

# Estimating Distributional Models with brms

Paul Bürkner

2019-05-22

- [Introduction](#)
- [A simple distributional model](#)
- [Zero-Inflated Models](#)
- [Additive Distributional Models](#)

## Introduction

This vignette provides an introduction on how to fit distributional regression models with **brms**. We use the term *distributional model* to refer to a model, in which we can specify predictor terms for all parameters of the assumed response distribution. In the vast majority of regression model implementations, only the location parameter (usually the mean) of the response distribution depends on the predictors and corresponding regression parameters. Other parameters (e.g., scale or shape parameters) are estimated as auxiliary parameters assuming them to be constant across observations. This assumption is so common that most researchers applying regression models are often (in my experience) not aware of the possibility of relaxing it. This is understandable insofar as relaxing this assumption drastically increase model complexity and thus makes models hard to fit. Fortunately, **brms** uses **Stan** on the backend, which is an incredibly flexible and powerful tool for estimating Bayesian models so that model complexity is much less of an issue.

Suppose we have a normally distributed response variable. Then, in basic linear regression, we specify a predictor term  $\eta_\mu$  for the mean parameter  $\mu$  of the normal distribution. The second parameter of the normal distribution – the residual standard deviation  $\sigma$  – is assumed to be constant across observations. We estimate  $\sigma$  but do not try to *predict* it. In a distributional model, however, we do exactly this by specifying a predictor term  $\eta_\sigma$  for  $\sigma$  in addition to the predictor term  $\eta_\mu$ . Ignoring group-level effects for the moment, the linear predictor of a parameter  $\theta$  for observation  $n$  has the form

$$\eta_{\theta n} = \sum_{i=1}^{K_\theta} b_{\theta i} x_{\theta i n}$$

where  $x_{\theta i n}$  denotes the value of the  $i$ th predictor of parameter  $\theta$  for observation  $n$  and  $b_{\theta i}$  is the  $i$ th regression coefficient of parameter  $\theta$ . A distributional normal model with response variable  $y$  can then be written as

$$y_n \sim \mathcal{N}(\eta_{\mu n}, \exp(\eta_{\sigma n}))$$

We used the exponential function around  $\eta_\sigma$  to reflect that  $\sigma$  constitutes a standard deviation and thus only takes on positive values, while a linear predictor can be any real number.

## A simple distributional model

Unequal variance models are possibly the most simple, but nevertheless very important application of distributional models. Suppose we have two groups of patients: One group receives a treatment (e.g., an antidepressive drug) and another group receives placebo. Since the treatment may not work equally well for all patients, the symptom variance of the treatment group may be larger than the symptom variance of the placebo group after some weeks of treatment. For simplicity, assume that we only investigate the post-treatment values.

```
group <- rep(c("treat", "placebo"), each = 30)
symptom_post <- c(rnorm(30, mean = 1, sd = 2), rnorm(30, mean = 0, sd = 1))
dat1 <- data.frame(group, symptom_post)
head(dat1)
```

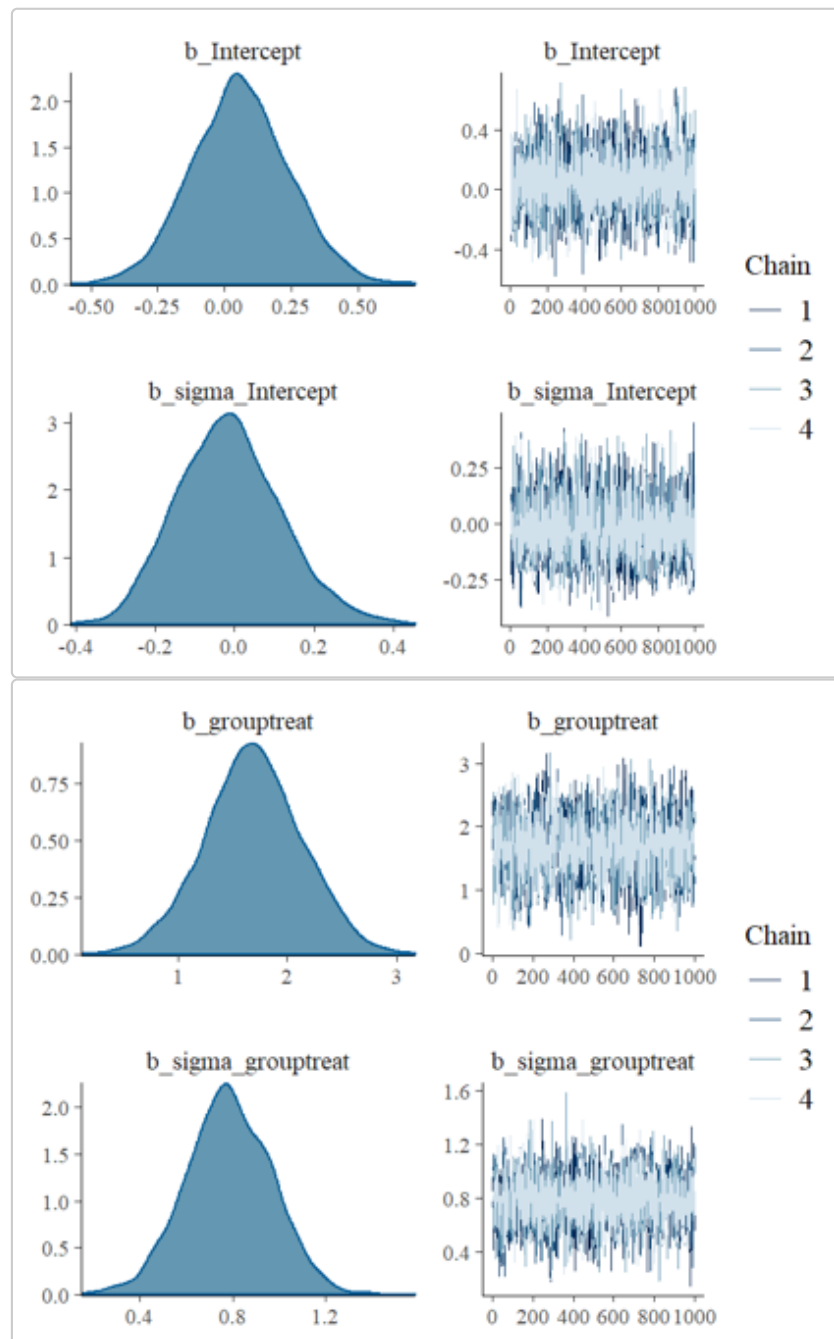
```
  group symptom_post
1 treat      2.5520771
2 treat      3.7890842
3 treat      1.2501676
4 treat     -0.8195757
5 treat      2.1445745
6 treat      1.3127191
```

The following model estimates the effect of group on both the mean and the residual standard deviation of the normal response distribution.

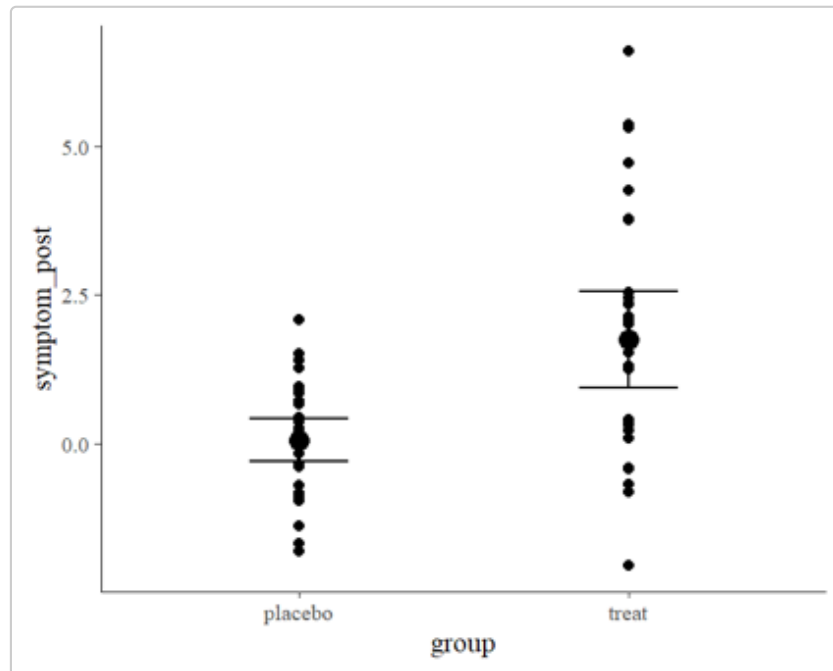
```
fit1 <- brm(bf(symptom_post ~ group, sigma ~ group),
            data = dat1, family = gaussian())
```

Useful summary statistics and plots can be obtained via

```
summary(fit1)
plot(fit1, N = 2, ask = FALSE)
```



```
plot(marginal_effects(fit1), points = TRUE)
```



The population-level effect `sigma_grouptreat`, which is the contrast of the two residual standard deviations on the log-scale, reveals that the variances of both groups are indeed different. This impression is confirmed when looking at the `marginal_effects` of `group`. Going one step further, we can compute the residual standard deviations on the original scale using the `hypothesis` method.

```
hyp <- c("exp(sigma_Intercept) = 0",
        "exp(sigma_Intercept + sigma_grouptreat) = 0")
hypothesis(fit1, hyp)
```

Hypothesis Tests for class b:

	Hypothesis	Estimate	Est.Error	CI.Lower	CI.Upper	Evid.Ratio	Post.Prob
Star							
1	(exp(sigma_Interc... = 0	1.00	0.13	0.78	1.30	NA	NA
*							
2	(exp(sigma_Interc... = 0	2.18	0.29	1.70	2.85	NA	NA
*							

---

'CI': 90%-CI for one-sided and 95%-CI for two-sided hypotheses.

'\*': For one-sided hypotheses, the posterior probability exceeds 95%;

for two-sided hypotheses, the value tested against lies outside the 95%-CI.

Posterior probabilities of point hypotheses assume equal prior probabilities.

We may also directly compare them and plot the posterior distribution of their difference.

```
hyp <- "exp(sigma_Intercept + sigma_grouptreat) > exp(sigma_Intercept)"
(hyp <- hypothesis(fit1, hyp))
```

Hypothesis Tests for class b:

	Hypothesis	Estimate	Est.Error	CI.Lower	CI.Upper	Evid.Ratio	Post.Prob
Star							

```
1 (exp(sigma_Interc... > 0      1.18      0.32      0.7      1.73      Inf      1
*
```

```
---
```

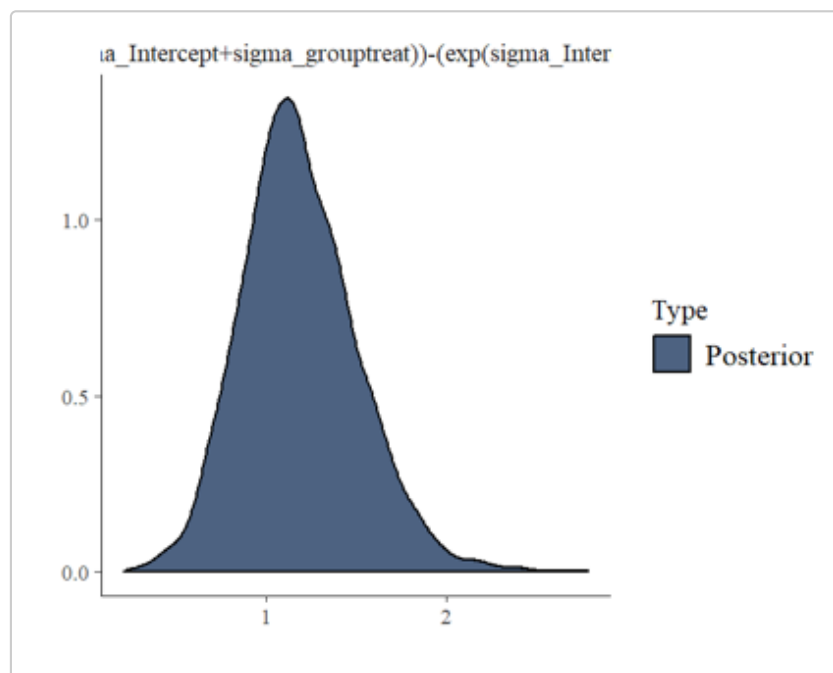
'CI': 90%-CI for one-sided and 95%-CI for two-sided hypotheses.

'\*': For one-sided hypotheses, the posterior probability exceeds 95%;

for two-sided hypotheses, the value tested against lies outside the 95%-CI.

Posterior probabilities of point hypotheses assume equal prior probabilities.

```
plot(hyp, chars = NULL)
```



Indeed, the residual standard deviation of the treatment group seems to be larger than that of the placebo group. Moreover, the magnitude of this difference is pretty similar to what we expected due to the values we put into the data simulations.

## Zero-Inflated Models

Another important application of the distributional regression framework are so-called zero-inflated models. These models are helpful whenever there are more zeros in the response variable than one would naturally expect. For example, if one seeks to predict the number of cigarettes people smoke per day and also includes non-smokers, there will be a huge amount of zeros which, when not modeled appropriately, can seriously distort parameter estimates. Here, we consider an example dealing with the number of fish caught by various groups of people. On the UCLA website (<https://www.ucla.edu/~biostatistics/>), the data are described as follows: “The state wildlife biologists want to model how many fish are being caught by fishermen at a state park. Visitors are asked how long they stayed, how many people were in the group, were there children in the group and how many fish were caught. Some visitors do not fish, but there is no data on whether a person fished or not. Some visitors who did fish did not catch any fish so there are excess zeros in the data because of the people that did not fish.”

```
zinp <- read.csv("https://paul-buerkner.github.io/data/fish.csv")
head(zinp)
```

	nofish	livebait	camper	persons	child	xb	zg	count
1	1	0	0	1	0	-0.8963146	3.0504048	0
2	0	1	1	1	0	-0.5583450	1.7461489	0
3	0	1	0	1	0	-0.4017310	0.2799389	0
4	0	1	1	2	1	-0.9562981	-0.6015257	0
5	0	1	0	1	0	0.4368910	0.5277091	1
6	0	1	1	4	2	1.3944855	-0.7075348	0

As predictors we choose the number of people per group, the number of children, as well as whether the group consists of campers. Many groups may not even try catching any fish at all (thus leading to many zero responses) and so we fit a zero-inflated Poisson model to the data. For now, we assume a constant zero-inflation probability across observations.

```
fit_zinb1 <- brm(count ~ persons + child + camper,
  data = zinb, family = zero_inflated_poisson())
```

Again, we summarize the results using the usual methods.

```
summary(fit_zinb1)
```

```
Family: zero_inflated_poisson
Links: mu = log; zi = identity
Formula: count ~ persons + child + camper
Data: zinb (Number of observations: 250)
Samples: 4 chains, each with iter = 2000; warmup = 1000; thin = 1;
         total post-warmup samples = 4000
```

Population-Level Effects:

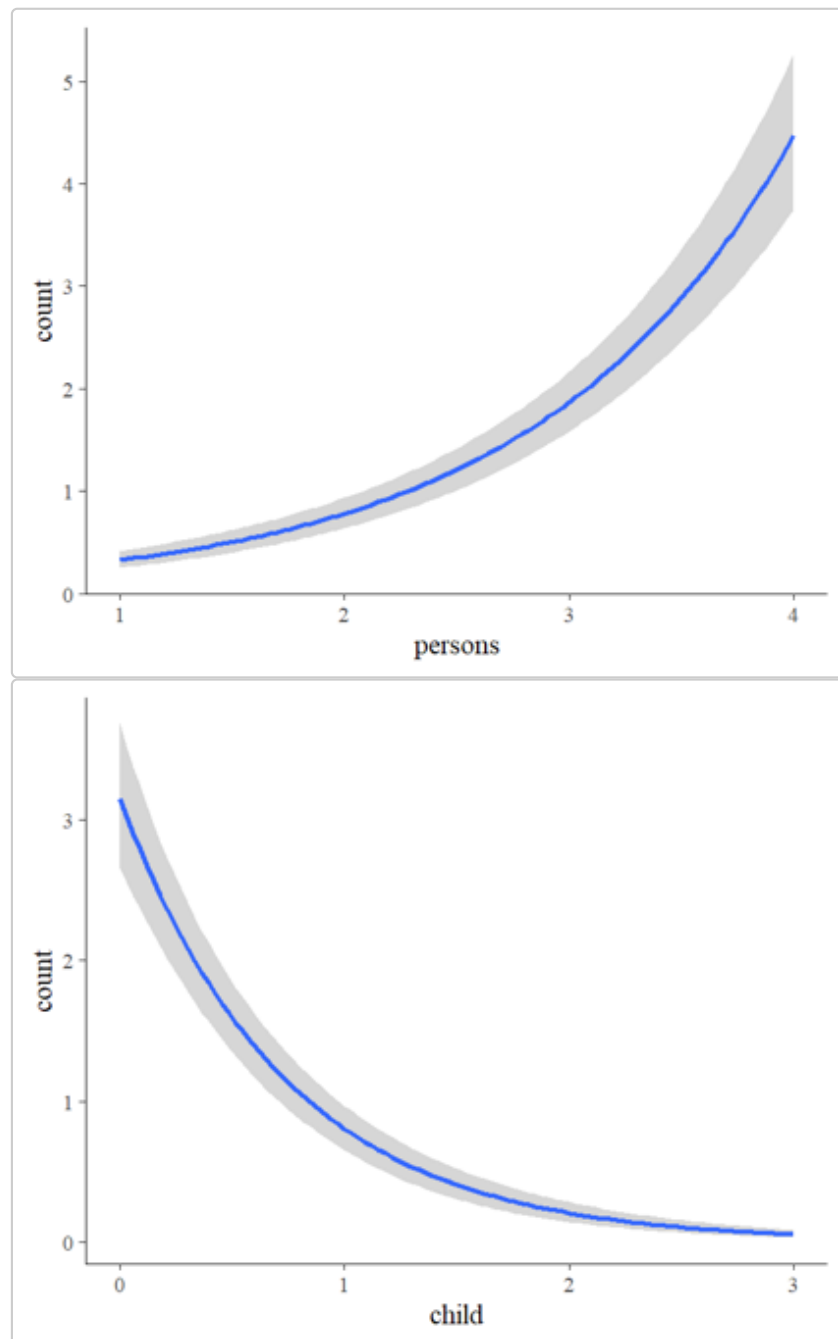
	Estimate	Est.Error	l-95% CI	u-95% CI	Eff.Sample	Rhat
Intercept	-1.01	0.18	-1.37	-0.67	2810	1.00
persons	0.87	0.05	0.78	0.97	2780	1.00
child	-1.36	0.09	-1.55	-1.19	2786	1.00
camper	0.80	0.10	0.61	0.99	3453	1.00

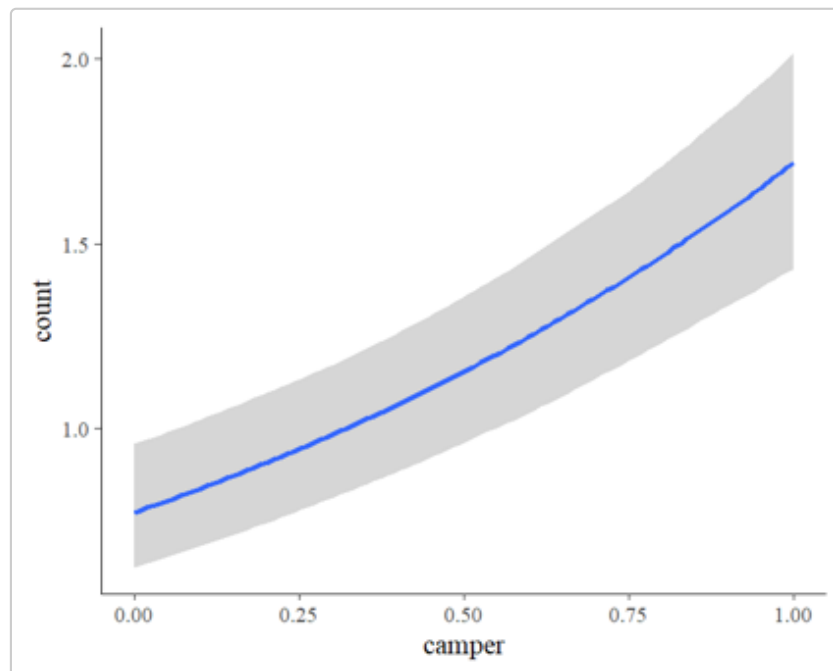
Family Specific Parameters:

	Estimate	Est.Error	l-95% CI	u-95% CI	Eff.Sample	Rhat
zi	0.41	0.04	0.32	0.49	2825	1.00

Samples were drawn using sampling(NUTS). For each parameter, Eff.Sample is a crude measure of effective sample size, and Rhat is the potential scale reduction factor on split chains (at convergence, Rhat = 1).

```
plot(marginal_effects(fit_zinb1), ask = FALSE)
```





According to the parameter estimates, larger groups catch more fish, campers catch more fish than non-campers, and groups with more children catch less fish. The zero-inflation probability  $z_i$  is pretty large with a mean of 41%. Please note that the probability of catching no fish is actually higher than 41%, but parts of this probability are already modeled by the poisson distribution itself (hence the name *zero-inflation*). If you want to treat all zeros as originating from a separate process, you can use hurdle models instead (not shown here).

Now, we try to additionally predict the zero-inflation probability by the number of children. The underlying reasoning is that we expect groups with more children to not even try catching fish. Most children are just terribly bad at waiting for hours until something happens. From a purely statistical perspective, zero-inflated (and hurdle) distributions are a mixture of two processes and predicting both parts of the model is natural and often very reasonable to make full use of the data.

```
fit_zinb2 <- brm(bf(count ~ persons + child + camper, zi ~ child),
  data = zinb, family = zero_inflated_poisson())
```

```
summary(fit_zinb2)
```

```
Family: zero_inflated_poisson
Links: mu = log; zi = logit
Formula: count ~ persons + child + camper
         zi ~ child
Data: zinb (Number of observations: 250)
Samples: 4 chains, each with iter = 2000; warmup = 1000; thin = 1;
         total post-warmup samples = 4000
```

Population-Level Effects:

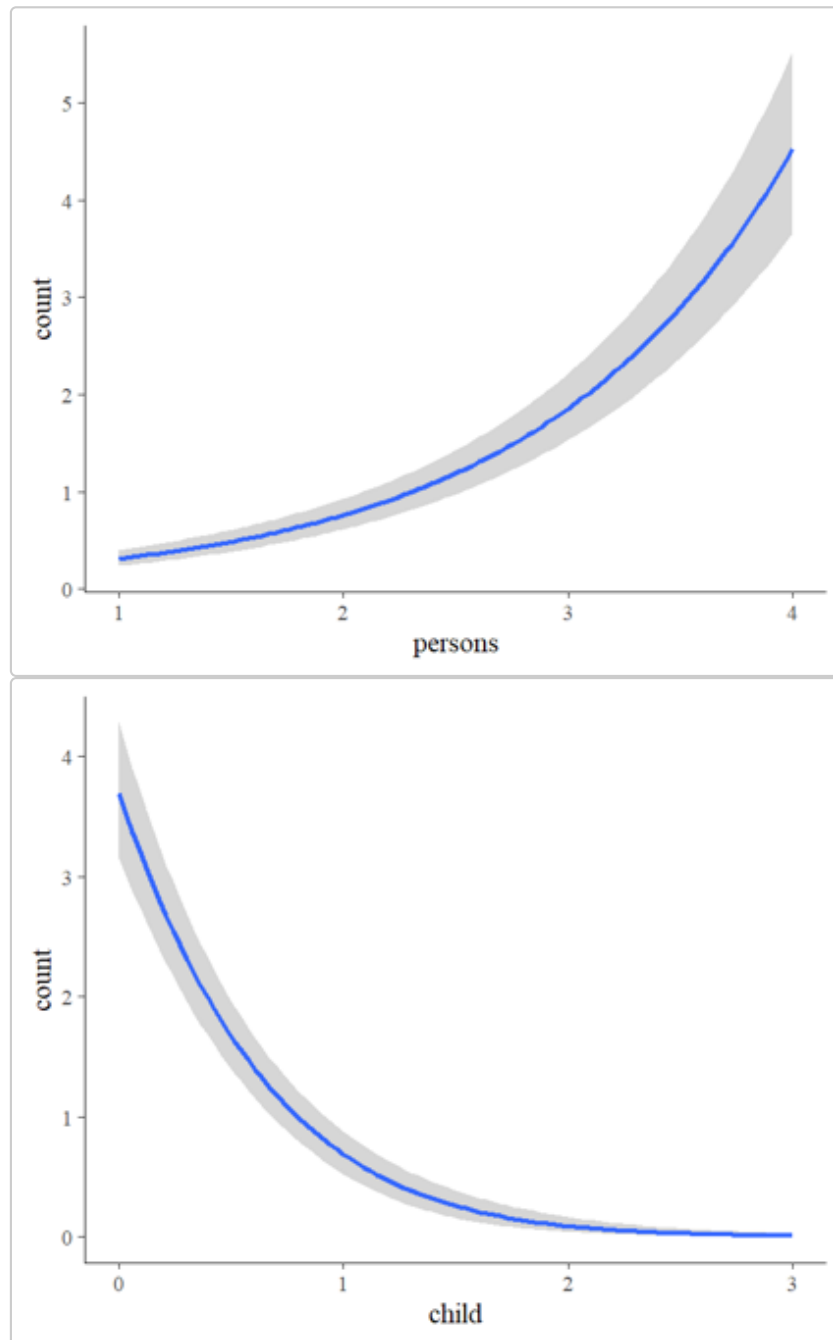
	Estimate	Est.Error	l-95% CI	u-95% CI	Eff.Sample	Rhat
Intercept	-1.08	0.18	-1.43	-0.72	3247	1.00
zi_Intercept	-0.96	0.26	-1.48	-0.48	3895	1.00
persons	0.89	0.05	0.80	0.98	3336	1.00

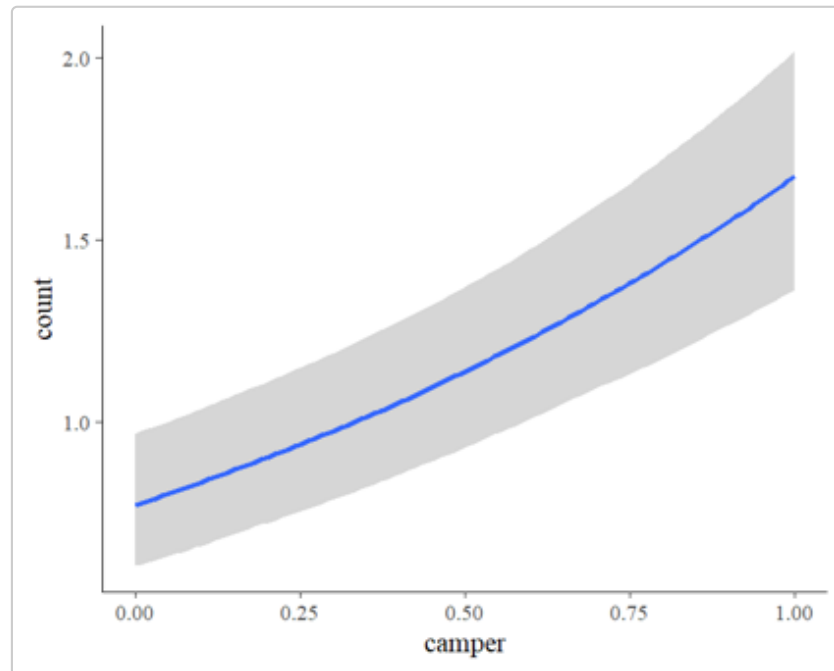


child	-1.17	0.10	-1.37	-0.99	3028	1.00
camper	0.78	0.09	0.60	0.96	3560	1.00
zi_child	1.22	0.28	0.67	1.79	3716	1.00

Samples were drawn using `sampling(NUTS)`. For each parameter, `Eff.Sample` is a crude measure of effective sample size, and `Rhat` is the potential scale reduction factor on split chains (at convergence, `Rhat` = 1).

```
plot(marginal_effects(fit_zinb2), ask = FALSE)
```





To transform the linear predictor of  $z_i$  into a probability, **brms** applies the logit-link:

$$\text{logit}(z_i) = \log\left(\frac{z_i}{1 - z_i}\right) = \eta_{z_i}$$

The logit-link takes values within  $[0, 1]$  and returns values on the real line. Thus, it allows the transition between probabilities and linear predictors.

According to the model, trying to fish with children not only decreases the overall number fish caught (as implied by the Poisson part of the model) but also drastically increases your change of catching no fish at all (as implied by the zero-inflation part) most likely because groups with more children are not even trying.

## Additive Distributional Models

In the examples so far, we did not have multilevel data and thus did not fully use the capabilities of the distributional regression framework of **brms**. In the example presented below, we will not only show how to deal with multilevel data in distributional models, but also how to incorporate smooth terms (i.e., splines) into the model. In many applications, we have no or only a very vague idea how the relationship between a predictor and the response looks like. A very flexible approach to tackle this problems is to use splines and let them figure out the form of the relationship. For illustration purposes, we simulate some data with the **mgcv** package, which is also used in **brms** to prepare smooth terms.

```
dat_smooth <- mgcv::gamSim(eg = 6, n = 200, scale = 2, verbose = FALSE)
```

```
Gu & Wahba 4 term additive model
```

```
head(dat_smooth[, 1:6])
```

	y	x0	x1	x2	x3	f
1	14.97398	0.5232871	0.70629564	0.2262371	0.9690609	18.01244
2	15.96906	0.1123638	0.78129331	0.5841745	0.3972866	14.56625
3	14.20430	0.2455518	0.03149645	0.4091435	0.9495137	15.46879
4	20.53585	0.6689515	0.58721612	0.5755327	0.1757639	20.00760
5	13.02828	0.4920482	0.94637205	0.1115531	0.7530496	15.89057
6	20.18388	0.4832472	0.95064483	0.8235346	0.2171470	15.40557

The data contains the predictors  $x_0$  to  $x_3$  as well as the grouping factor  $fac$  indicating the nested structure of the data. We predict the response variable  $y$  using smooth terms of  $x_1$  and  $x_2$  and a varying intercept of  $fac$ . In addition, we assume the residual standard deviation  $\sigma$  to vary by a smoothing term of  $x_0$  and a varying intercept of  $fac$ .

```
fit_smooth1 <- brm(
  bf(y ~ s(x1) + s(x2) + (1|fac), sigma ~ s(x0) + (1|fac)),
  data = dat_smooth, family = gaussian(),
  chains = 2, control = list(adapt_delta = 0.95)
)
```

```
summary(fit_smooth1)
```

```
Family: gaussian
Links: mu = identity; sigma = log
Formula: y ~ s(x1) + s(x2) + (1 | fac)
         sigma ~ s(x0) + (1 | fac)
Data: dat_smooth (Number of observations: 200)
Samples: 2 chains, each with iter = 2000; warmup = 1000; thin = 1;
         total post-warmup samples = 2000
```

Smooth Terms:

	Estimate	Est.Error	l-95% CI	u-95% CI	Eff.Sample	Rhat
sds(sx1_1)	4.32	2.94	0.99	12.16	715	1.01
sds(sx2_1)	23.88	6.44	14.46	39.44	775	1.00
sds(sigma_sx0_1)	1.13	1.28	0.03	4.70	617	1.00

Group-Level Effects:

~fac (Number of levels: 4)

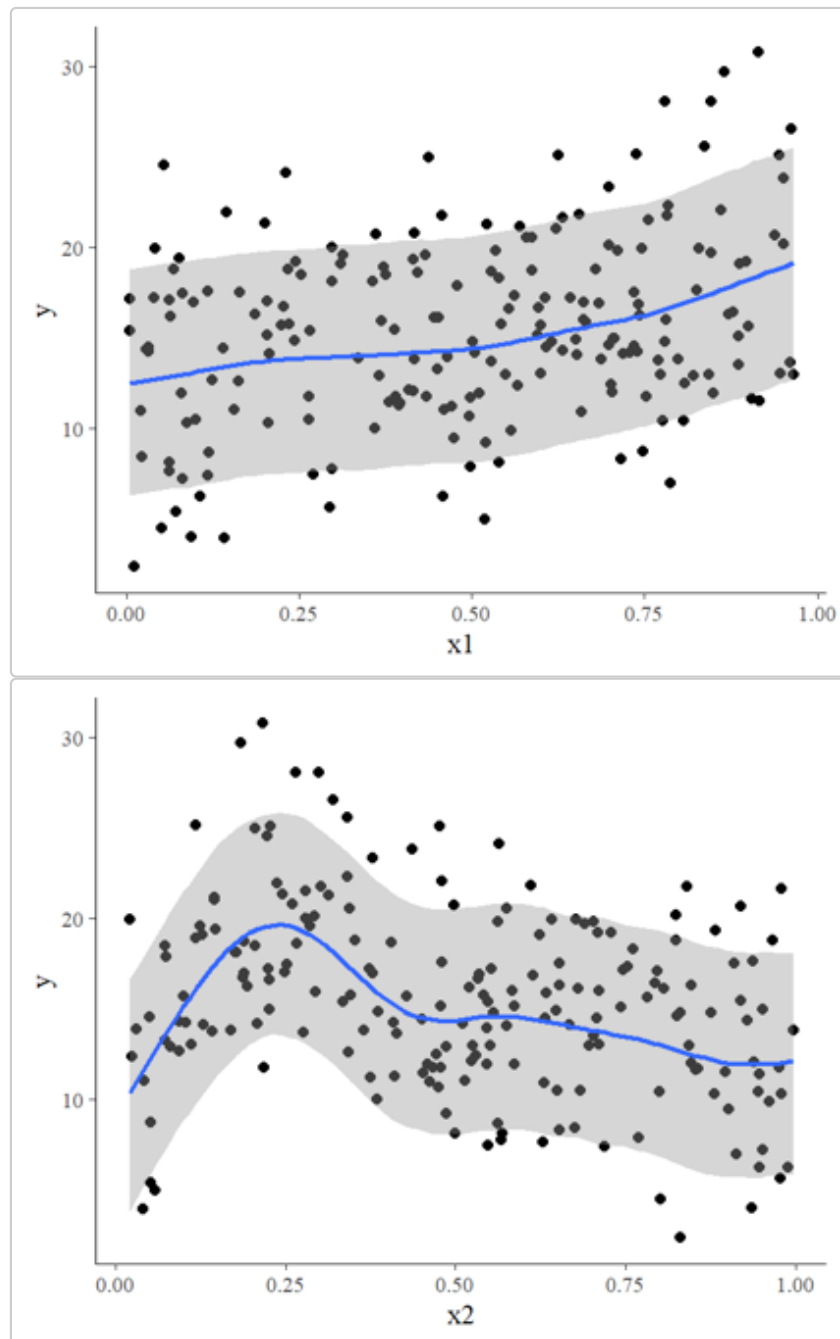
	Estimate	Est.Error	l-95% CI	u-95% CI	Eff.Sample	Rhat
sd(Intercept)	6.01	2.99	2.58	14.31	945	1.00
sd(sigma_Intercept)	0.19	0.25	0.01	0.79	662	1.00

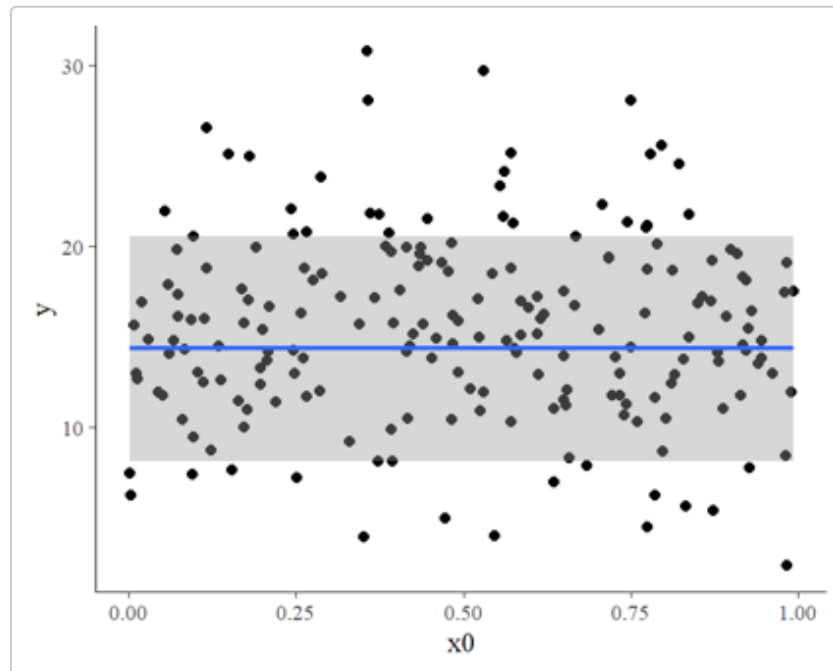
Population-Level Effects:

	Estimate	Est.Error	l-95% CI	u-95% CI	Eff.Sample	Rhat
Intercept	15.51	2.96	9.29	21.72	840	1.00
sigma_Intercept	0.70	0.15	0.41	0.99	696	1.00
sx1_1	15.89	6.98	3.08	31.71	1423	1.00
sx2_1	67.28	18.40	32.81	104.35	1698	1.00
sigma_sx0_1	0.88	2.08	-3.21	5.67	1008	1.00

Samples were drawn using `sampling(NUTS)`. For each parameter, `Eff.Sample` is a crude measure of effective sample size, and `Rhat` is the potential scale reduction factor on split chains (at convergence, `Rhat` = 1).

```
plot(marginal_effects(fit_smooth1), points = TRUE, ask = FALSE)
```





This model is likely an overkill for the data at hand, but nicely demonstrates the ease with which one can specify complex models with **brms** and to fit them using **Stan** on the backend.