 2: Model fit configurations.

Factor	Levels
Double-bounded Likelihoods	beta, Kumaraswamy, simplex, transformed-normal, normal
Double-bounded Links	logit, cauchit, cloglog, (identity)
Single-bounded Likelihoods	gamma, Weibull, transformed-normal, Fréchet, beta-prime, Gompertz, normal
Single-bounded Links	log, softplus, (identity)
Formulas (right-hand side)	$x + z_1 + z_2$, $x + z_2$, $x + z_1$, $x + z_1 + z_2 + z_3$, $x + z_1 + z_2 + z_4$

Note: The identity link is listed in parenthesis as the transformed normal models technically use an identity link while fitting, as the link transformation is part of the likelihood. In addition, the normal likelihood is also fit with an identity link.

2.4. Model-Based Metrics

To investigate the model utilities introduced in Section 1.1, we calculated several metrics per fitted model for convergence, PR, and PP, most of which are implemented in the `loo` [124] and `posterior` [22] packages as well as `bayesim` itself [105]. Causal consistency is directly based on the known relation of the assumed formula with the true DGP so requires no further metrics.

Convergence We computed the number of divergent transitions (DT) [14], as well as state-of-the-art versions of \hat{R} , bulk- and tail effective sample sizes ESS_{bulk} and ESS_{tail} [123]. In this paper we report the ratio of effective samples to total number of post-warmup samples, $\text{REFF} := \text{ESS}/S$, as $\text{REFF}_{\text{bulk}}$ and $\text{REFF}_{\text{tail}}$ respectively.

Estimation Speed We computed the fitting time needed per effective sample both with and without the warmup time included:

$$\text{TESS}_{\text{total}} := \frac{t_{\text{warmup}} + t_{\text{sampling}}}{\text{ESS}}, \quad (4)$$

$$\text{TESS}_{\text{sampling}} := \frac{t_{\text{sampling}}}{\text{ESS}}, \quad (5)$$

For ESS_{bulk} we call this $\text{TESS}_{\text{total}}^{\text{bulk}}$ and $\text{TESS}_{\text{sampling}}^{\text{bulk}}$, for ESS_{tail} , $\text{TESS}_{\text{total}}^{\text{tail}}$, and $\text{TESS}_{\text{sampling}}^{\text{tail}}$ respectively. As estimation speed is not a focus of this paper and strongly implementation dependent, we present the results only in the online appendix [104].

Parameter recoverability We calculated the posterior bias and RMSE of the model’s estimation of β_{xy} . Given a true parameter value $\tilde{\beta}$ and posterior samples for the estimation of said parameter $\beta^{(s)}$, we define (a sampling-based approximation of) the posterior bias, and RMSE of the estimate as

$$\text{bias}(\beta) := \mathbb{E} [\beta^{(s)}] - \tilde{\beta}, \quad (6)$$

$$\text{RMSE}(\beta) := \sqrt{\mathbb{E}[(\beta^{(s)} - \tilde{\beta})^2]}, \quad (7)$$

where $\mathbb{E}[\cdot]$ denotes expectations over the posterior distribution $p(\theta \mid \tilde{y})$ for training data \tilde{y} , which is approximated by the arithmetic mean over posterior samples. As we are interested in the size of the bias more than the direction, we will present the absolute bias $|\text{bias}(\beta)|$ in our results. The above are reasonable measures for comparing models only if the assumed link coincides with the true link of the DGP, as the link determines the scale of linear predictor and thus the scale of the regression coefficients. To compare models using different links, we calculated the false positive rate (FPR; also known as Type I- or alpha-error rate) and the true positive rate (TPR; also known as statistical power, which is the inverse of the Type II-error rate) as they are not scale dependent. In addition, we also present FPR and TPR in their combined form via receiver operating characteristic (ROC) curves [139].

Predictive Performance In the absence of any case-specific arguments for a particular predictive metric, log-probability scores are recommended as a general-purpose choice [122]. For this reason, we computed the expected log pointwise predictive density (ELPD, [120, 122]) as our main predictive metric, which is defined as

$$\text{ELPD}(y^*) := \sum_{i=1}^{N^*} \log p(y_i^*) = \sum_{i=1}^{N^*} \log \mathbb{E}[p(y_i^* \mid \theta)], \quad (8)$$

where y_i^* denotes test data (previously unseen by the model), which was generated from the same true DGP configuration as the training dataset. We refer to the above metrics as $\text{ELPD}_{\text{test}}$. In addition, we also calculated ELPD via leave-one-out cross-validation (LOO-CV) as approximated via Pareto-smoothed importance sampling (PSIS) [122, 124]. We refer to this metric as ELPD_{loo} . Approximate LOO-CV metrics have the advantage that they are readily available for real data inference, at the expense of being only an approximation that might fail to estimate out-of-sample PP accurately if there are influential observations [122]. As the $\text{ELPD}_{\text{test}}$ would usually not be available during a real world scenario, we mainly used it to double-check the reliability of the ELPD_{loo} results. As results of both metrics were highly comparable, conditional on passing the Pareto- \hat{k} diagnostic, so we do not report $\text{ELPD}_{\text{test}}$ results for brevity.

In our study we also used the ELPD_{loo} difference to the best performing model from the group of all models that were fit on the same dataset. This representation of relative PP has the benefit of being available in real world analysis scenario where the ground truth is unknown.

$$\Delta \text{ELPD}_{\text{loo}}(m_i) := \text{ELPD}_{\text{loo}}(m_i) - \max(\text{ELPD}_{\text{loo}}(M)), \quad (9)$$

where $M := m_1, \dots, m_n$ are all models, or a subset of those, that were fit on the same dataset. Unless stated otherwise, M will be restricted to models that used the same formula and in the case of scale dependent PR metrics the same link as the true DGP.

2.5. Statistical Analysis

To study the relationship between PR and PR, we fit Bayesian multilevel models (BMMs) to the simulation results in order to investigate the (predictive) relationship of $\Delta \text{ELPD}_{\text{loo}}$

with $|\text{bias}(\beta_{xy})|$, $\text{RMSE}(\beta_{xy})$, FPR, and TPR.

Before fitting the BMMS, we filtered out results from models that did not converge well enough, and thus would not be advisable to use during a real world analysis. Specifically, we treated models as converged if they had less than 10 divergent transitions, a $\hat{R} < 1.01$ and both $\text{ESS}_{\text{bulk}} > 400$ and $\text{ESS}_{\text{tail}} > 400$ as recommended by [123, 1]. In addition, we set an upper threshold of ≤ 5 Pareto- \hat{k} values above 0.7 as PSIS becomes unreliable with Pareto- \hat{k} values above 0.7 [121, 122]. This is relevant as the reliability of the ELPD_{loo} approximation depends on PSIS. Only very few converged models had any high Pareto- \hat{k} values at all, so the above threshold on their number was not practically relevant. Finally, we set a lower threshold for the $\Delta\text{ELPD}_{\text{loo}}$ of -100 to remove models that had exceptionally bad predictions compared to the best performing model (within the set of compared models). Some convergence problems can be fixed via hyper-parameter tuning (i.e., increased *adapt-delta* to reduce the number of divergent transitions) or longer sampling time. However, we expect models that had problems in the here-considered, rather simple scenarios to have even worse convergence in more complex cases, such that simple hyper-parameter tuning would not be an available remedy any more. As those models would not be viable in more complex scenarios and because refitting specific models within the simulations would have required a lot of manual interventions, we decided against pursuing convergence for all models.

We then fit BMMS with brms [19] using the following formula with $(\log) \text{RMSE}(\beta_{xy})$ as outcome:

```
log(rmse) ~ 1 + formula*data_link*data_shape +
  delta_elpd_loo:formula:data_link:data_shape +
  (1 + formula + delta_elpd_loo:formula | dataset)
```

The formula for $(\log) |\text{bias}(\beta)|$ was defined analogously.

```
zero_in_95_ci ~ 1 + formula*data_link*data_shape*beta_xy_category +
  delta_elpd_loo:formula:data_link:data_shape:beta_xy_category +
  (1 + formula + delta_elpd_loo:formula | dataset)
```

The relevant modeling terms were the overall ("fixed") effects of $\Delta\text{ELPD}_{\text{loo}}$ in interaction with the formula, the data generating link, and data generating shape as well as the varying ("random") effects of $\Delta\text{ELPD}_{\text{loo}}$ in interaction with the formula, where the dataset served as grouping factor. For the model on (binary) CI zero overlap, we used a Bernoulli likelihood and added a categorical term on whether or not $\beta_{xy} = 0$. We could directly model the zero-overlap as for $\beta_{xy} = 0$ it equals the FPR and for $\beta_{xy} \neq 0$ it equals the TPR. We included the formula in the interaction, as we do not aim for comparison across formulas. Due to causal misspecifications they can have a strong influence on PR and thus should not simply be averaged over. The grouping by dataset is necessary as comparison of models is only fair if they are all trained on the same dataset. The data generating shape and link influence the parameter scales which is why we also included them in the interactions with $\Delta\text{ELPD}_{\text{loo}}$. As both are unique per dataset, modeling their effects as varying across datasets would have been non-sensible.

In case of the $\text{RMSE}(\beta_{xy})$ as outcome, we used a log-normal likelihood expressed as a normal likelihood on $\log(\text{RMSE}(\beta_{xy}))$ to gain the speed improvements of highly optimized linear regression functions in Stan [1]. We ran a single MCMC chain with 500 warmup and

1000 post-warmup samples in threading mode (within-chain parallelization) on 20 CPU cores [127]. The use of only using a single chain and within-chain-parallelization was necessary to reduce the initial runtimes of several days, which exceeded the maximum job run time on our computing cluster, to something manageable (and feasible on our available hardware) of less than one day. All convergence metrics indicated sufficient convergence. As we were mainly interested in the qualitative patterns, rather than high resolution numerical results, we considered 1000 post-warmup samples leading to few hundred ESS as sufficient.

3. Simulation Results

In this section, we present a selection of results from our simulation study that are representative for the overarching results patterns. Additional results and the raw simulation data are available in our online appendix [104]. We split this section into three parts. First, we present summaries of the descriptive statistics for the double-bounded and lower-bounded results, as these regard individual likelihoods and links. Thereafter, we present the results from our investigation of the PP4PR question as these are less dependent on individual likelihoods and links.

3.1. *Double-bounded Results*

Generally, we found that most of the double-bounded likelihoods and links performed similar across many of the metrics we calculated (see Table 3). In terms of convergence, the average \hat{R} -values were far below the threshold of 1.01 for all combinations. The beta and transformed-normal likelihoods showed the best sampling behaviour with no convergence problems and the highest sampling efficiency (measured by $\text{RESS}_{\text{bulk}}$ and $\text{RESS}_{\text{tail}}$) of all likelihoods. On the other end of the spectrum, the cloglog link was especially problematic, causing higher numbers of DTs for the simplex (22 on average) and Kumaraswamy likelihoods (8 on average). In addition, the simplex likelihood generally had the worst sampling behaviour of all likelihoods, while the normal and Kumaraswamy likelihoods were in the middle.

For PR measured via $\text{RMSE}(\beta_{xy})$ and $|\text{bias}(\beta_{xy})|$ the same patterns showed for both metrics. The most obvious observation was the stark difference between the bathtub shape and both the symmetric and asymmetric shapes for both recovery metrics, as exemplarily shown in Figure 7 for the $\text{RMSE}(\beta_{xy})$ and logit link. For both the symmetric and asymmetric shapes, likelihood choice seems to have had little influence on $\text{RMSE}(\beta_{xy})$ and $|\text{bias}(\beta_{xy})|$. The beta likelihood had the best average $\text{RMSE}(\beta_{xy})$ and $|\text{bias}(\beta_{xy})|$, however the differences were heavily influenced by the bathtub shape. The normal and Kumaraswamy likelihoods achieved similar recovery performance with less overall consistency. The cauchit-normal was the exception from the general trend of similarity, as it performed considerably worse than all other likelihoods.

In terms of error rates, the main take away is that causal consistency was the most important aspect for FPR and TPR. Figure 8 shows that biased formulas had very high FPR and varying TPR compared to the unbiased formulas. While choice of likelihood and link mattered for individual DGPs, the only differences apparent when averaging

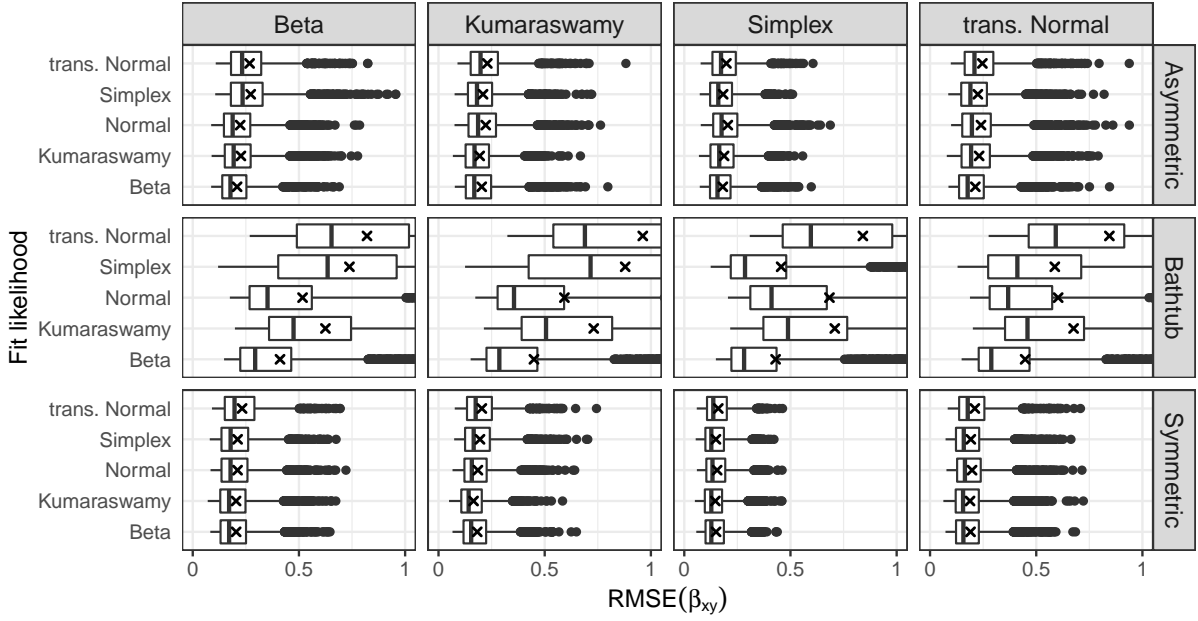


Figure 7: $\text{RMSE}(\beta_{xy})$ performance by data generating likelihood and shape for the logit link. Only models with fit link = data link are included. The x indicates the mean.

across all DGPs (as shown in Figure 8) were lower TPR for the simplex and cauchit-normal likelihoods. Notably, the normal-identity models had error rates similar to the well performing canonical likelihood and link combinations.

For $\Delta\text{ELPD}_{\text{loo}}$ we exemplarily present Figure 9 for data generated with a logit link. The PP of the different likelihoods and links was very similar in many cases but the model using the same likelihood and link as the true DGP was always among the best performing models. The beta likelihood again performed very consistent for all DGPs besides data from a combination of the simplex likelihood and bathtub shape which only the simplex likelihood could predict well. The logit-normal and Kumaraswamy likelihoods performed similarly to the beta, however with less consistency. The normal likelihood struggled with data from the asymmetric and bathtub shapes but achieved good $\Delta\text{ELPD}_{\text{loo}}$ for symmetric data. The simplex and cauchit-normal likelihoods had the worst performance outside of the self recovery scenarios. Similarly to the before mentioned utilities, the cloglog and cauchit link were not as reliable as the logit link regarding $\Delta\text{ELPD}_{\text{loo}}$ performance.

3.2. Lower-bounded Results

Similar to the double-bounded likelihoods, we found that many of the lower-bounded likelihoods performed similar across many of the metrics we calculated (see Table 4). As before, the average \hat{R} -values for the lower bounded likelihoods were far below the threshold of 1.01 for all combinations. The transformed-normal likelihoods showed the best sampling behaviour with no convergence problems and the highest sampling efficiency (measured by $\text{RESS}_{\text{bulk}}$ and $\text{RESS}_{\text{tail}}$) of all likelihoods. Close but not quite as good as the transformed-normal likelihoods were the beta prime, gamma and Weibull likelihoods. On the other end of the spectrum, the Gompertz and Fréchet had higher numbers of

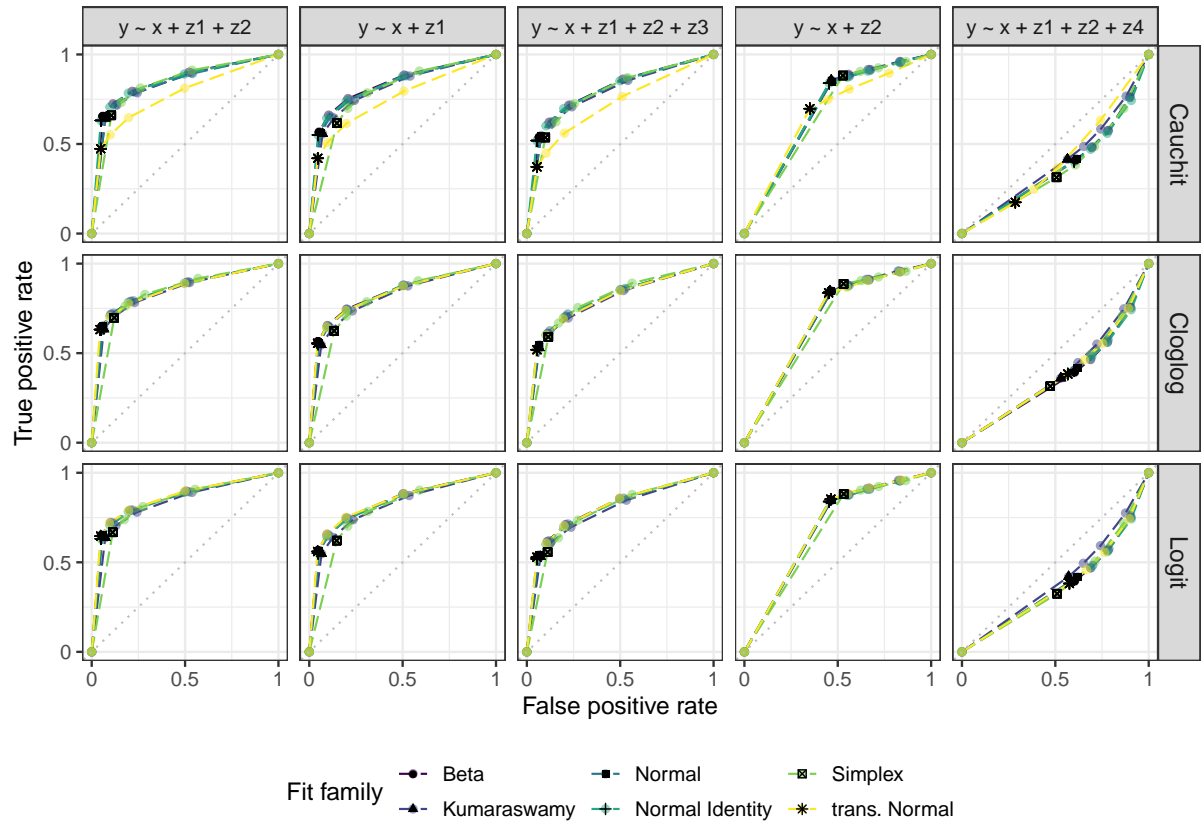


Figure 8: ROC for fit likelihoods by formula and fit link averaged over all double-bounded DGPs. The points are calculated from the 50%, 80%, 90% and 95% CIs, with the 95% CI highlighted in black.

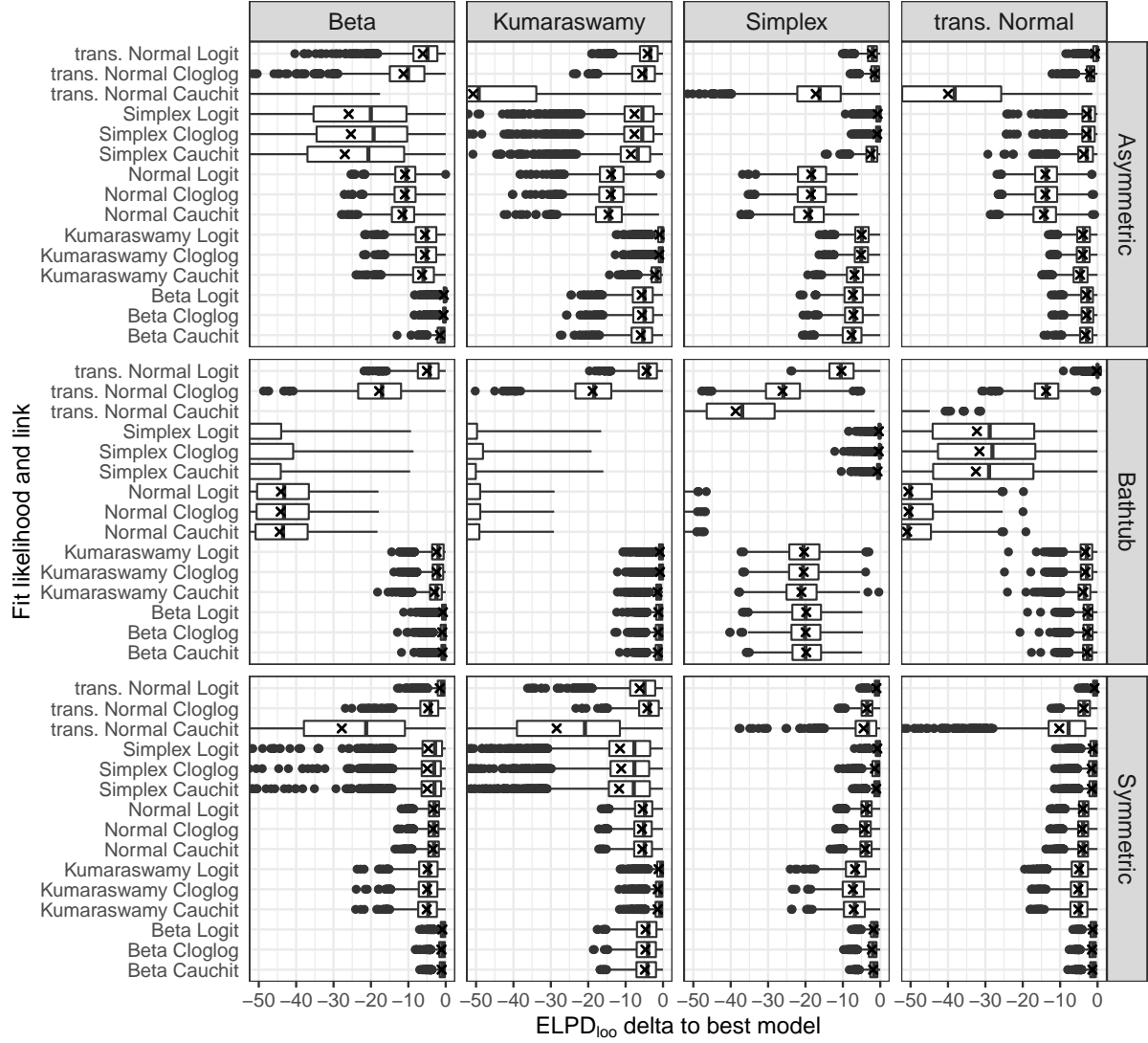


Figure 9: ΔELPD_{100} performance for fit likelihood and link combinations over data generating likelihoods and shapes for logit data. The data was truncated to a difference of 100 or less and the plot was truncated to a difference of 50 for clarity. The x indicates the mean.

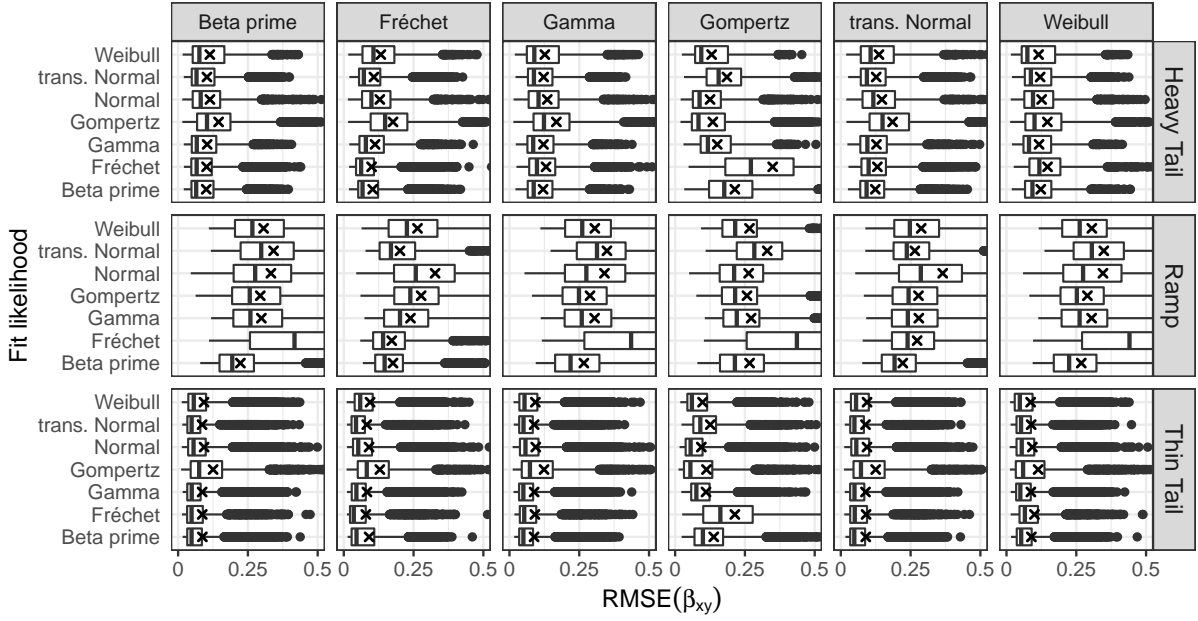


Figure 10: $\text{RMSE}(\beta_{xy})$ performance by data generating likelihood and shape for the log link. Only models with fit link = data link are included. The x indicates the mean.

DTs for the log and softplus links respectively. The normal likelihood was somewhat in the middle with higher DTs for the log than for the softplus link. Generally, the log link showed better sampling behaviour than softplus for all likelihoods besides the normal.

For PR measured via $\text{RMSE}(\beta_{xy})$ and $|\text{bias}(\beta_{xy})|$ the same patterns showed for both metrics. The most obvious observation was the stark difference between the ramp shape and both the thin- and heavy tail shapes for the log link as shown exemplarily in Figure 10 for the $\text{RMSE}(\beta_{xy})$ and log link. For the softplus link all three shapes behaved more similar to each other compared to the log link. Overall the best choice of likelihood depended on the true DGP. The Weibull, normal, transformed-normal, Gamma and beta prime likelihoods all performed well across many scenarios. The beta prime had the best median and the gamma the best mean values across all DGPs. Notable exceptions to those general trends were worse recovery for the beta prime likelihood on all Gompertz data and softplus-normal heavy tail data. In addition, both the normal and transformed-normal likelihoods had worse recovery on ramp shaped data. The Gompertz and Fréchet likelihoods generally lacked consistency and had worse recovery than the other likelihoods outside of a few favourable scenarios.

In terms of error rates, the main take-away is that causal consistency again was the most important aspect for FPR and TPR. Figure 11 shows that biased formulas had very high FPR and varying TPR compared to the unbiased formulas. While choice of likelihood and link mattered for individual DGPs, the only differences apparent when averaging across all DGPs (as shown in Figure 11) were lower TPR for the normal-log combination and higher FPR for the Fréchet likelihood. As for the double-bounded likelihoods, the normal-identity models had similar error rates to the well performing but structurally faithful lower-bounded likelihood and link combinations.

For $\Delta\text{ELPD}_{\text{loo}}$ we exemplarily present Figure 12 for data generated with a log link. The

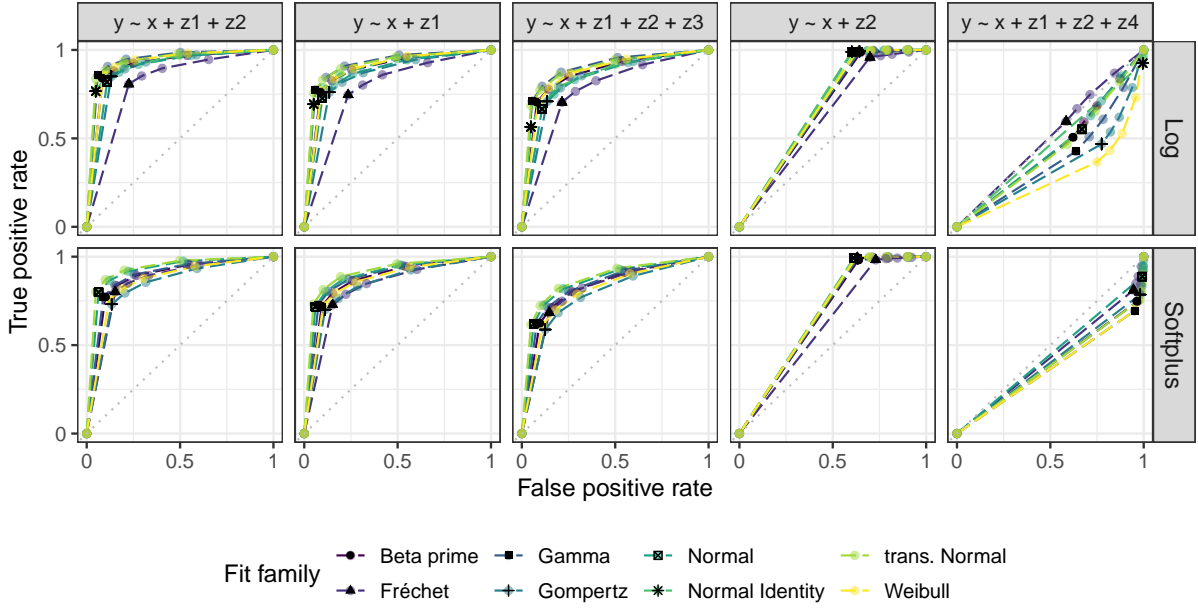


Figure 11: ROC for fit likelihoods by formula and fit link averaged over all lower-bounded DGPs. The points are calculated from the 50%, 80%, 90% and 95% CIs, with the 95% CI highlighted in black.

PP of the different likelihoods and links was very similar in many cases but the model using the same likelihood and link as the true DGP was among the best performing models as would be expected. For the log data, models fitted with a softplus link generally performed worse than models fitted with a log link. The ramp data stands out here as softplus models performed similar or sometimes even better than the log models. For softplus data, models with a softplus link performed only slightly better than models with a log link. The best ΔELPD_{100} performance was achieved by the gamma, log-normal, beta prime and Weibull likelihoods. Of those four, only the log-normal showed a big improvement when using the softplus link (i.e., the softplus-normal likelihood) for softplus data. The other three likelihoods showed little improvement from just using the log link on softplus data. This is interesting as we calibrated our DGPs to have latent means in the more linear parts of the softplus link, where we would expect bigger differences between the log and softplus. The normal, Fréchet, and Gompertz likelihoods all had considerably worse average ΔELPD_{100} values for both links outside of a few scenarios.

3.3. Predictive performance as Proxy for Parameter Recovery

We present results from the models discussed in Section 2.5 that predicted $\text{RMSE}(\beta_{xy})$ and zero-overlap of 95% CIs (giving us FPR and TPR depending on β_{xy}) with ΔELPD_{100} . We do not present $|\text{bias}(\beta_{xy})|$ results for brevity, as they show the same qualitative trends as for the $\text{RMSE}(\beta_{xy})$. See the online appendix [104] for $|\text{bias}(\beta_{xy})|$ results.

Figures 13 and 14 show the conditional effects of ΔELPD_{100} on $\text{RMSE}(\beta_{xy})$ for the double- and lower-bounded data, respectively, split by data generating link and shape. The most important observation is that all slopes of unbiased formulas for the thin tail,

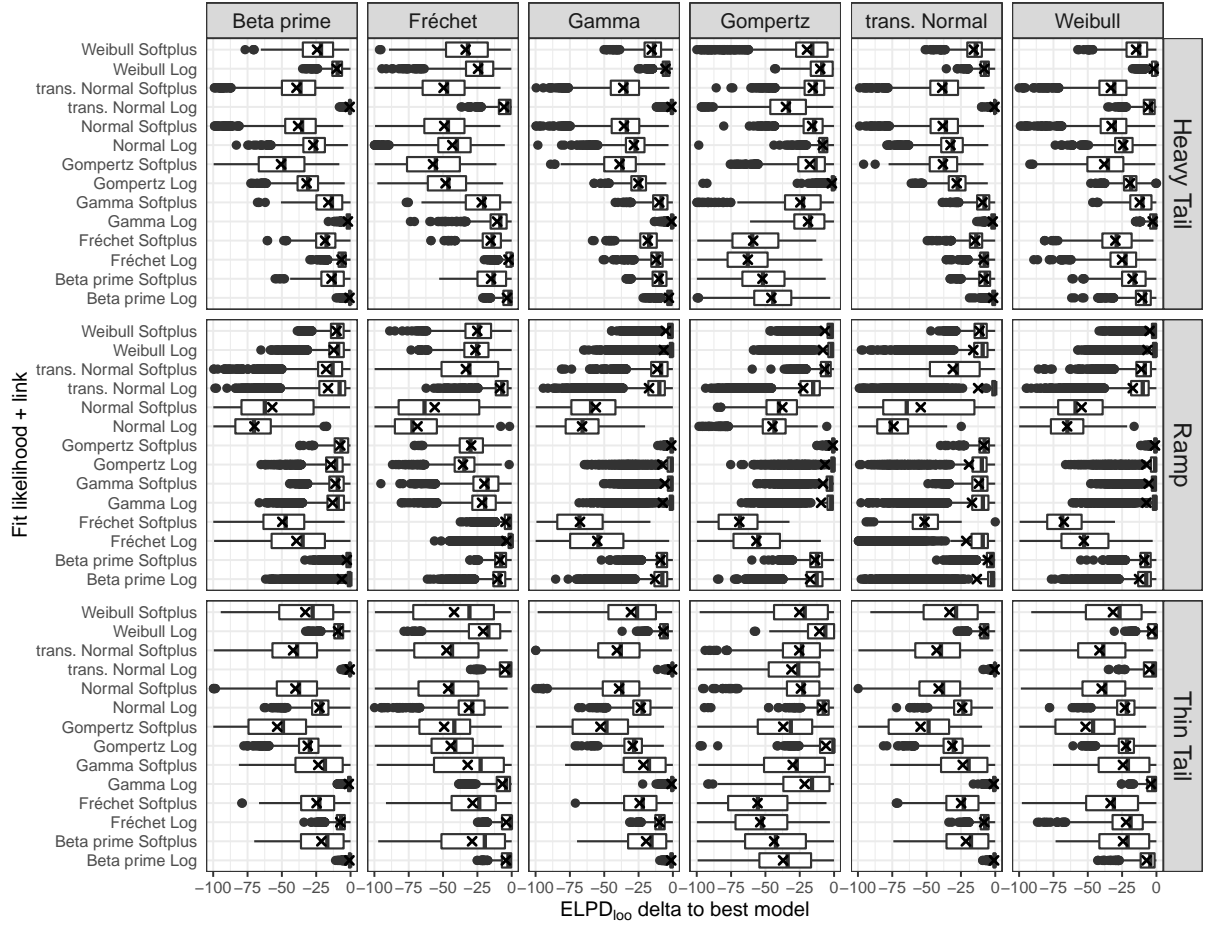


Figure 12: ΔELPD_{100} performance for fit likelihood and link combinations over data generating likelihoods and shapes for log data. The data was truncated to a difference of 100 or less and the plot was truncated to a difference of 50 for clarity. The x indicates the mean.

heavy tail and the ramp shapes are negative. This means that for a given dataset and fit formula, there was a clear trend of lower $\text{RMSE}(\beta_{xy})$ for models that ranked higher in ΔELPD_{100} . For the bathtub shape, ΔELPD_{100} performance did not predict $\text{RMSE}(\beta_{xy})$ for the unbiased formulas with slopes near zero. The slopes for the two biased formulas are inconsistent across shapes and links with both positive and negative cases.

Generally we can see that, given causal consistency (i.e., the three unbiased formulas) and only varying the fit likelihood (as fit link = true link here), we could rarely do harm by selecting models based on PP, even if when PR is the goal. In most cases, using ELPD_{100} performance as a proxy for PR turned out to be useful, as it reduced $\text{RMSE}(\beta_{xy})$ considerably. When using a causal inconsistent model, the main source of bias came from the causal inconsistency itself. Using ΔELPD_{100} for model selection did not have a reliable positive or negative effect on PR for the biased formulas.

Figures 15 and 16 show the conditional effects of ΔELPD_{100} on the zero-overlap of the 95% CIs for the double- and lower-bounded data, respectively, split by data generating link and shape. Depending on the true β_{xy} , this is either FPR ($\beta_{xy} = 0$) or TPR ($\beta_{xy} \neq 0$).

For the TPR, we can again see a clear trend in the thin and heavy tail shapes that, for all unbiased formulas, better ΔELPD_{100} predicts higher TPR. The unbiased slopes for the ramp shape already start with a TPR close to 1 but are still positive while the unbiased slopes for the bathtub shape stay very close to 0 but again still point upwards. We can also nicely see the higher precision of the true formula compared to the other two unbiased ones, as the true formula reached higher TPR earlier. For the biased formulas, there is no clear trend as they mostly stuck to one of the boundaries or had slopes with inconsistent signs.

For the FPR, the slopes for the unbiased formulas stayed close to 0 in all cases but the lower-bounded data with true softplus link. In the true softplus link case, the FPR decreased (improved) for growing ΔELPD_{100} , that is, there where badly predicting models using unbiased formulas that had substantially inflated FPR. The biased formulas either had positive slopes, that is, higher FPR with increasing ΔELPD_{100} , or stayed at one of the boundaries of 0 and 1.

4. Discussion

In this section we connect the results of our research to the main aims laid out in Section 1.5, discuss some limitations, and give an outlook how to possibly further advance the understanding of the posed questions.

4.1. *Using Prediction as a Proxy for Recovery*

Regarding the "predictive performance as proxy for parameter recoverability" (PP4PR) question, we observed one consistent trend over almost all investigated scenarios: When comparing models using the same causally unbiased formula, and in the case of $|\text{bias}(\beta)|$ and $\text{RMSE}(\beta)$ the same link, better ΔELPD_{100} was predictive of better PR for all investigated metrics. In the case of $|\text{bias}(\beta)|$ and $\text{RMSE}(\beta)$, these trends were over all likelihoods and in the case of FPR and TPR over all likelihood and link combinations. The trends were also consistent for individual datasets, as exemplarily shown in Figure 17 which indicates that the proxy can be reliably used in practice, where usually only a single dataset

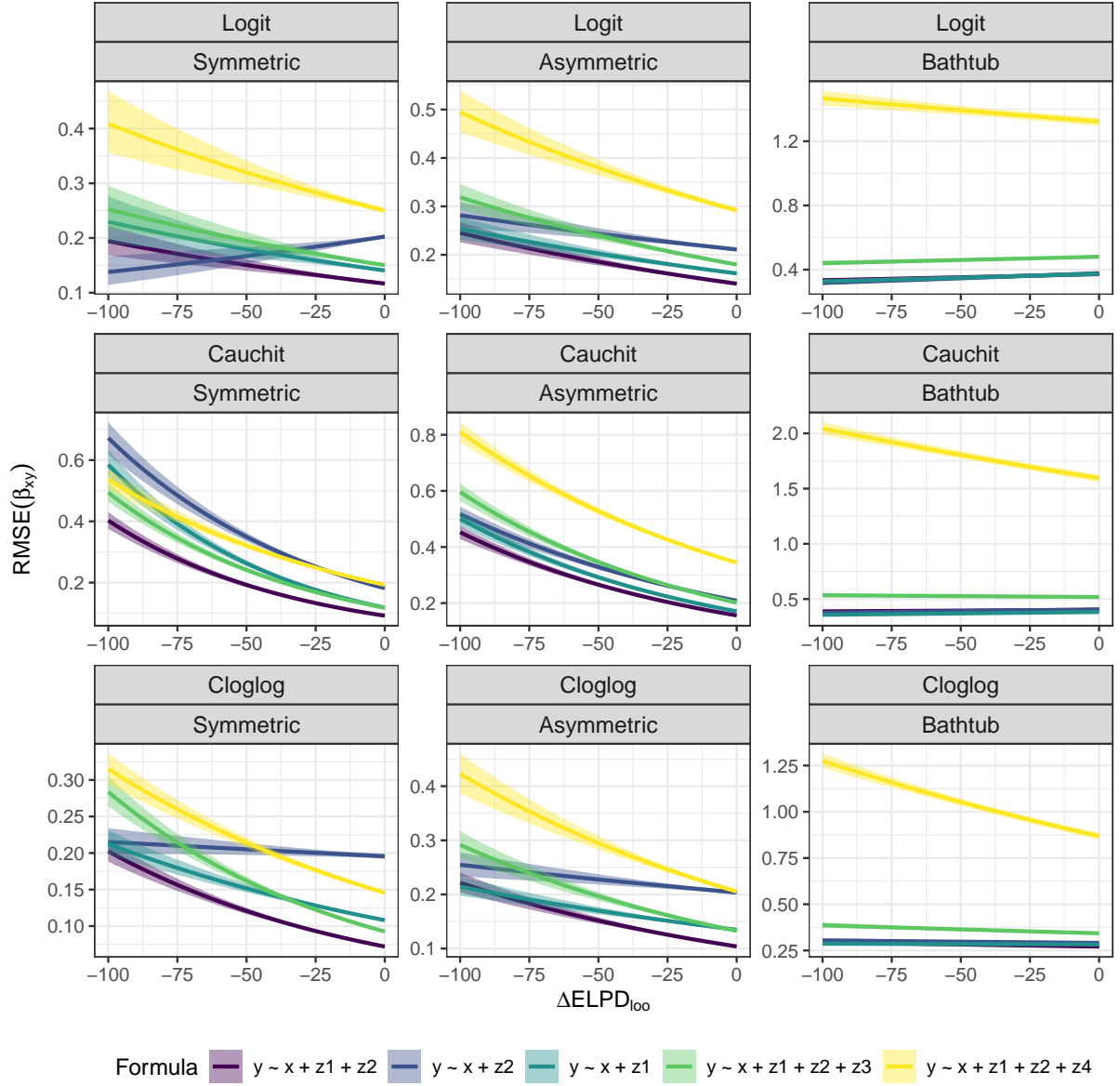


Figure 13: Conditional effects of ΔELPD_{100} on $\text{RMSE}(\beta_{xy})$ for double-bounded data and models. Split by data generating link and shape. See Section 2.5 for model details.

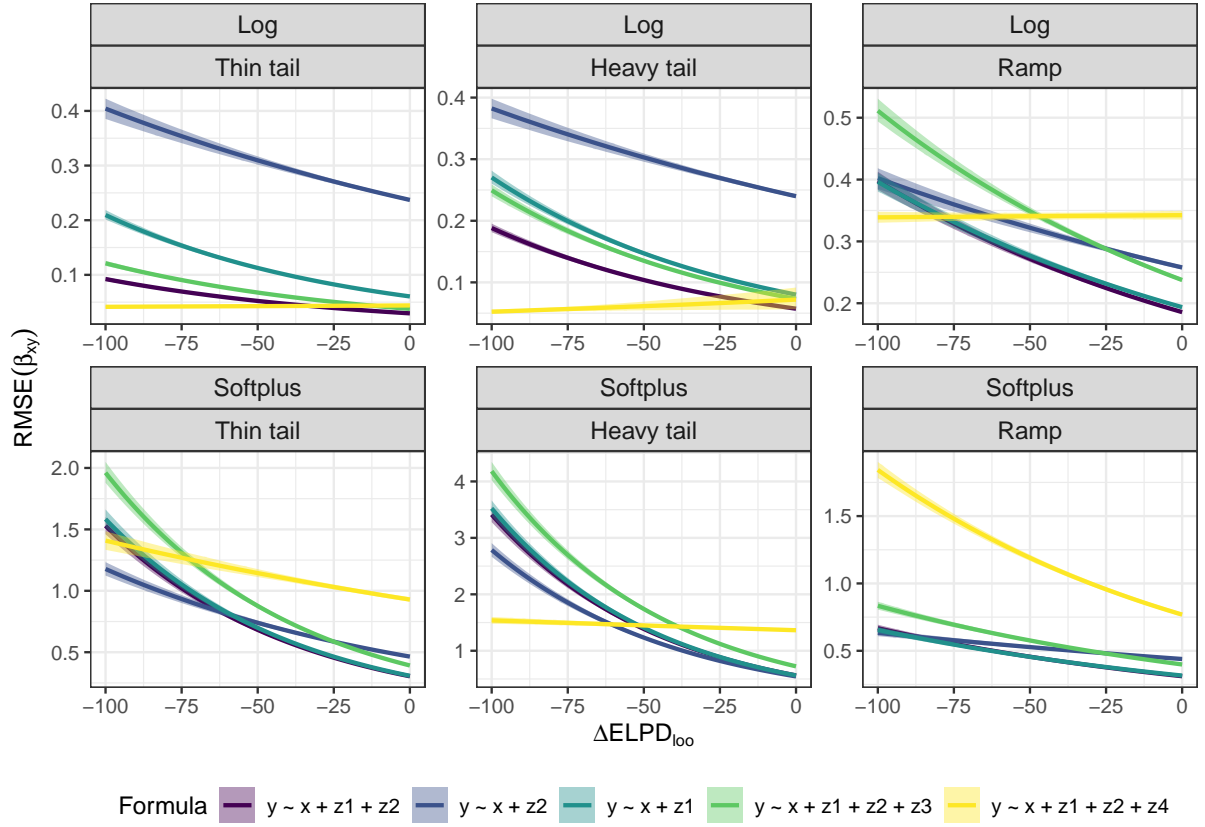


Figure 14: Conditional effects of $\Delta\text{ELPD}_{\text{loo}}$ on $\text{RMSE}(\beta_{xy})$ for lower-bounded data and models. Split by data generating link and shape. See Section 2.5 for model details.

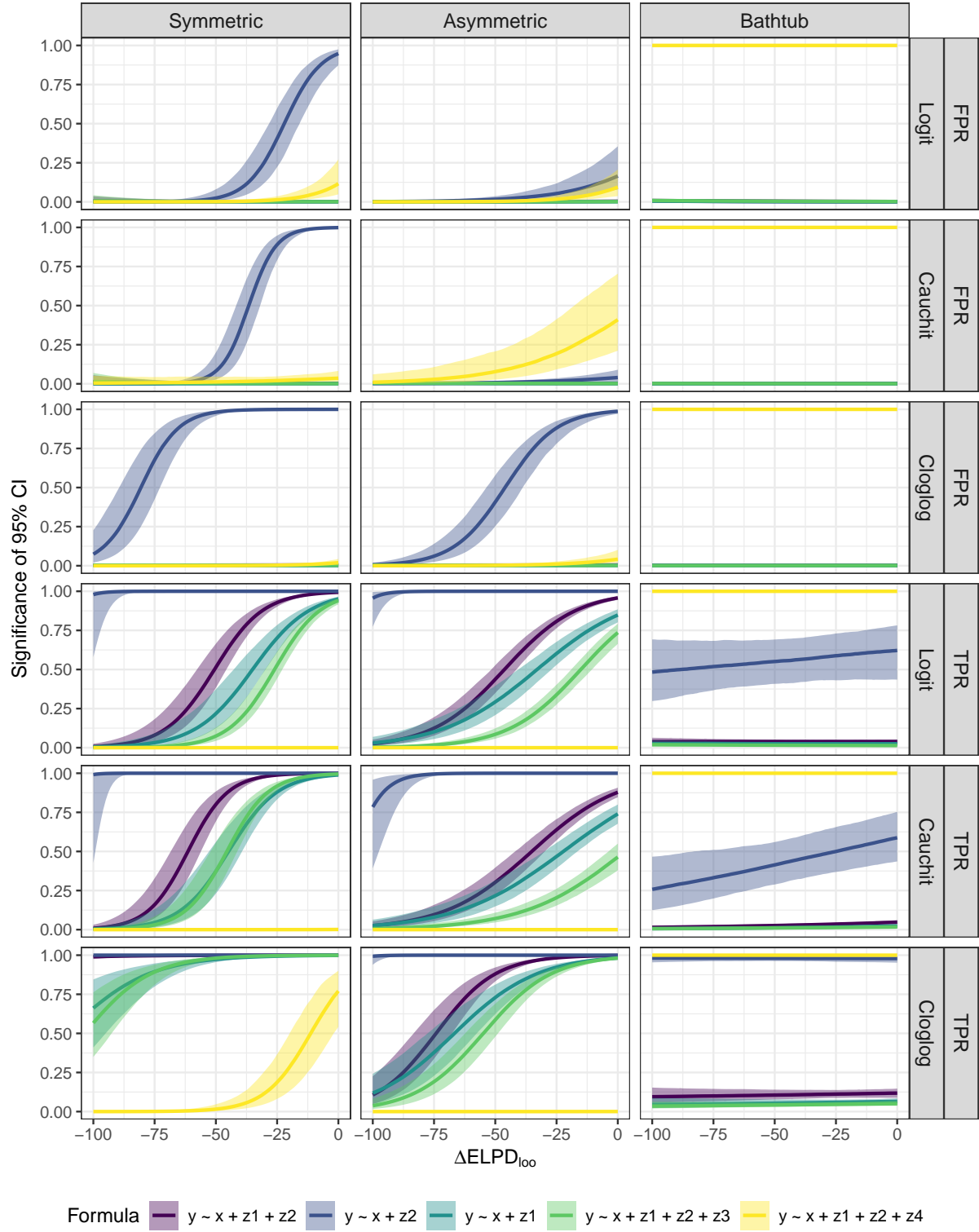


Figure 15: Conditional effects of ΔELPD_{100} on the zero overlap of the 95% CI (FPR and TPR depending on β_{xy}) for double-bounded data and models. Split by data generating link and shape. See Section 2.5 for model details.

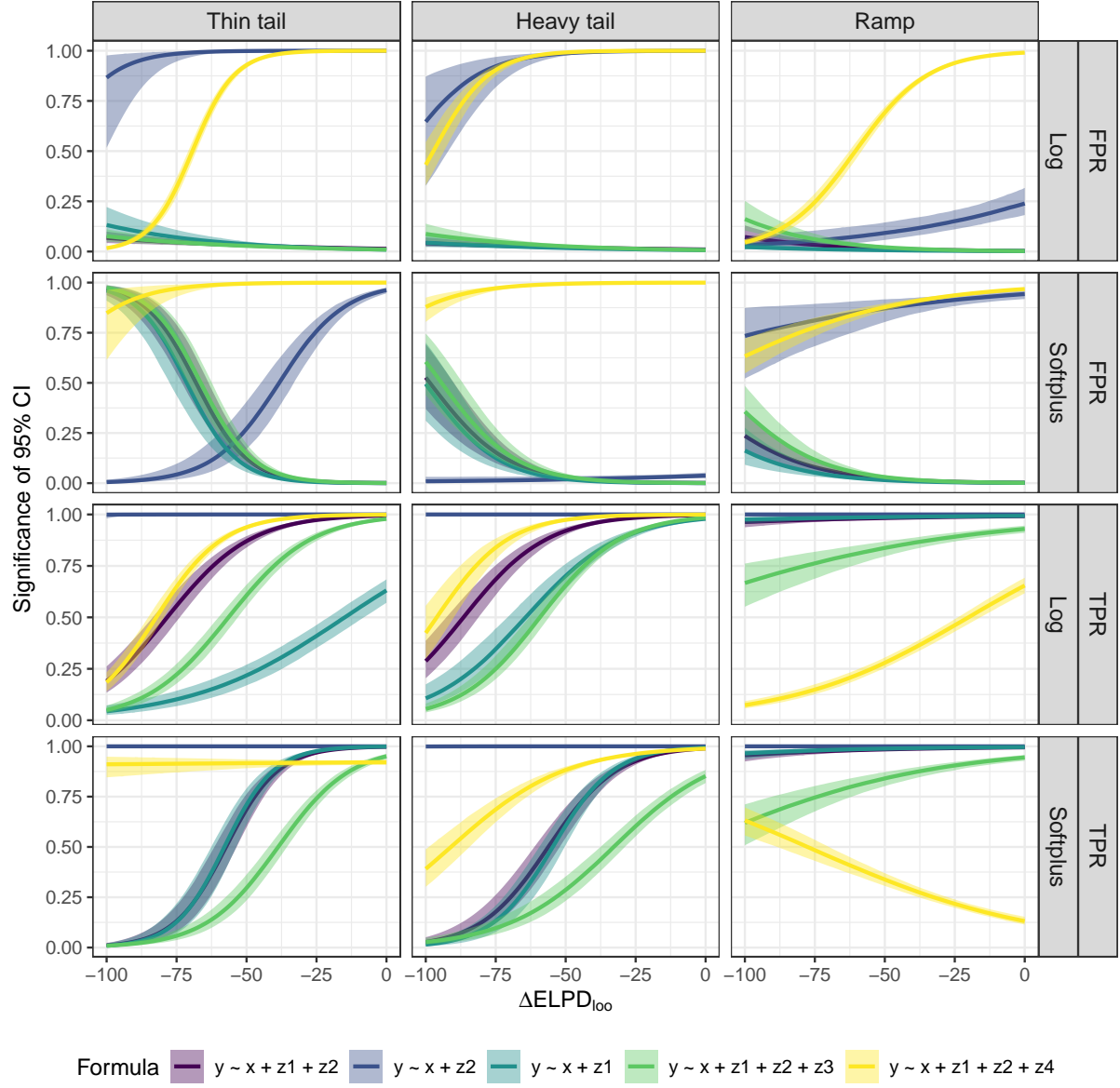


Figure 16: Conditional effects of ΔELPD_{100} on the zero overlap of the 95% CI (FPR and TPR depending on β_{xy}) for lower-bounded data and models. Split by data generating link and shape. See Section 2.5 for model details.

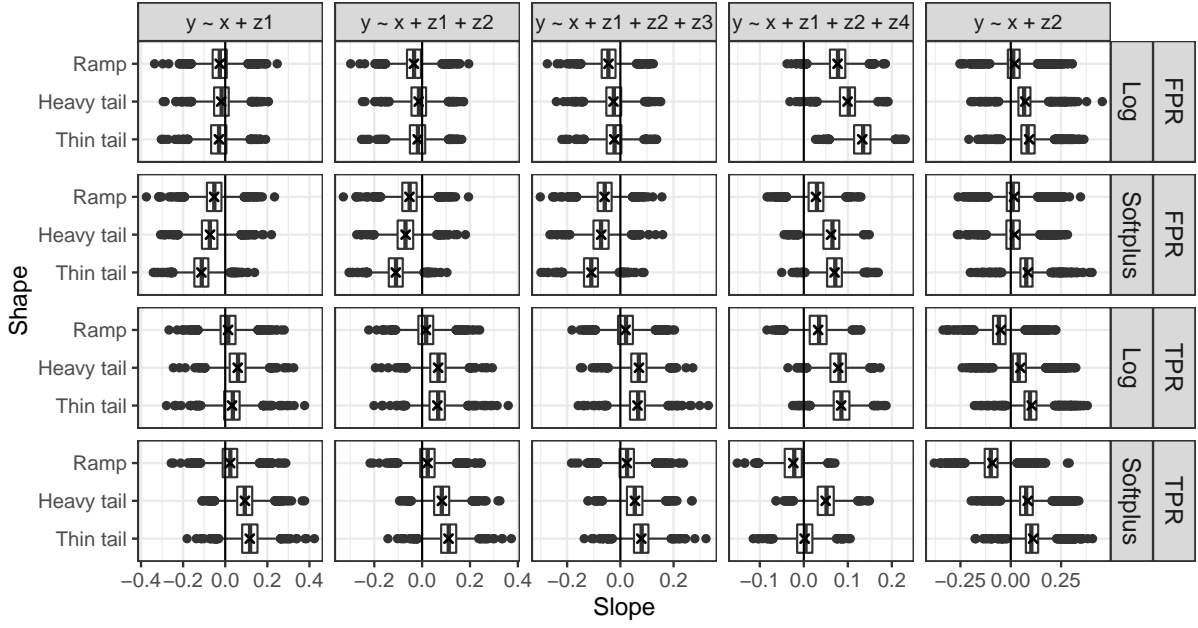


Figure 17: Slopes of individual datasets for the conditional effects of $\Delta\text{ELPD}_{\text{loo}}$ on the zero overlap of the 95% CI (FPR and TPR depending on β_{xy}) for lower-bounded data and models. Split by fit formula, data generating link and shape. See Section 2.5 for model details.

is available. In the few cases without a clear trend, for causally unbiased formulas, the slope was effectively flat. While better predicting models did not have better PR in those cases, using prediction as a proxy did also not reduce PR on average.

We conclude that, given a causally consistent model, PP can be safely used as a proxy for PR in model comparison for the kind of models we investigated, when the comparison is done for a fixed formula (i.e., for a fixed predictor structure). That said, causal consistency remains an even more important indicator for PR and little can be saved if the former is violated. This was evident in the generally worse recovery for the biased formulas and partially inverted slopes of PP predicting PR, that is, where better PP implied consistently worse recovery.

4.2. Default Recommendations

The second aim for this paper was to make practical recommendations for the choice of both likelihoods and link functions in the context of Bayesian GLMs for lower-bounded and double-bounded continuous data.

Double-bounded data. For double-bounded data, we found the beta likelihood with a logit link to be the most robust in terms of sampling behaviour, PR, and PP. While it is not the optimal solution in every scenario, as discussed in more detail in the online appendix [104], it was a good starting point for most scenarios. The logit-normal likelihood and Kumaraswamy likelihood with a logit link are alternatives worth considering. Both are less consistent than the beta likelihood overall but can perform better than the beta in some scenarios. They also offer median parameterizations as an alternative to the common

mean parameterization of the beta likelihood. The simplex, cauchit- and cloglog-normal likelihoods were all limited in their utility for datasets generated from other likelihoods, as were the cauchit and cloglog link in general. While it might be worthwhile to remember these options in case of bad performance of the aforementioned recommendations, their overall performance warrants caution.

Lower-bounded data. For lower-bounded data we found that the beta prime, gamma, transformed-normal, and Weibull likelihoods all performed well across our metrics with minimal differences. The biggest difference between those four likelihoods manifested during sampling, where the transformed-normal models were the most efficient. The Fréchet and Gompertz likelihoods had the worst performance across all metrics, besides a few favourable scenarios. The log link turned out to be more flexible than the softplus link in our scenarios, as log models could predict softplus data better than the softplus models could predict log data.

Normal models. The normal likelihood combined with appropriate links, while not a structurally faithful choice for bounded outcomes, showed good PR for both double- and lower-bounded data, similar to the better performing structurally faithful likelihoods. However, the PP of the normal likelihood was worse for non-symmetric data, such that we would not recommend it if prediction was the goal. This is due to a combination of the fact that the normal likelihood does not respect the outcome boundaries and is restricted to its symmetric shape. In addition to the recommendations above, we found that models using a normal likelihood with an identity link had similar FPR and TPR to the more canonical alternatives. This indicates that such normal-identity models can be valid alternatives if frequentist calibration is the only objective. This can be practically relevant especially as the estimation speed of normal-identity models can be magnitudes faster than other models, due to the availability of highly optimized implementations.

The good calibration of the normal-identity models is also noteworthy as they are a commonly used staple in many scientific fields even if not structurally faithful (see also Section 1.3). There, they are most often used in conjunction with null-hypothesis significance testing without any discussion of alternative likelihoods and links. Our results show that, at least for the simple GLMs we analyzed, this process can be viable and achieve reliable calibration. This is reassuring for the validity of many scientific results more generally as it shows that the practice of using linear regression can have good calibration even if it is not structurally faithful.

4.3. *Limitations of Scope and Outlook*

As is natural in simulation studies, their generalizability beyond the studied scenarios remains unclear. While an analytic investigation supplementing the simulations would have been desirable, the lack of available closed-form posteriors outside of a few conjugate model families prevented such an investigation. The limitations in generalizability apply specifically for the results regarding the relationship between PP and PR, and touch essentially every aspect of our study.

We have no concrete reason to expect other likelihood classes (e.g., for unbounded continuous or count data) to show qualitatively different behaviour regarding this rela-

tionship. However, we cannot rule it out based on the available evidence. One cautionary observation we made was the big influence of different likelihood shapes on both the general performance of a likelihood or link and on the PP4PR question. The mechanisms by which different shapes imply different downstream results, especially with regard to the latter relationship, remain unclear to us to some degree. An extension of this work that would investigate different data generating shapes more systematically, their properties, and how to best deal with them during modeling workflows could improve said understanding.

With regard to causal assumptions, we chose our DGPs and the four misspecified formulas in an attempt to span the most important classes of controls presented in [29]. Still, there are many more possible DGPs and misspecifications that we did not include, for example, unobserved variables or variables with measurement error [93]. An extension of this work to span different kinds of DGPs and misspecifications would ultimately add to the generalizability of results. That said, we currently do not see how a causally unbiased formula would have to look like in order to render the relationship of PP and PR negative (i.e., better PP implied worse PR), as long as all compared models shared the same (unbiased) formula.

As discussed in Section 1.3, we limited the scope of this paper to models with simple linear predictor terms. We would expect to observe the same general trends for more complicated additive predictor terms where individual components may be non-linear in the model parameters as the DAG-based misspecifications are agnostic to the model structure [94]. However, extending our work in this direction may be interesting insofar as to not only allow to include or exclude certain variables but also to study the implications of (mis-)specifying their individual terms (e.g., modeling an effect as linear while it is truly non-linear in some way). This would more strongly play into the question of overfitting in sparse scenarios with (relatively) small data yet complicated non-linear relationships to be approximated. Similar questions would arise when adding multilevel structure especially in sparse data scenarios where the (true) relevance of multilevel terms had to be balanced against their weak identification implied by the given data [23].

In this study, we followed a Bayesian perspective on model building although with flat priors. The latter choice was primarily made to avoid some models having an "unfair" advantage due to incidentally more suitable prior choices (see Section 2.3). As a side effect of this choice, we think that our main results are likely to hold as well in case of frequentist (maximum likelihood) estimation of all models, and corresponding model-based metrics. We see it as unlikely that our likelihood-link recommendations would have changed in the light of applying (weakly-)informative priors [1, 82, 43]. However, there may be one subtle place where priors (or regularization more generally) may influence the relationship of PP and PR: Cross-validation (CV) estimates of out-of-sample PP produce correlated folds due to overlap in training data and have an intricate relationship with different true quantities they can be seen to approximate [11]. These properties can be altered via regularization, especially in sparse regimes where the number of model parameters is high relative to the number of observations in the data [11]. As such the relationship of CV estimates with PR may also vary in such scenarios, which would be interesting to study in the future.

Acknowledgements

This work was partially funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) under Germany’s Excellence Strategy – EXC-2075 - 390740016 (the Stuttgart Cluster of Excellence SimTech).

This work was performed on the computational resource bwUniCluster funded by the Ministry of Science, Research and the Arts Baden-Württemberg and the Universities of the State of Baden-Württemberg, Germany, within the framework program bwHPC.

Data Availability Statement

The data that support the findings of this study are openly available in OSF at <http://doi.org/10.17605/OSF.IO/XGKZV>.

Disclosure Statement

The authors report there are no competing interests to declare.

A. *Review of Likelihoods and Link Functions*

In this section, we present an overview of likelihoods and link functions used in GLMs of continuous lower-bounded and double-bounded data. This is a subset of all options that our literature search revealed, obtained from the inclusion criteria discussed below. The full overview can be found in our online appendix [104].

Inclusion criteria for likelihoods to be used in our simulation study were as follows. First, we only included distributions that we could find to have been previously applied as likelihoods in regression models. Second, we required them to either have a mean or median parameterization, such that the parameter μ is always the mean or median. We used mean parameterizations when available and median parameterizations otherwise. Third, we only included likelihoods that are no special cases of other to-be-included likelihoods (e.g., as is exponential a special case of both Gamma and Weibull), as their comparisons would be trivial given sufficient data. Fourth, we only considered likelihoods with two distributional parameters, that is one auxiliary parameter in addition to μ . This facilitates achieving the third inclusion criteria, without excluding many practically relevant likelihoods: Three or more-parameter likelihoods appear to be rare in practical applications of the considered GLM classes. According to our own experience and the reports of users of software we maintain, one likely reason is that their estimation tends to be much harder, showing more convergence problems and often requiring additional regularization (see [62, 26] for examples).

Inclusion criteria for link functions were somewhat stricter in that we only included widely applied links or, for less common links, only those with qualitatively different shapes. What is more, we only considered links that are fixed, fully known functions. This means that we did not include parametric links that have additional parameters learned from data during model fitting [24, 102, 103, 69, 57, 75]. The reasons for this choice are three-fold. First, parametric links increase the model complexity by adding

further distributional parameters to the model, in a way similar to using more complex likelihoods, with the same drawbacks and benefits. Second, the scale of the linear predictor changes with different values of the link function’s parameters making interpretation harder, specifically in a Bayesian setting where linear predictor’s scale would accordingly vary across posterior samples. Third, we haven’t found them to be commonly applied in practice, which might partially be due to the lack of general-purpose software support. While parametric link functions are an important development in the context of GLMs, we feel that the present paper is not the right place to study them.

A.1. *Double-Bounded Likelihoods*

Based on the discussed inclusion criteria, we included six distributions for double-bounded data in our analysis as detailed below. Figure 2 shows some exemplary densities for each of them, illustrating qualitatively different kinds of shapes they can accommodate. The three distinct shapes are uni-modal symmetric and asymmetric shapes as well as a bi-modal bathtub shape. We refer to these shapes as symmetric, asymmetric, and bathtub, respectively.

Beta The beta distribution is common and thoroughly studied [71, 50, 35], often considered as *the* distribution for unit interval regression [38, 30], and regularly used as a baseline to compare other unit interval distributions against [66, 15, 77]. We use a common mean parameterization, given by

$$f(y | \mu, \phi) = \frac{1}{B(\mu\phi, (1-\mu)\phi)} y^{\mu\phi-1} (1-y)^{(1-\mu)\phi-1},$$

with $0 < \mu < 1, \phi > 0$, mean $\mathbb{E}[y] = \mu$, and variance $\text{Var}(y) = \frac{\mu(1-\mu)}{\phi+1}$.

Kumaraswamy The Kumaraswamy distribution originated in the field of hydrology [72] and is a special case of the generalized beta distribution [90]. Compared to the beta distribution, its closed forms for the PDF, CDF and quantile functions are computationally cheaper, which seems to be the reason its popularity [60, 90]. While not as common as the beta distribution, the Kumaraswamy distribution has been studied in-depth [39, 60, 86], and is used in regression applications [87, 52, 97]. We use a median parameterization proposed by [87]

$$f(y | \mu, p) = \frac{p \log(0.5)}{\log(1-\mu^p)} y^{p-1} (1-y^p)^{\frac{\log(0.5)}{\log(1-\mu^p)}-1},$$

with $0 < \mu < 1, p > 0$, mean $\mathbb{E}[y] = qB(1 + \frac{1}{p}, q)$, and variance $\text{Var}(y) = qB(1 + \frac{2}{p}, q) - [qB(1 + \frac{1}{p}, q)]^2$ where $q = \frac{\log(0.5)}{\log(1-\mu^p)}$.

Simplex The simplex distribution was proposed by [8]. It seems to be mostly used as a comparison to the beta and sometimes Kumaraswamy distribution [77, 66] and is so far rarely used for regression [77, 9]. We use the mean parameterization proposed by [109]

$$f(y | \mu, \sigma^2) = [2\pi\sigma^2\{y(1-y)\}^3]^{-\frac{1}{2}} \exp\left(-\frac{1}{2\sigma^2}d(y; \mu)\right),$$

with $0 < \mu < 1$ and mean $\mathbb{E}[y] = \mu$. We omit the variance formula for brevity but interested readers can find it in [61, 110].

Transformed Normal One option to create novel distributions is the use of data-transformations. The most common transformation-based distribution for unit interval values is the logit-normal distribution [7], the distribution of a variable whose logit follows a normal distribution [68]. It has been studied thoroughly, [83, 3, 41] and is commonly used in regression [119, 6, 96, 114]. We use the standard median parameterization

$$f(y | \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \frac{1}{y(1-y)} \exp\left(-\frac{(\text{logit}(y) - \mu)^2}{2\sigma^2}\right).$$

In addition to the logit-normal we also use the cauchit- and cloglog-normal distributions that work similar to the logit-normal but uses the respective link functions for transformations (see Appendix A.2 for justification of those three options). The density for the cauchit-normal is

$$f(y | \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \pi \left(\frac{1}{\cos(\pi(y - \frac{1}{2}))} \right)^2 \exp\left(-\frac{(\text{cauchit}(y) - \mu)^2}{2\sigma^2}\right),$$

and the density for the cloglog-normal is

$$f(y | \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \frac{1}{\log(1-y)(y-1)} \exp\left(-\frac{(\text{cloglog}(y) - \mu)^2}{2\sigma^2}\right).$$

For all three transformed normal distributions, the median is equal to $\text{med}(y) = \text{inv_link}(\mu)$, where inv_link is the inverse of the respective transformation. No analytical solutions for mean or variance are available for these distributions.

A.2. Double-Bounded Link Functions

In terms of link functions for double-bounded variables, the symmetric logit link

$$\text{logit}(x) = \log(x) - \log(1-x) \quad \text{inv_logit}(x) = \frac{1}{1 + \exp(-x)},$$

symmetric probit link

$$\text{probit}(x) = \frac{1}{2} \left(1 + \text{erf}\left(\frac{x}{\sqrt{2}}\right) \right) \quad \text{inv_probit}(x) = \sqrt{2} \text{erf}^{-1}(2x - 1),$$

where erf is the error function [58], and the right-skewed cloglog link

$$\text{cloglog}(x) = \log(-\log(1-x)) \quad \text{inv_cloglog}(x) = 1 - \exp(-\exp(x)),$$

are, in that order, by far the most common choices [135, 57, 75, 31, 36, 47, 98]. Since logit and probit yield almost indistinguishable results due to the similar shapes of logistic and normal distribution [36, 47, 98], we decided to exclude the probit link from our simulations. Instead, as a third link function, we include the symmetric cauchit link

$$\text{cauchit}(x) = \tan(\pi(x - 0.5)) \quad \text{inv_cauchit}(x) = \frac{1}{\pi} \arctan(x) + 0.5,$$

as it has much wider tails than logit and probit and thus may show qualitatively different behavior [88, 69, 75]. Figure 3 illustrates the included link functions and their corresponding response functions.

A.3. Lower-Bounded Likelihoods

Based on the discussed inclusion criteria, we included eight distributions for lower-bounded data in our analysis as detailed below. Figure 4 shows some exemplary densities for each of them, illustrating qualitatively different kinds of shapes they can accommodate. The three distinct shapes are uni-modal thin tail and heavy tail shapes as well as a ramp shape. We refer to these shapes as thin tail, heavy tail, and ramp respectively.

Gamma The gamma distribution is one of the most common lower-bounded distributions with application, for example, in hydrology [4], meteorology [116, 129], medicine [12], and inventory control [18]. We use a mean parameterization as suggested by [65]

$$f(y | \mu, \alpha) = \frac{\left(\frac{\alpha}{\mu}\right)^\alpha}{\Gamma(\alpha)} y^{\alpha-1} \exp\left(-y \frac{\alpha}{\mu}\right),$$

where Γ is the gamma function, $\mu > 0, \alpha > 0$, mean $\mathbb{E}(y) = \mu$, and variance $\text{Var}(y) = \frac{\mu^2}{\alpha}$.

Weibull The Weibull distribution [101] is commonly used in survival analysis [137], to model failure data [73] and in reliability research in general [89]. We use a mean parameterization as implemented in [19]

$$f(y | \mu, k) = \frac{k}{\lambda} \left(\frac{y}{\lambda}\right)^{k-1} \exp\left(-\left(\frac{y}{\lambda}\right)^k\right),$$

with $\mu > 0, k > 0$, mean $\mathbb{E}[y] = \mu$, and variance $\text{Var}(y) = \lambda^2 \left[\Gamma\left(1 + \frac{2}{k}\right) - \left(\Gamma\left(1 + \frac{1}{k}\right)\right)^2\right]$, where $\lambda = \frac{\mu}{\Gamma\left(1 + \frac{1}{k}\right)}$.

Fréchet The Fréchet distribution is also known as the inverse Weibull distribution [32]. It is used for regression in Hydrology [49, 136] and reliability modeling [64]. We use a mean parameterization following [100]

$$f(y | \mu, \nu) = \frac{\nu}{s} \left(\frac{y}{s}\right)^{-1-\nu} \exp\left(-\left(\frac{y}{s}\right)^{-\nu}\right),$$

with $\mu > 0, \nu > 1$, mean $\mathbb{E}[y] = \mu$, and variance $\text{Var}(y) = s^2 \left(\Gamma\left(1 - \frac{2}{\nu}\right) - \left(\Gamma\left(1 - \frac{1}{\nu}\right)\right)^2\right)$, where $s = \frac{\mu}{\Gamma\left(1 - \frac{1}{\nu}\right)}$ for $\nu > 2$ and $\text{Var}(y) = \infty$ otherwise.

Inverse Gaussian The inverse Gaussian distribution is also known as the Wald distribution. It is used for regression in many different fields [106] and used for different applications like lifetime and reliability models [28], especially in cases where there is a burn-in period [40]. It also found use in modeling word frequencies and water reservoir levels [27]. An advantage over similar distributions is the especially tractable sampling theory of the inverse Gaussian [40, 106]. We use a mean parameterization

$$f(y | \mu, \lambda) = \sqrt{\frac{\lambda}{2\pi y^3}} \exp\left(-\frac{\lambda(y - \mu)^2}{2\mu^2 y}\right),$$

with $\mu > 0, \lambda > 0$, mean $\mathbb{E}[y] = \mu$, and variance $\text{Var}(y) = \frac{\mu^3}{\lambda}$.

Beta Prime The beta prime distribution [63, 80] has been used only little in regression so far [118, 84, 81], but still satisfies our inclusion criteria. We use a mean parameterization following [16]

$$f(y | \mu, \phi) = \frac{y^{\mu(\phi+1)-1}(1+y)^{-(\mu(\phi+1)+\phi+2)}}{B(\mu(\phi+1), \phi+2)},$$

with $\mu > 0, \phi > 0$, mean $\mathbb{E}[y] = \mu = \frac{\alpha}{\beta-1}, \beta > 1$, and variance $\text{Var}(x) = \frac{\alpha(\alpha+\beta-1)}{(\beta-2)(\beta-1)^2}, \beta > 2$, where B is the beta function, $\alpha = \frac{\mu}{\phi+1}$ and $\beta = \phi + 2$.

Gompertz The Gompertz distribution is used in fields connected to research of mortality, life expectancy and incidence rates [48, 55] and applied to regressions in those fields [67, 56]. We use a median parameterization

$$f(y | \mu, \eta) = \eta b \exp(\eta + by - \eta \exp(by)),$$

with $\mu > 0, \eta > 0$, where $b = \frac{\log(1 - \frac{1}{\eta} \log(0.5))}{\mu}$. We omit the mean and variance for brevity here but interested readers can find them here [58].

Transformed Normal As for the unit interval distributions, we use transformed normal distributions with the lower-bounded links as well. This leads to the log-normal distribution, that is used in mechanical testing [117], risk analysis [20], hydrology [70], and ecology [33]. While there exists a mean parameterization, we use a much more commonly used median parameterization

$$f(y | \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \frac{1}{y} \exp\left(-\frac{(\log(y) - \mu)^2}{2\sigma^2}\right),$$

with $\mu \in \mathbb{R}, \sigma > 0$, median $\text{med}(y) = \exp(\mu)$, mean $\mathbb{E}[y] = \exp\left(\mu + \frac{\sigma^2}{2}\right)$, and variance $\text{Var}(y) = (\exp(\sigma^2) - 1) \exp(2\mu + \sigma^2)$.

The softplus-normal likelihood [128] has density

$$f(y | \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \frac{\exp(x)}{\exp(x) - 1} \exp\left(-\frac{(\text{softplus}(y) - \mu)^2}{2\sigma^2}\right),$$

with $\mu \in (-\infty, +\infty), \sigma > 0$, and median $\text{med}(y) = \text{inv_softplus}(\mu)$, but no available analytic mean or variance.

A.4. Lower-Bounded Link Functions

The default recommendation, and almost exclusively applied link function for lower-bounded variables is the log link [37] with inverse function $\text{inv_log}(x) = \exp(x)$. The log link implies a multiplicative relationship on the response scale, which is a sensible assumption for positive data, at least close to zero, as additivity in this range would imply impossible negative values. The implied exponential growth of the corresponding response function can however be problematic in practice, especially for predictions. For this reason, the softplus link

$$\text{softplus}(x) = \log(\exp(x) + 1) \quad \text{inv_softplus}(x) = \log(\exp(x) + 1)$$

was proposed as an alternative [138, 34]. It behaves almost linear for large enough values, and thus allows for quasi-additive interpretation, while it behaves like the log link for small values to prevent predictions to exceed the lower boundary. Figure 5 illustrates the included link functions and their corresponding response functions.

B. Descriptive statistics tables

B.1. *Double-bounded Likelihoods*

Table 3: Descriptive Statistics of Double-bounded Results

Link	Likelihood	DT	$P(\text{Conv.})$	$\text{RESS}_{\text{bulk}}$	$\text{RESS}_{\text{tail}}$	$\text{TESS}_{\text{total}}^{\text{bulk}}$
Beta	Cauchit	0 (0)	1	0.84 (0.18)	0.7 (0.073)	0.001 (0.001)
Beta	Cloglog	0 (0)	1	0.85 (0.17)	0.72 (0.066)	0.002 (0.002)
Beta	Logit	0 (0)	1	0.87 (0.17)	0.73 (0.065)	0.14 (1.4)
Kumaraswamy	Cauchit	0 (0.026)	1	0.73 (0.18)	0.62 (0.11)	0.001 (0.001)
Kumaraswamy	Cloglog	7.9 (105)	0.97	0.81 (0.17)	0.69 (0.066)	0.001 (0.001)
Kumaraswamy	Logit	0 (0.004)	1	0.78 (0.17)	0.67 (0.068)	0.001 (0.001)
Normal	Cauchit	0.012 (0.83)	0.99	0.78 (0.2)	0.63 (0.12)	0.003 (0.17)
Normal	Cloglog	0.13 (2.1)	0.95	0.83 (0.18)	0.68 (0.088)	0.62 (2.9)
Normal	Logit	0.018 (0.83)	0.99	0.84 (0.19)	0.69 (0.088)	0.03 (0.38)
Simplex	Cauchit	0.83 (14)	0.99	0.81 (0.2)	0.69 (0.087)	0.001 (0)
Simplex	Cloglog	22 (106)	0.88	0.77 (0.21)	0.67 (0.099)	0.001 (0)
Simplex	Logit	0.46 (5.4)	0.99	0.8 (0.2)	0.7 (0.08)	0.002 (0.001)
trans. Normal	Cauchit	0 (0)	1	0.88 (0.18)	0.72 (0.065)	0.11 (1.1)
trans. Normal	Cloglog	0 (0)	1	0.88 (0.18)	0.72 (0.065)	0.001 (0.001)
trans. Normal	Logit	0 (0)	1	0.88 (0.18)	0.72 (0.065)	0.001 (0)

Link	Likelihood	PKs	ELPD_{loo}	$\text{bias}(\beta_{xy})$	$\text{SD}(\beta_{xy})$	$\text{RMSE}(\beta_{xy})$
Beta	Cauchit	0.024 (0.16)	53 (28)	0.21 (0.28)	0.14 (0.06)	0.28 (0.27)
Beta	Cloglog	0.021 (0.14)	51 (29)	0.34 (0.59)	0.23 (0.21)	0.44 (0.61)
Beta	Logit	0.015 (0.12)	46 (36)	0.27 (0.43)	0.19 (0.15)	0.35 (0.44)
Kumaraswamy	Cauchit	0.25 (0.52)	48 (32)	0.38 (0.71)	0.23 (0.32)	0.47 (0.77)
Kumaraswamy	Cloglog	0.23 (0.47)	11 (81)	0.67 (1.2)	0.54 (0.78)	0.93 (1.4)
Kumaraswamy	Logit	0.34 (0.56)	53 (28)	0.16 (0.18)	0.1 (0.04)	0.2 (0.17)
Normal	Cauchit	0.015 (0.13)	52 (29)	0.18 (0.22)	0.14 (0.09)	0.24 (0.22)
Normal	Cloglog	0.012 (0.11)	46 (36)	0.2 (0.28)	0.13 (0.07)	0.25 (0.27)
Normal	Logit	0.0054 (0.074)	51 (28)	0.19 (0.22)	0.11 (0.05)	0.23 (0.21)
Simplex	Cauchit	0.88 (1.4)	51 (33)	0.25 (0.35)	0.18 (0.13)	0.34 (0.35)
Simplex	Cloglog	0.96 (1.4)	53 (28)	0.21 (0.25)	0.15 (0.06)	0.27 (0.24)
Simplex	Logit	0.91 (1.3)	52 (29)	0.27 (0.38)	0.19 (0.12)	0.36 (0.37)
trans. Normal	Cauchit	0.73 (0.65)	46 (36)	0.26 (0.37)	0.18 (0.1)	0.34 (0.36)
trans. Normal	Cloglog	0.068 (0.26)	49 (32)	0.29 (0.38)	0.17 (0.1)	0.36 (0.37)
trans. Normal	Logit	0.052 (0.23)	53 (28)	0.33 (0.46)	0.23 (0.16)	0.43 (0.46)

Note: Most results are presented as mean(sd). $P(\text{Conv.})$ is calculated as a proportion of all models. $\text{TESS}_{\text{total}}^{\text{bulk}}$ is also calculated using all models. Both the metrics for predictive performance and parameter recoverability are calculated only for models that had 10 or less divergent transitions and converged. In addition the parameter recoverability metrics are only shown for models using the true link.

B.2. Lower-bounded Likelihoods

Table 4: Descriptive Statistics of Lower-bounded Results

Link	Likelihood	DT	$P(\text{Conv.})$	RESS _{bulk}	RESS _{tail}	TESS _{total} ^{bulk}
Beta prime	Log	0.012 (2.2)	1	0.75 (0.18)	0.66 (0.082)	0.002 (0.018)
Beta prime	Softplus	0.004 (0.096)	1	0.69 (0.16)	0.64 (0.081)	0.018 (0.57)
Fréchet	Log	0.73 (27)	1	0.63 (0.16)	0.59 (0.089)	0.001 (0.065)
Fréchet	Softplus	298 (580)	0.72	0.58 (0.15)	0.56 (0.12)	0.72 (6.2)
Gamma	Log	0.02 (5.3)	1	0.85 (0.2)	0.71 (0.074)	0.46 (2.6)
Gamma	Softplus	0.023 (0.56)	1	0.72 (0.17)	0.65 (0.078)	0 (0)
Gompertz	Log	99 (487)	0.9	0.68 (0.16)	0.63 (0.078)	0.011 (0.47)
Gompertz	Softplus	7.3 (80)	0.95	0.63 (0.17)	0.61 (0.097)	0.002 (0.001)
Normal	Log	2.4 (20)	0.96	0.7 (0.18)	0.63 (0.098)	0.45 (3.3)
Normal	Softplus	0.68 (10)	0.97	0.81 (0.21)	0.67 (0.1)	0.001 (0.072)
trans. Normal	Log	0 (0)	1	0.85 (0.21)	0.71 (0.074)	0.3 (3.9)
trans. Normal	Softplus	0 (0)	1	0.86 (0.21)	0.71 (0.073)	0.4 (2.8)
Weibull	Log	0.46 (22)	1	0.8 (0.19)	0.69 (0.075)	0.001 (0)
Weibull	Softplus	0.34 (15)	1	0.67 (0.16)	0.62 (0.083)	0.016 (0.59)

Link	Likelihood	PKs	ELPD _{loo}	bias(β_{xy})	SD(β_{xy})	RMSE(β_{xy})
Beta prime	Log	0.31 (0.57)	-190 (100)	0.12 (0.13)	0.07 (0.04)	0.15 (0.12)
Beta prime	Softplus	0.33 (0.62)	-208 (84)	0.2 (0.29)	0.09 (0.06)	0.24 (0.28)
Fréchet	Log	1.1 (1)	-186 (97)	0.13 (0.14)	0.08 (0.06)	0.17 (0.14)
Fréchet	Softplus	0.95 (0.99)	-202 (107)	0.15 (0.15)	0.09 (0.06)	0.19 (0.14)
Gamma	Log	0.19 (0.44)	-201 (66)	0.15 (0.18)	0.08 (0.08)	0.18 (0.19)
Gamma	Softplus	0.52 (0.77)	-186 (95)	0.13 (0.14)	0.09 (0.07)	0.18 (0.15)
Gompertz	Log	0.75 (0.89)	-188 (99)	0.13 (0.14)	0.08 (0.06)	0.17 (0.14)
Gompertz	Softplus	0.97 (1.3)	-199 (106)	0.6 (0.68)	0.34 (0.26)	0.75 (0.67)
Normal	Log	0.71 (0.92)	-222 (87)	0.93 (1.2)	0.43 (0.4)	1.1 (1.2)
Normal	Softplus	0.36 (0.63)	-184 (109)	0.5 (0.47)	0.32 (0.18)	0.65 (0.43)
trans. Normal	Log	0.21 (0.46)	-194 (126)	0.57 (0.51)	0.32 (0.18)	0.71 (0.46)
trans. Normal	Softplus	0.25 (0.47)	-200 (86)	0.56 (0.52)	0.34 (0.18)	0.71 (0.47)
Weibull	Log	0.4 (0.95)	-187 (113)	0.53 (0.48)	0.33 (0.17)	0.68 (0.43)
Weibull	Softplus	0.87 (1.2)	-187 (112)	0.5 (0.46)	0.3 (0.15)	0.64 (0.41)

Note: Most results are presented as mean(sd). $P(\text{Conv.})$ is calculated as a proportion of all models. TESS_{total}^{bulk} is also calculated using all models. Both the metrics for predictive performance and parameter recoverability are calculated only for models that had 10 or less divergent transitions and converged. In addition the parameter recoverability metrics are only shown for models using the true link.

References

- [1] Stan modeling language users guide and reference manual, 2.30.0, 2022. URL <https://mc-stan.org>.
- [2] C. J. Adcock. Asset pricing and portfolio selection based on the multivariate extended skew-Student-t distribution. *Annals of Operations Research*, 176(1):221–234, April 2010. ISSN 0254-5330, 1572-9338. doi: 10.1007/s10479-009-0586-4.
- [3] J. Aitchison. The Statistical Analysis of Compositional Data. *Journal of the Royal Statistical Society: Series B (Methodological)*, 44(2):139–160, 1982. ISSN 2517-6161. doi: 10.1111/j.2517-6161.1982.tb01195.x.
- [4] Hafzullah Aksoy. Use of gamma distribution in hydrological analysis. *Turkish Journal of Engineering and Environmental Sciences*, 24(6):419–428, 2000.

- [5] RB Arellano-Valle, Heleno Bolfarine, and VH Lachos. Skew-normal linear mixed models. *Journal of Data Science*, 3(4):415–438, 2005.
- [6] Noah Arthurs, Ben Stenhaug, Sergey Karayev, and Chris Piech. *Grades Are Not Normal: Improving Exam Score Models Using the Logit-Normal Distribution*. International Educational Data Mining Society, July 2019.
- [7] J. Atchison and S.M. Shen. Logistic-normal distributions: Some properties and uses. *Biometrika*, 67(2):261–272, January 1980. ISSN 0006-3444. doi: 10.1093/biomet/67.2.261.
- [8] O. E. Barndorff-Nielsen and B. Jørgensen. Some parametric models on the simplex. *Journal of Multivariate Analysis*, 39(1):106–116, October 1991. ISSN 0047-259X. doi: 10.1016/0047-259X(91)90008-P.
- [9] C. P. Barros, Peter Wanke, Silvestre Dumbo, and Jose Pires Manso. Efficiency in angolan hydro-electric power station: A two-stage virtual frontier dynamic DEA and simplex regression approach. *Renewable and Sustainable Energy Reviews*, 78:588–596, October 2017. ISSN 1364-0321. doi: 10.1016/j.rser.2017.04.100.
- [10] Douglas Bates, Martin Mächler, Ben Bolker, and Steve Walker. Fitting Linear Mixed-Effects Models Using `\pkg{lme4}`. *Journal of Statistical Software*, 67(1):1–48, 2015.
- [11] Stephen Bates, Trevor Hastie, and Robert Tibshirani. Cross-validation: what does it estimate and how well does it do it?, July 2022.
- [12] Aleksey V Belikov. The number of key carcinogenic events can be predicted from cancer incidence. *Scientific reports*, 7(1):1–8, 2017.
- [13] Betancourt. Towards A Principled Bayesian Workflow, April 2020. URL https://betanalpha.github.io/assets/case_studies/principled_bayesian_workflow.html.
- [14] Michael Betancourt. Diagnosing suboptimal cotangent disintegrations in hamiltonian monte carlo. *arXiv preprint arXiv:1604.00695*, 2016.
- [15] Wagner Hugo Bonat, Paulo Justiniano RIBEIRO Jr, and Walmes Marques Zeviani. Regression models with responses on the unity Interval: Specification, estimation and comparison. *Biometric Brazilian Journal*, 30(4):18, 2013.
- [16] Marcelo Bourguignon, Manoel Santos-Neto, and Mário de Castro. A new regression model for positive data. *arXiv:1804.07734 [stat]*, April 2018. arXiv: 1804.07734.
- [17] Leo Breiman. Statistical Modeling: The Two Cultures (with comments and a rejoinder by the author). *Statistical Science*, 16(3):199–231, August 2001. ISSN 0883-4237, 2168-8745. doi: 10.1214/ss/1009213726. Publisher: Institute of Mathematical Statistics.
- [18] TA Burgin. The gamma distribution and inventory control. *Journal of the Operational Research Society*, 26(3):507–525, 1975.
- [19] Paul-Christian Bürkner. brms: An r package for bayesian multilevel models using stan. *Journal of statistical software*, 80:1–28, 2017.
- [20] David E Burmaster and Delores A Hull. Using lognormal distributions and lognormal probability plots in probabilistic risk assessments. *Human and Ecological Risk Assessment*, 3(2):235–255, 1997.
- [21] Paul-Christian Bürkner. Bayesian Item Response Modelling in R with brms and Stan. *Journal of Statistical Software*, pages 1–54, 2021.
- [22] Paul-Christian Bürkner, Jonah Gabry, Matthew Kay, and Aki Vehtari. posterior: Tools for working with posterior distributions, 2022. URL <https://mc-stan.org/posterior/>. R package version 1.3.0.

- [23] Paul-Christian Bürkner, Maximilian Scholz, and Stefan T. Radev. What makes a good bayesian model? towards a unified model taxonomy. *arXiv preprint*, 2022.
- [24] Diego Ramos Canterle and Fábio Mariano Bayer. Variable dispersion beta regressions with parametric link functions. *Statistical Papers*, 60(5):1541–1567, October 2019. ISSN 1613-9798. doi: 10.1007/s00362-017-0885-9.
- [25] Bob Carpenter, Andrew Gelman, Matthew D Hoffman, Daniel Lee, Ben Goodrich, Michael Betancourt, Marcus Brubaker, Jiqiang Guo, Peter Li, and Allen Riddell. Stan: A probabilistic programming language. *Journal of statistical software*, 76(1), 2017.
- [26] Daniela Castro-Camilo, Raphaël Huser, and Håvard Rue. Practical strategies for GEV-based regression models for extremes. May 2022. arXiv:2106.13110 [stat].
- [27] Raj Chhikara. *The inverse Gaussian distribution: theory: methodology, and applications*, volume 95. CRC Press, 1988.
- [28] RS Chhikara and JL Folks. The inverse gaussian distribution as a lifetime model. *Technometrics*, 19(4):461–468, 1977.
- [29] Carlos Cinelli, Andrew Forney, and Judea Pearl. A Crash Course in Good and Bad Controls. *SSRN Electronic Journal*, 2020. ISSN 1556-5068. doi: 10.2139/ssrn.3689437.
- [30] Francisco Cribari-Neto and Achim Zeileis. Beta Regression in *R*. *Journal of Statistical Software*, 34(2), 2010. ISSN 1548-7660. doi: 10.18637/jss.v034.i02.
- [31] S A Damisa, F Musa, and S Sani. On the Comparison of Some Link Functions In Binary Response Analysis. 1:6, 2017.
- [32] Felipe RS De Gusmao, Edwin MM Ortega, and Gauss M Cordeiro. The generalized inverse weibull distribution. *Statistical Papers*, 52(3):591–619, 2011.
- [33] Brian Dennis and GP Patil. Applications in ecology. In *Lognormal distributions*, pages 303–330. Routledge, 2018.
- [34] Charles Dugas, Yoshua Bengio, François Bédille, Claude Nadeau, and René Garcia. Incorporating second-order functional knowledge for better option pricing. *Advances in neural information processing systems*, 13, 2000.
- [35] P. Espinheira, S. Ferrari, and Francisco Cribari-Neto. On beta regression residuals. 2008. doi: 10.1080/02664760701834931.
- [36] Ludwig Fahrmeir, Gerhard Tutz, Wolfgang Hennevogl, and Eliane Salem. *Multivariate statistical modelling based on generalized linear models*, volume 425. Springer, 1994.
- [37] Ludwig Fahrmeir, Thomas Kneib, Stefan Lang, and Brian D Marx. Regression models. In *Regression*, pages 23–84. Springer, 2021.
- [38] Silvia Ferrari and Francisco Cribari-Neto. Beta Regression for Modelling Rates and Proportions. *Journal of Applied Statistics*, 31(7):799–815, August 2004. ISSN 0266-4763, 1360-0532. doi: 10.1080/0266476042000214501.
- [39] S. G. Fletcher and K. Ponnambalam. Estimation of reservoir yield and storage distribution using moments analysis. *Journal of Hydrology*, 182(1):259–275, July 1996. ISSN 0022-1694. doi: 10.1016/0022-1694(95)02946-X.
- [40] J Leroy Folks and Raj S Chhikara. The inverse gaussian distribution and its statistical application—a review. *Journal of the Royal Statistical Society: Series B (Methodological)*, 40(3):263–275, 1978.

- [41] Patrizio Frederic and Frank Lad. Two Moments of the Logitnormal Distribution. *Communications in Statistics - Simulation and Computation*, 37(7):1263–1269, August 2008. ISSN 0361-0918. doi: 10.1080/03610910801983178.
- [42] Jonah Gabry, Daniel Simpson, Aki Vehtari, Michael Betancourt, and Andrew Gelman. Visualization in Bayesian workflow. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, 182(2):389–402, February 2019. ISSN 09641998. doi: 10.1111/rssa.12378.
- [43] Andrew Gelman, John B. Carlin, Hal S. Stern, David B. Dunson, Aki Vehtari, and Donald B. Rubin. *Bayesian Data Analysis*. Chapman and Hall/CRC, New York, 3 edition, July 2013. ISBN 978-0-429-11307-9. doi: 10.1201/b16018.
- [44] Andrew Gelman, Daniel Simpson, and Michael Betancourt. The prior can often only be understood in the context of the likelihood. *Entropy*, 19(10):555–567, 2017. Publisher: Multidisciplinary Digital Publishing Institute.
- [45] Andrew Gelman, Jennifer Hill, and Aki Vehtari. *Regression and other stories*. Cambridge University Press, 2020.
- [46] Andrew Gelman, Aki Vehtari, Daniel Simpson, Charles C. Margossian, Bob Carpenter, Yuling Yao, Lauren Kennedy, Jonah Gabry, Paul-Christian Bürkner, and Martin Modrák. Bayesian Workflow. *arXiv:2011.01808 [stat]*, November 2020. arXiv: 2011.01808.
- [47] Jeff Gill, Jefferson M Gill, Michelle Torres, and Silvia Michelle Torres Pacheco. *Generalized linear models: a unified approach*, volume 134. Sage Publications, 2001.
- [48] Benjamin Gompertz. On the nature of the function expressive of the law of human mortality, and on a new mode of determining the value of life contingencies. *Philosophical Transactions of the Royal Society of London*, 115:513–583, January 1825. doi: 10.1098/rstl.1825.0026.
- [49] Emil Julius Gumbel. *Statistical theory of extreme values and some practical applications: a series of lectures*, volume 33. US Government Printing Office, 1954.
- [50] Arjun K. Gupta and Saralees Nadarajah, editors. *Handbook of Beta Distribution and Its Applications*. CRC Press, Boca Raton, April 2014. ISBN 978-0-429-15288-7. doi: 10.1201/9781482276596.
- [51] Michael R Hagerty and V Srinivasan. Comparing the predictive powers of alternative multiple regression models. *Psychometrika*, 56(1):77–85, 1991.
- [52] Soudabeh Hamed-Shahraki, Aliakbar Rasekhi, Mir Saeed Yekaninejad, Mohammad Reza Eshraghian, and Amir H. Pakpour. Kumaraswamy regression modeling for Bounded Outcome Scores. *Pakistan Journal of Statistics and Operation Research*, pages 79–88, March 2021. ISSN 2220-5810. doi: 10.18187/pjsor.v17i1.3411.
- [53] Frank E Harrell. *Regression modeling strategies: with applications to linear models, logistic and ordinal regression, and survival analysis*. Springer, 2015.
- [54] Ottar Hellevik. Linear versus logistic regression when the dependent variable is a dichotomy. *Quality & Quantity*, 43(1):59–74, January 2009. ISSN 0033-5177, 1573-7845. doi: 10.1007/s11135-007-9077-3.
- [55] Jong-Hyeon Jeong and Jason Fine. Direct parametric inference for the cumulative incidence function. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 55(2):187–200, April 2006. ISSN 0035-9254, 1467-9876. doi: 10.1111/j.1467-9876.2006.00532.x.
- [56] Jong-Hyeon Jeong and Jason P. Fine. Parametric regression on cumulative incidence function. *Biostatistics*, 8(2):184–196, April 2007. ISSN 1465-4644. doi: 10.1093/biostatistics/kxj040.

- [57] Xun Jiang, Dipak K. Dey, Rachel Prunier, Adam M. Wilson, and Kent E. Holsinger. A new class of flexible link functions with application to species co-occurrence in cape floristic region. *The Annals of Applied Statistics*, 7(4):2180–2204, 2013. ISSN 1932-6157.
- [58] Norman L. Johnson, Samuel Kotz, and Narayanaswamy Balakrishnan. *Continuous Univariate Distributions, Volume 2*. John Wiley & Sons, May 1995. ISBN 978-0-471-58494-0.
- [59] Norman L Johnson, Samuel Kotz, and Adrienne W Kemp. *Univariate discrete distributions*. John Wiley & Sons, 2005.
- [60] M. C. Jones. Kumaraswamy’s distribution: A beta-type distribution with some tractability advantages. *Statistical Methodology*, 6(1):70–81, January 2009. ISSN 1572-3127. doi: 10.1016/j.stamet.2008.04.001.
- [61] Bent Jorgensen. *The Theory of Dispersion Models*. CRC Press, June 1997. ISBN 978-0-412-99711-2.
- [62] Miguel A. Juárez and Mark F. J. Steel. Model-Based Clustering of Non-Gaussian Panel Data Based on Skew- t Distributions. *Journal of Business & Economic Statistics*, 28(1):52–66, January 2010. ISSN 0735-0015, 1537-2707. doi: 10.1198/jbes.2009.07145.
- [63] Ernest Sydney Keeping. *Introduction to Statistical Inference*. Courier Corporation, 1995.
- [64] M Shuaib Khan, GR Pasha, and Ahmed Hesham Pasha. Theoretical analysis of inverse weibull distribution. *WSEAS Transactions on Mathematics*, 7(2):30–38, 2008.
- [65] Morteza Khodabina and Alireza Ahmadabadib. Some properties of generalized gamma distribution. 2010.
- [66] Robert Kieschnick and B D McCullough. Regression analysis of variates observed on $(0, 1)$: percentages, proportions and fractions. *Statistical Modelling*, 3(3):193–213, October 2003. ISSN 1471-082X. doi: 10.1191/1471082X03st053oa.
- [67] Hee-Cheul Kim. A Performance Analysis of Software Reliability Model using Lomax and Gompertz Distribution Property. *Indian Journal of Science and Technology*, 9(20), May 2016. ISSN 0974-5645, 0974-6846. doi: 10.17485/ijst/2016/v9i20/94668.
- [68] Seongho Kim, Elisabeth Heath, and Lance Heilbrun. Sample size determination for logistic regression on a logit-normal distribution. *Statistical Methods in Medical Research*, 26(3):1237–1247, June 2017. ISSN 0962-2802. doi: 10.1177/0962280215572407. Publisher: SAGE Publications Ltd STM.
- [69] Roger Koenker and Jungmo Yoon. Parametric links for binary choice models: A Fisherian–Bayesian colloquy. *Journal of Econometrics*, 152(2):120–130, October 2009. ISSN 0304-4076. doi: 10.1016/j.jeconom.2009.01.009.
- [70] Ken’ichirou Kosugi. Lognormal distribution model for unsaturated soil hydraulic properties. *Water Resources Research*, 32(9):2697–2703, 1996.
- [71] Włodzimierz Kryszicki. On some new properties of the beta distribution. *Statistics & Probability Letters*, 42(2):131–137, April 1999. ISSN 0167-7152. doi: 10.1016/S0167-7152(98)00197-7.
- [72] P. Kumaraswamy. A generalized probability density function for double-bounded random processes. *Journal of Hydrology*, 46(1):79–88, March 1980. ISSN 0022-1694. doi: 10.1016/0022-1694(80)90036-0.
- [73] Chin-Diew Lai, DN Murthy, and Min Xie. Weibull distributions and their applications. In *Springer Handbooks*, pages 63–78. Springer, 2006.
- [74] Ben Lambert and Aki Vehtari. \mathcal{R}^* : A robust MCMC convergence diagnostic with uncertainty using decision tree classifiers. *arXiv:2003.07900 [stat]*, November 2020. arXiv: 2003.07900.

- [75] Artur J. Lemonte and Jorge L. Bazán. New links for binary regression: an application to coca cultivation in Peru. *TEST*, 27(3):597–617, September 2018. ISSN 1133-0686, 1863-8260. doi: 10.1007/s11749-017-0563-1.
- [76] Torrin M Liddell and John K Kruschke. Analyzing ordinal data with metric models: What could possibly go wrong? *Journal of Experimental Social Psychology*, 79:328–348, 2018. Publisher: Elsevier.
- [77] Freddy Omar López. A Bayesian Approach to Parameter Estimation in Simplex Regression Model: A Comparison with Beta Regression. *Revista Colombiana de Estadística*, page 21, 2013.
- [78] Stephen R Martin and Donald R Williams. Outgrowing the procrustean bed of normality: the utility of Bayesian modeling for asymmetrical data analysis. *PsyArXiv preprint*, 2017.
- [79] P. McCullagh. *Generalized Linear Models*. Routledge, New York, 2 edition, January 2019. ISBN 978-0-203-75373-6. doi: 10.1201/9780203753736.
- [80] James B. McDonald. Model selection: some generalized distributions. *Communications in Statistics - Theory and Methods*, June 2007. ISSN 1049-1074. doi: 10.1080/03610928708829422.
- [81] James B McDonald and Richard J Butler. Regression models for positive random variables. *Journal of Econometrics*, 43(1):227–251, January 1990. ISSN 0304-4076. doi: 10.1016/0304-4076(90)90118-D.
- [82] Richard McElreath. *Statistical Rethinking: A Bayesian Course with Examples in R and Stan*. Chapman and Hall/CRC, New York, 2 edition, March 2020. ISBN 978-0-429-02960-8. doi: 10.1201/9780429029608.
- [83] R. Mead. A Generalised Logit-Normal Distribution. *Biometrics*, 21(3):721–732, 1965. ISSN 0006-341X. doi: 10.2307/2528553.
- [84] Francisco M. C. Medeiros, Mariana C. Araújo, and Marcelo Bourguignon. Improved estimators in beta prime regression models. *Communications in Statistics - Simulation and Computation*, pages 1–14, October 2021. ISSN 0361-0918, 1532-4141. doi: 10.1080/03610918.2021.1990322.
- [85] Petrus Mikkola, Osvaldo A Martin, Suyog Chandramouli, Marcelo Hartmann, Oriol Abril Pla, Owen Thomas, Henri Pesonen, Jukka Corander, Aki Vehtari, Samuel Kaski, Bürkner, Paul-Christian, and Klami, Arto. Prior knowledge elicitation: The past, present, and future. *arXiv preprint arXiv:2112.01380*, 2021.
- [86] Pablo A. Mitnik. New Properties of the Kumaraswamy Distribution. *Communications in Statistics - Theory and Methods*, 42(5):741–755, March 2013. ISSN 0361-0926. doi: 10.1080/03610926.2011.581782.
- [87] Pablo A. Mitnik and Sunyoung Baek. The Kumaraswamy distribution: median-dispersion reparameterizations for regression modeling and simulation-based estimation. *Statistical Papers*, 54(1):177–192, February 2013. ISSN 1613-9798. doi: 10.1007/s00362-011-0417-y.
- [88] Byron JT Morgan and DM Smith. A note on Wadley’s problem with overdispersion. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 41(2):349–354, 1992. Publisher: Wiley Online Library.
- [89] DN Prabhakar Murthy, Min Xie, and Renyan Jiang. *Weibull models*. John Wiley & Sons, 2004.
- [90] Saralees Nadarajah. On the distribution of Kumaraswamy. *Journal of Hydrology*, 348(3):568–569, January 2008. ISSN 0022-1694. doi: 10.1016/j.jhydrol.2007.09.008.
- [91] Danielle J. Navarro. Between the Devil and the Deep Blue Sea: Tensions Between Scientific Judgement and Statistical Model Selection. *Computational Brain & Behavior*, 2(1):28–34, March 2019. ISSN 2522-087X. doi: 10.1007/s42113-018-0019-z.

- [92] John Ashworth Nelder and Robert WM Wedderburn. Generalized linear models. *Journal of the Royal Statistical Society: Series A (General)*, 135(3):370–384, 1972. Publisher: Wiley Online Library.
- [93] Judea Pearl. *Causality*. Cambridge university press, 2009.
- [94] Judea Pearl. Causal inference in statistics: An overview. *Statistics Surveys*, 3(none):96–146, January 2009. ISSN 1935-7516. doi: 10.1214/09-SS057. Publisher: Amer. Statist. Assoc., the Bernoulli Soc., the Inst. Math. Statist., and the Statist. Soc. Canada.
- [95] Judea Pearl. The do-calculus revisited. *arXiv preprint arXiv:1210.4852*, 2012.
- [96] P. Pinson. Very-short-term probabilistic forecasting of wind power with generalized logit–normal distributions. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 61(4):555–576, 2012. ISSN 1467-9876. doi: 10.1111/j.1467-9876.2011.01026.x.
- [97] Shima Pirmohammadi and Hamid Bidram. On the Liu estimator in the beta and Kumaraswamy regression models: A comparative study. *Communications in Statistics - Theory and Methods*, 0(0):1–26, April 2021. ISSN 0361-0926. doi: 10.1080/03610926.2021.1900254.
- [98] Daniel Powers and Yu Xie. *Statistical methods for categorical data analysis*. Emerald Group Publishing, 2008.
- [99] R Core Team. R: A language and environment for statistical computing. 2022. URL <https://www.R-project.org/>.
- [100] Pedro L Ramos, Francisco Louzada, Eduardo Ramos, and Sanku Dey. The fréchet distribution: Estimation and application-an overview. *Journal of Statistics and Management Systems*, 23(3): 549–578, 2020.
- [101] Horst Rinne. *The Weibull distribution: a handbook*. Chapman and Hall/CRC, 2008.
- [102] A. Scallan. Some Aspects of Parametric Link Functions. In Robert Gilchrist, editor, *GLIM 82: Proceedings of the International Conference on Generalised Linear Models*, Lecture Notes in Statistics, pages 122–127, New York, NY, 1982. Springer. ISBN 978-1-4612-5771-4. doi: 10.1007/978-1-4612-5771-4_12.
- [103] A. Scallan, R. Gilchrist, and M. Green. Fitting parametric link functions in generalised linear models. *Computational Statistics & Data Analysis*, 2(1):37–49, June 1984. ISSN 0167-9473. doi: 10.1016/0167-9473(84)90031-8.
- [104] Maximilian Scholz and Paul-Christian Bürkner. Prediction can be safely used as a proxy for explanation in causally consistent bayesian generalized linear models, September 2022. URL <https://osf.io/xgkzv/>.
- [105] Maximilian Scholz, Yannick Dzubba, and Paul-Christian Bürkner. Bayesim, Sep 2022. URL <https://github.com/sims1253/bayesim>. R package version 0.29.5.9000.
- [106] Venkata Seshadri. *The inverse Gaussian distribution: statistical theory and applications*, volume 137. Springer Science & Business Media, 2012.
- [107] Galit Shmueli. To Explain or to Predict? *Statistical Science*, 25(3), August 2010. ISSN 0883-4237. doi: 10.1214/10-STS330.
- [108] Galit Shmueli, Thomas P Minka, Joseph B Kadane, Sharad Borle, and Peter Boatwright. A useful distribution for fitting discrete data: revival of the Conway–Maxwell–Poisson distribution. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 54(1):127–142, 2005. doi: 10.1111/j.1467-9876.2005.00474.x. Publisher: Wiley Online Library.

- [109] Peter X.-K. Song, Zhenguo Qiu, and Ming Tan. Modelling Heterogeneous Dispersion in Marginal Models for Longitudinal Proportional Data. *Biometrical Journal*, 46(5):540–553, 2004. ISSN 1521-4036. doi: 10.1002/bimj.200110052.
- [110] Peter Xue-Kun Song and Ming Tan. Marginal Models for Longitudinal Continuous Proportional Data. *Biometrics*, 56(2):496–502, 2000. ISSN 1541-0420. doi: 10.1111/j.0006-341X.2000.00496.x.
- [111] David J Spiegelhalter, Nicola G Best, Bradley P Carlin, and A Van der Linde. Bayesian deviance, the effective number of parameters, and the comparison of arbitrarily complex models. Technical report, Citeseer, 1998.
- [112] D Mikis Stasinopoulos and Robert A Rigby. Generalized additive models for location scale and shape (GAMLSS) in R. *Journal of Statistical Software*, 23(7):1–46, 2007.
- [113] Walter W Stroup. Rethinking the analysis of non-normal data in plant and soil science. *Agronomy journal*, 107(2):811–827, 2015.
- [114] E. Sunandi, A. Kurnia, K. Sadik, and K. A. Notodiputro. A Bayesian Logit-Normal Model in Small Area Estimation. *Journal of Physics: Conference Series*, 1863(1):012039, March 2021. ISSN 1742-6596. doi: 10.1088/1742-6596/1863/1/012039. Publisher: IOP Publishing.
- [115] Sean Talts, Michael Betancourt, Daniel Simpson, Aki Vehtari, and Andrew Gelman. Validating bayesian inference algorithms with simulation-based calibration. doi: 10.48550/ARXIV.1804.06788. URL <https://arxiv.org/abs/1804.06788>.
- [116] Herbert CS Thom. A note on the gamma distribution. *Monthly weather review*, 86(4):117–122, 1958.
- [117] RJ Torrent. The log-normal distribution: A better fitness for the results of mechanical testing of materials. *Matériaux et Construction*, 11(4):235–245, 1978.
- [118] Alexander Tulupyeu, Alena Suvorova, Jennifer Sousa, and Daniel Zelterman. Beta prime regression with application to risky behavior frequency screening. *Statistics in Medicine*, 32(23):4044–4056, 2013. ISSN 1097-0258. doi: 10.1002/sim.5820.
- [119] J. C. S. Vasconcelos, F. Pratavia, E. M. M. Ortega, and G. M. Cordeiro. An extended logit-normal regression with application to human development index data. *Communications in Statistics - Simulation and Computation*, 0(0):1–12, March 2022. ISSN 0361-0918. doi: 10.1080/03610918.2022.2045497.
- [120] Aki Vehtari and Janne Ojanen. A survey of Bayesian predictive methods for model assessment, selection and comparison. *Statistics Surveys*, 6(none):142–228, January 2012. ISSN 1935-7516. doi: 10.1214/12-SS102.
- [121] Aki Vehtari, Daniel Simpson, Andrew Gelman, Yuling Yao, and Jonah Gabry. Pareto smoothed importance sampling. *arXiv preprint arXiv:1507.02646*, 2015.
- [122] Aki Vehtari, Andrew Gelman, and Jonah Gabry. Practical bayesian model evaluation using leave-one-out cross-validation and waic. *Statistics and computing*, 27(5):1413–1432, 2017.
- [123] Aki Vehtari, Andrew Gelman, Daniel Simpson, Bob Carpenter, and Paul-Christian Bürkner. Rank-normalization, folding, and localization: An improved \widehat{R} for assessing convergence of MCMC. *Bayesian Analysis*, 16(2), June 2021. ISSN 1936-0975. doi: 10.1214/20-BA1221. arXiv: 1903.08008.
- [124] Aki Vehtari, Jonah Gabry, Mans Magnusson, Yuling Yao, Paul-Christian Bürkner, Topi Paananen, and Andrew Gelman. loo: Efficient leave-one-out cross-validation and waic for bayesian models, 2022. URL <https://mc-stan.org/loo/>. R package version 2.5.0.

- [125] Sumio Watanabe. *Algebraic geometry and statistical learning theory*. Number 25. Cambridge university press, 2009.
- [126] Sumio Watanabe and Manfred Opper. Asymptotic equivalence of Bayes cross validation and widely applicable information criterion in singular learning theory. *Journal of machine learning research*, 11(12), 2010.
- [127] Sebastian Weber and Paul-Christian Bürkner. Running brms models with within-chain parallelization, Sep 2022. URL https://cran.r-project.org/web/packages/brms/vignettes/brms_threading.html.
- [128] Paul FV Wiemann, Thomas Kneib, and Julien Homburgers. Using the softplus function to construct alternative link functions in generalized linear models and beyond. *arXiv preprint arXiv:2111.14207*, 2021.
- [129] Daniel S Wilks. Maximum likelihood estimation for the gamma distribution using data containing zeros. *Journal of climate*, pages 1495–1501, 1990.
- [130] Donald R Williams and Stephen R Martin. Rethinking robust statistics with modern Bayesian methods. *PsyArXiv preprint*, 2017.
- [131] Bodo Winter and Paul-Christian Bürkner. Poisson regression for linguists: A tutorial introduction to modelling count data with brms. *Language and Linguistics Compass*, 15(11):e12439, 2021. Publisher: Wiley Online Library.
- [132] Yuling Yao, Aki Vehtari, Daniel Simpson, and Andrew Gelman. Using stacking to average Bayesian predictive distributions (with discussion). *Bayesian Analysis*, 13(3):917–1007, 2018. Publisher: International Society for Bayesian Analysis.
- [133] Tal Yarkoni and Jacob Westfall. Choosing Prediction Over Explanation in Psychology: Lessons From Machine Learning. *Perspectives on Psychological Science*, 12(6):1100–1122, November 2017. ISSN 1745-6916, 1745-6924. doi: 10.1177/1745691617693393.
- [134] Thomas W Yee. The VGAM package for categorical data analysis. *Journal of Statistical Software*, 32(10):1–34, 2010. Publisher: Citeseer.
- [135] Shuang Yin, Dipak K. Dey, Emiliano A. Valdez, Guojun Gan, and Jeyaraj Vadiveloo. Skewed link regression models for imbalanced binary response with applications to life insurance. *arXiv:2007.15172 [stat]*, July 2020.
- [136] Haitham M Yousof, Ahmed Z Afify, N Ebrahim Abd El Hadi, Gholamhossein G Hamedani, and Nadeem Shafique Butt. On six-parameter fréchet distribution: properties and applications. *Pakistan Journal of Statistics and Operation Research*, pages 281–299, 2016.
- [137] Zhongheng Zhang. Parametric regression model for survival data: Weibull regression model as an example. *Annals of translational medicine*, 4(24), 2016.
- [138] Hao Zheng, Zhanlei Yang, Wenju Liu, Jizhong Liang, and Yanpeng Li. Improving deep neural networks using softplus units. In *2015 International joint conference on neural networks (IJCNN)*, pages 1–4. IEEE, 2015.
- [139] Mark H Zweig and Gregory Campbell. Receiver-operating characteristic (roc) plots: a fundamental evaluation tool in clinical medicine. *Clinical chemistry*, 39(4):561–577, 1993.