

## 03 Linear modeling

Shravan Vasishth and Bruno Nicenboim

25-29 March 2019

# Linear modeling

Suppose  $y$  is a vector of continuous responses; assume for now that it is coming from a normal distribution:

$$y \sim \text{Normal}(\mu, \sigma)$$

This is the simple linear model:

$$y = \mu + \varepsilon \text{ where } \varepsilon \sim \text{Normal}(0, \sigma)$$

There are two parameters,  $\mu, \sigma$ , so we need priors on these. We expand on this simple model next.

# Linear modeling

Recall from the foundations lecture that the way we will conduct data analysis is as follows.

- Given data, specify a *likelihood function*.
- Specify *prior distributions* for model parameters.
- Evaluate whether model makes sense, using fake-data simulation, *prior predictive* and *posterior predictive* checks, and (if you want to claim a discovery) calibrating true and false discovery rates.
- Using software, derive *marginal posterior distributions* for parameters given likelihood function and prior density. I.e., simulate parameters to get *samples from posterior distributions* of parameters using some *Markov Chain Monte Carlo (MCMC) sampling algorithm*.
- Check that the model converged using *model convergence* diagnostics,
- Summarize *posterior distributions* of parameter samples and make your scientific decision.

We will now work through some specific examples to illustrate how the data

## Example 1: A single subject pressing a button repeatedly

As a first example, we will fit a simple linear model to some reaction time data.

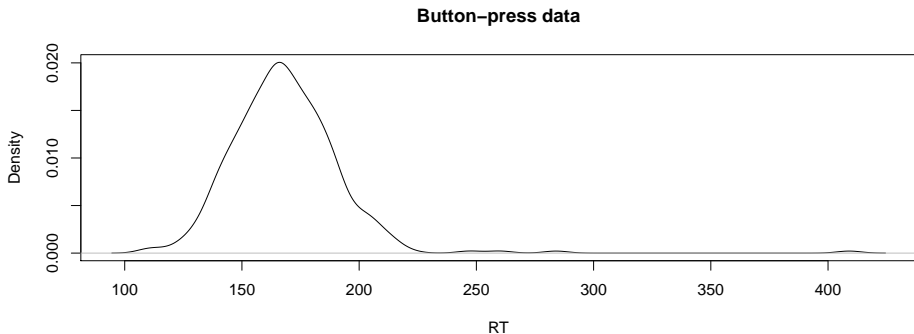
The file `button_press.dat` contains data of a subject pressing the space bar without reading in a self-paced reading experiment.

### Preprocessing of the data

```
##          type item wordn          word    y
## 356 filler      3      0    Vielleicht 214
## 357 filler      3      1         haben 182
## 358 filler      3      2  die_Zahnärztin 179
## 359 filler      3      3     aus_Bonn 177
## 360 filler      3      4  die_Patienten 183
## 361 filler      3      5    verklagt. 162
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##    110.0   156.0   166.0   168.6   181.0   409.0
```

# Visualizing the data

It is a good idea to look at the distribution of the data before doing anything else. See Figure 1.



**Figure 1:** Visualizing the data.

The data looks a bit skewed, but we ignore this for the moment.

# Define the likelihood function

Let's model the data with the following assumptions:

- There is a true underlying time,  $\mu$ , that the participant needs to press the space-bar.
- There is some noise in this process.
- The noise is normally distributed (this assumption is questionable given the skew but; we fix this assumption later).

# Define the likelihood function

This means that the likelihood for each observation  $i$  will be:

$$y_i \sim \text{Normal}(\mu, \sigma) \quad (1)$$

where  $i = 1 \dots N$ .

This is just the simple linear model:

$$y = \mu + \varepsilon \text{ where } \varepsilon \sim \text{Normal}(0, \sigma) \quad (2)$$

# Define the priors for the parameters

We are going to use the following priors for the two parameters in this model:

$$\begin{aligned}\mu &\sim \text{Normal}(0, 2000) \\ \sigma &\sim \text{Normal}(0, 500) \text{ truncated so that } \sigma > 0\end{aligned}\tag{3}$$



# Define the priors for the parameters

In order to decide on a prior for the parameters, always visualize them first. See Figure 2.



**Figure 2:** Visualizing the priors for example 1.

# Prior predictive checks

With these priors, we are going to generate something called the **prior predictive distribution**. This helps us check whether the priors make sense.

Formally, we want to know the density  $f(\cdot)$  of data points  $y_1, \dots, y_n$ , given a vector of priors  $\Theta$ . In our example,  $\Theta = \langle \mu, \sigma \rangle$ . The prior predictive density is:

$$f(y_1, \dots, y_n) = \int f(y_1) \cdot f(y_2) \cdots f(y_n) f(\Theta) d\Theta \quad (4)$$

# Prior predictive checks

In essence, we integrate out the parameters. Here is one way to do it in R:

- Take one sample from each of the prior distributions
- Generate data using those samples
- Repeat until you have a substantial amount of data
- Plot the prior predictive density of the generated data

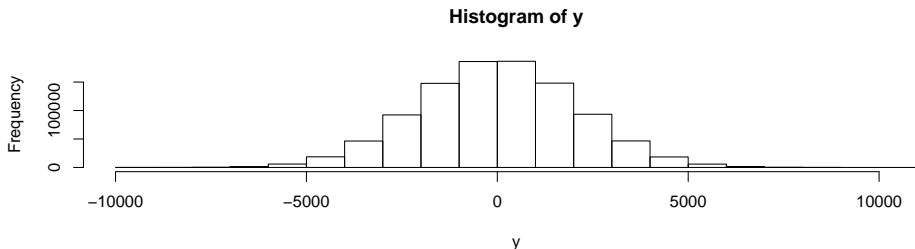
# Prior predictive checks

```
nsim<-1000000
y<-rep(NA,nsim)
mu<-rnorm(nsim,mean=0,sd=2000)
sigma<-abs(rnorm(nsim,mean=0,sd=500))

for(i in 1:nsim){
  y[i]<-rnorm(1,mean=mu[i],sd=sigma[i])
}

hist(y)
```

# Prior predictive checks



**Figure 3:** First attempt at prior predictive distribution of the data, model m1.

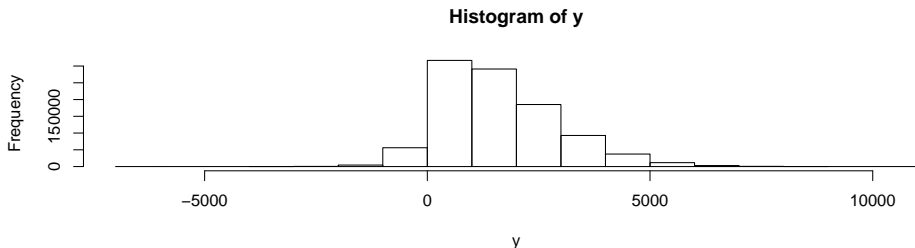
# Prior predictive checks

We can try to redefine the prior for  $\mu$  to have only positive values, and then check again (Figure 4). We still get some negative values, but that is because we are assuming that

$$y \sim \text{Normal}(\mu, \sigma)$$

which will have negative values for small  $\mu$  and large  $\sigma$ .

# Prior predictive checks



**Figure 4:** Second attempt at prior predictive distribution of the data, model m1.

This prior predictive distribution in Figure 4 looks reasonable for now.

# Prior predictive checks

We can generate a prior predictive distribution using Stan as follows.

First, we define a Stan model that defines the priors and defines how the data are to be generated.

Documentation on Stan is available at [mc-stan.org](https://mc-stan.org).



# Prior predictive checks

```
priorpred<-"data {  
  int N;  
}  
parameters {  
  real<lower=0> mu;  
  real<lower=0> sigma;  
}  
model {  
  mu ~ normal(0,2000);  
  sigma ~ normal(0,500);  
}  
generated quantities {  
  vector[N] y_sim;  
  for(i in 1:N) {  
    y_sim[i] = normal_rng(mu,sigma);  
  }  
}
```

# Prior predictive checks

Load RStan and brms.

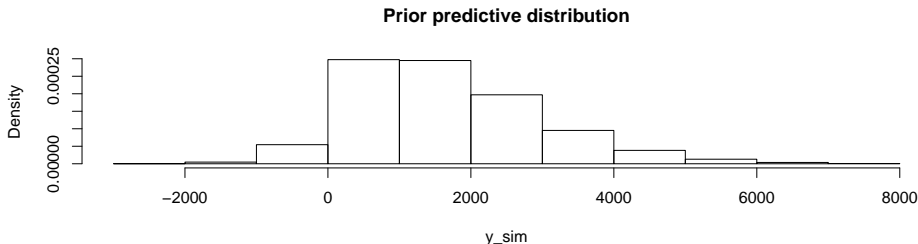
# Prior predictive checks

Then we generate the data:

```
## generate 100 data sets  
dat<-list(N=100)  
  
## fit model:  
m1priorpred<-stan(model_code=priorpred,  
                   data=dat,  
                   chains = 4,  
                   iter = 2000)
```

# Prior predictive checks

```
## extract and plot one of the data-sets:  
y_sim<-extract(m1priorpred,pars="y_sim")  
hist(y_sim$y_sim[,1],  
      main="Prior predictive distribution",  
      xlab="y_sim",freq=FALSE)
```



# Prior predictive checks

Having satisfied ourselves that the priors mostly make sense, we now fit the model to fake data. The goal here is to ensure that the model recovers the true underlying parameters.

# Fake-data simulation and modeling

Next, we write the Stan model, adding a likelihood in the model block:

# Fake-data simulation and modeling

```
m1<-"data {  
  int N;  
  real y[N]; // data  
}  
parameters {  
  real<lower=0> mu;  
  real<lower=0> sigma;  
}  
model {  
  mu ~ normal(0,2000);  
  sigma ~ normal(0,500);  
  y ~ normal(mu,sigma);  
}  
generated quantities {  
  vector[N] y_sim;  
  for(i in 1:N) {  
    y_sim[i] = normal_rng(mu,sigma);  
  }}  
"
```

# Fake-data simulation and modeling

Then generate fake data with known parameter values (we decide what these are):

```
set.seed(123)
N <- 500
true_mu <- 400
true_sigma <- 125
y <- rnorm(N, true_mu, true_sigma)

y <- round(y)
fake_data <- data.frame(y=y)
dat<-list(y=y,N=N)
```

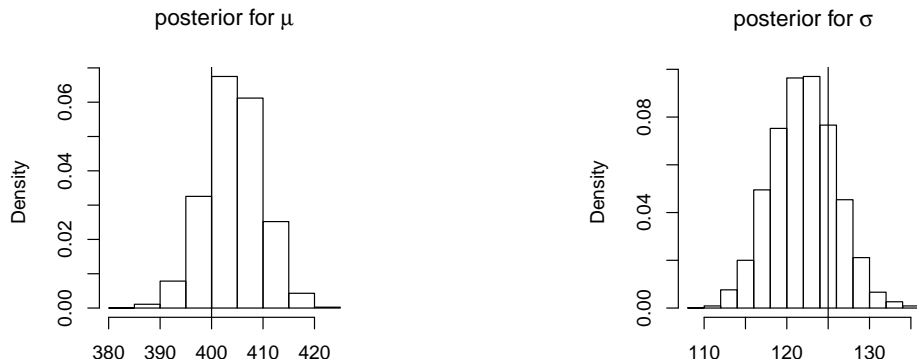


# Fake-data simulation and modeling

Finally, we fit the model:

```
## fit model:  
m1rstan<-stan(model_code=m1,  
              data=dat,  
              chains = 4,  
              iter = 2000)  
  
## extract posteriors:  
posteriors<-extract(m1rstan,pars=c("mu","sigma"))
```

# Fake-data simulation and modeling



**Figure 5:** Posteriors from fake data, model m1. Vertical lines show the true values of the parameters.

# Posterior predictive checks

Once we have the posterior distribution  $f(\Theta | y)$ , we can derive the predictions based on this posterior distribution:

$$p(y_{pred} | y) = \int p(y_{pred}, \Theta | y) d\Theta = \int p(y_{pred} | \Theta, y) p(\Theta | y) d\Theta \quad (5)$$

# Posterior predictive checks

Assuming that past and future observations are conditionally independent given  $\Theta$ , i.e.,  $p(y_{pred} | \Theta, y) = p(y_{pred} | \Theta)$ , we can write:

$$p(y_{pred} | y) = \int p(y_{pred} | \Theta) p(\Theta | y) d\Theta \quad (6)$$

Note that we are conditioning  $y_{pred}$  only on  $y$ , we do not condition on what we don't know ( $\Theta$ ); **we integrate out the unknown parameters.**

# Posterior predictive checks

This posterior predictive distribution is different from the frequentist approach, which gives only a predictive distribution of  $y_{pred}$  given our estimate of  $\theta$  (a point value).

In the Stan code above, we have already generated the posterior predictive distribution, in the generated quantities block.

# Implementing model in brms

This model is expressed in brms in the following way. First, define the priors:

```
priors <- c(set_prior("normal(0, 2000)",  
                     class = "Intercept"),  
            set_prior("normal(0, 500)",  
                     class = "sigma"))
```

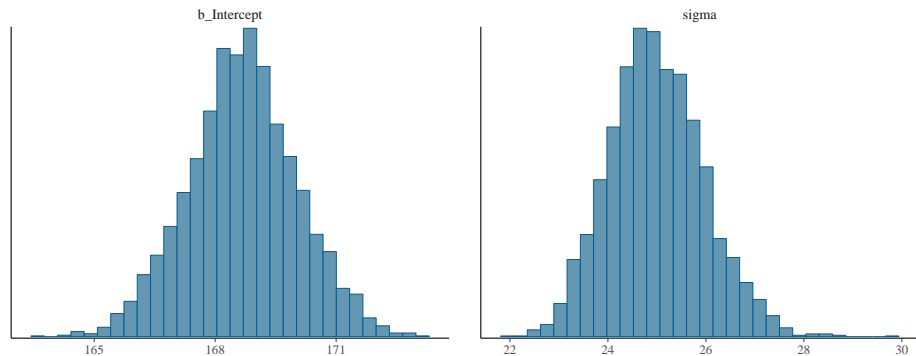
# Implementing model in brms

Then, define the generative process assumed:

```
m1brms<-brm(y~1,noreading_data,prior = priors,  
  iter = 2000,  
  warmup = 1000,  
  chains = 4,  
  family = gaussian(),  
  control = list(adapt_delta = 0.99))
```

# Summarizing the posteriors, and convergence diagnostics

A graphical summary of posterior distributions of model m1 is shown in Figure 6:

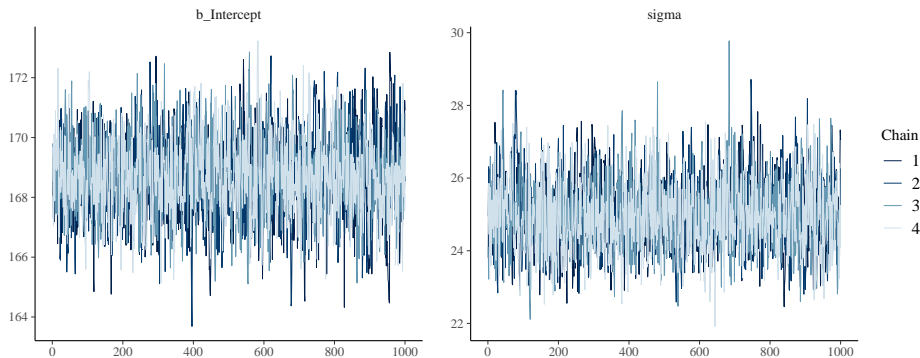


**Figure 6:** Posterior distributions of the parameters in model m1.



# Summarizing the posteriors, and convergence diagnostics

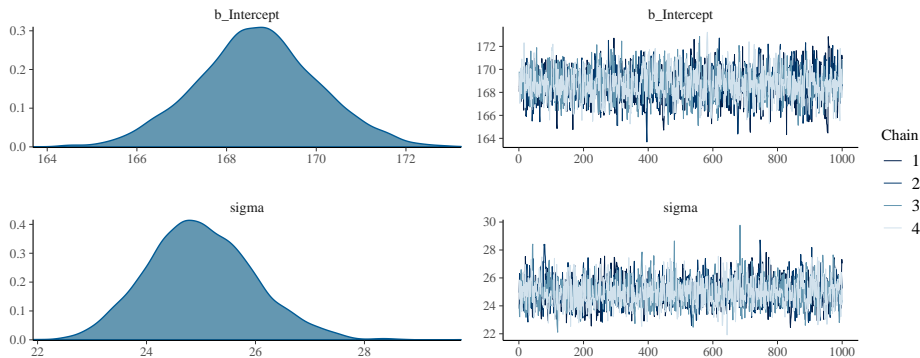
The trace plots in Figure 7 show how well the four chains are mixing:



**Figure 7:** Trace plots in model m1.

# Summarizing the posteriors, and convergence diagnostics

An alternative way to plot is shown in Figure 8.



**Figure 8:** Posterior distributions and trace plots in model m1.

# Fitting the brms model on fake data

```
m1_fakebrms<-brm(y~1,fake_data,prior = priors,  
  iter = 2000, chains = 4,family = gaussian(),  
  control = list(adapt_delta = 0.99))
```

## Summarizing the posterior distribution: posterior probabilities and the credible interval

We are assuming that there's a true underlying time it takes to press the space bar,  $\mu$ , and there is normally distributed noise with distribution  $\text{Normal}(0, \sigma)$  that generates the different RTs. All this is encoded in our likelihood by assuming that RTs are distributed with an unknown true mean  $\mu$  (and an unknown standard deviation  $\sigma$ ).

## Summarizing the posterior distribution: posterior probabilities and the credible interval

The objective of the Bayesian model is to learn about the plausible values of  $\mu$ , or in other words, to get a distribution that encodes what we know about the true mean of the distribution of RTs, and about the true standard deviation,  $\sigma$ , of the distribution of RTs.

# Summarizing the posterior distribution: posterior probabilities and the credible interval

Our model allows us to answer questions such as:

**What is the probability that the underlying value of the mindless press of the space bar would be over, say 170 ms?**

# Summarizing the posterior distribution: posterior probabilities and the credible interval

As an example, consider this model that we ran above.

```
priors <- c(set_prior("normal(0, 2000)",  
                    class = "Intercept"),  
           set_prior("normal(0, 500)",  
                    class = "sigma"))  
  
m1brms<-brm(y~1,noreading_data,prior = priors,  
            iter = 2000,  
            warmup = 1000,  
            chains = 4,  
            family = gaussian(),  
            control = list(adapt_delta = 0.99))
```

# Summarizing the posterior distribution: posterior probabilities and the credible interval

We now compute the posterior probability  $Prob(\mu > 170)$ :

```
mu_post<-posterior_samples(m1brms,  
                             pars=c("b_Intercept"))$b_Intercept  
mean(mu_post>170)  
  
## [1] 0.14875
```



# Summarizing the posterior distribution: posterior probabilities and the credible interval

## The credible interval

The 95% credible interval can be extracted for  $\mu$  as follows:

```
posterior_interval(m1brms, pars=c("b_Intercept"))
```

```
##                2.5%    97.5%  
## b_Intercept 166.1305 171.1494
```

This type of interval is also known as a *credible interval*.

A credible interval demarcates the range within which we can be certain with a certain probability that the “true value” of a parameter lies given the data and the model.

This is very different from the frequentist confidence interval!

# Summarizing the posterior distribution: posterior probabilities and the credible interval

The percentile interval is a type of credible interval (the most common one), where we assign equal probability mass to each tail.

We generally report 95% credible intervals. But we can extract any interval, a 73% interval, for example, leaves 13.5% of the probability mass on each tail, and we can calculate it like this:

```
round(quantile(mu_post,prob=c(0.135,0.865)))
```

```
## 13.5% 86.5%
```

```
##    167    170
```

# Influence of priors and sensitivity analysis

$$\begin{aligned}\mu &\sim \text{Uniform}(0, 5000) \\ \sigma &\sim \text{Uniform}(0, 500)\end{aligned}\tag{7}$$

```
priors <- c(set_prior("uniform(0, 5000)",  
                     class = "Intercept"),  
            set_prior("normal(0, 500)",  
                     class = "sigma"))
```

# Influence of priors and sensitivity analysis

```
m2<-brm(y~1,noreading_data,prior = priors,  
        iter = 2000, chains = 4,family = gaussian(),  
        control = list(adapt_delta = 0.99))
```

# Influence of priors and sensitivity analysis

```
summary(m2)
```

```
## Family: gaussian
## Links: mu = identity; sigma = identity
## Formula: y ~ 1
## Data: noreading_data (Number of observations: 361)
## Samples: 4 chains, each with iter = 2000; warmup = 1000; th
##           total post-warmup samples = 4000
##
## Population-Level Effects:
##           Estimate Est.Error l-95% CI u-95% CI Eff.Sample R
## Intercept    168.68      1.29   166.20   171.23     2176 1
##
## Family Specific Parameters:
##           Estimate Est.Error l-95% CI u-95% CI Eff.Sample Rhat
## sigma        24.98      0.93    23.25    26.89     2165 1.00
##
```

# Influence of priors and sensitivity analysis

In general, we don't want our priors to have too much influence on our posterior.

This is unless we have *very* good reasons for having informative priors, such as a very small sample and a lot of prior information; an example would be if we have data from an impaired population, which makes it hard to increase our sample size.

# Influence of priors and sensitivity analysis

We usually center the priors on 0 and we let the likelihood dominate in determining the posterior.

This type of prior is called *weakly informative prior*. Notice that a uniform prior is not a weakly informative prior, it assumes that every value is equally likely, zero is as likely as 5000.

You should always do a *sensitivity analysis* to check how influential the prior is: try different priors and verify that the posterior doesn't change drastically.

## Example 2: Investigating adaptation effects

More realistically, we might have run the small experiment to find out whether the participant tended to speedup (practice effect) or slowdown (fatigue effect) while pressing the space bar.



## Example 2: Investigating adaptation effects

### Preprocessing the data

- We need to have data about the number of times the space bar was pressed for each observation, and add it to our list.
- It's a good idea to center the number of presses (a covariate) to have a clearer interpretation of the intercept.
- In general, centering predictors is always a good idea, for interpretability and for computational reasons.
- See @SchadEtAlcontrasts for details on this point.

## Example 2: Investigating adaptation effects

### Preprocessing the data

```
# create a new vector representing trial id in the data frame  
noreading_data$presses <- 1:nrow(noreading_data)  
# center the vector  
noreading_data$c_presses <- noreading_data$presses -  
  mean(noreading_data$presses)
```

## Example 2: Investigating adaptation effects

### Probability model

Our model changes, because we have a new parameter.

$$y_i \sim \text{Normal}(\alpha + \text{presses}_i \cdot \beta, \sigma) \quad (8)$$

where  $i = 1 \dots N$

And we are going to use the following priors:

$$\begin{aligned} \alpha &\sim \text{Normal}(0, 2000) \\ \beta &\sim \text{Normal}(0, 500) \\ \sigma &\sim \text{Normal}(0, 500) \text{ truncated so that } \sigma > 0 \end{aligned} \quad (9)$$

## Example 2: Investigating adaptation effects

### Probability model

We are basically fitting a linear model,  $\alpha$  represents the intercept (namely, the grand mean of the RTs), and  $\beta$  represents the slope.

What information are the priors encoding?

Do the priors make sense?

## Example 2: Investigating adaptation effects

### Probability model

We'll write this in brms as follows.

```
priors <- c(set_prior("normal(0, 2000)",  
                    class = "Intercept"),  
            set_prior("normal(0, 500)",  
                    class = "b",  
                    coef="presses"),  
            set_prior("normal(0, 500)",  
                    class = "sigma"))
```

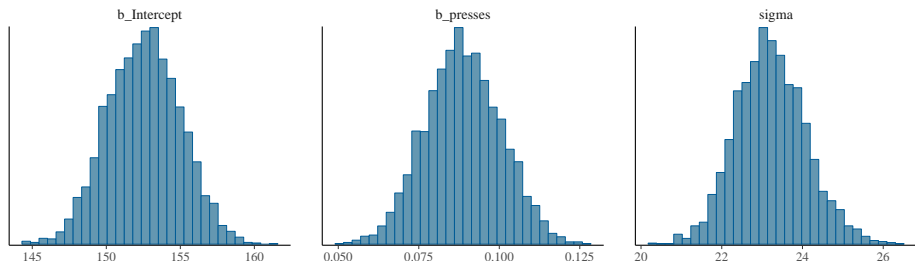
## Example 2: Investigating adaptation effects

### Probability model

```
m2<-brm(y~1+presses,noreading_data,prior = priors,  
        iter = 2000, chains = 4,family = gaussian(),  
        control = list(adapt_delta = 0.99))  
  
## Compiling the C++ model  
## Start sampling
```

## Example 2: Investigating adaptation effects

### Posteriors



## Example 2: Investigating adaptation effects

### Summarizing the posterior and inference

We'll need to examine what happens with  $\beta$ . The summary gives us the relevant information:

```
m2_post_samp_b <- posterior_samples(m2, "^b")
beta_samples <- m2_post_samp_b$b_presses
beta_mean <- mean(beta_samples)
quantiles_beta <- quantile(beta_samples,
                           prob=c(0.025, 0.975))
beta_low <- quantiles_beta[1]
beta_high <- quantiles_beta[2]
```



## Example 2: Investigating adaptation effects

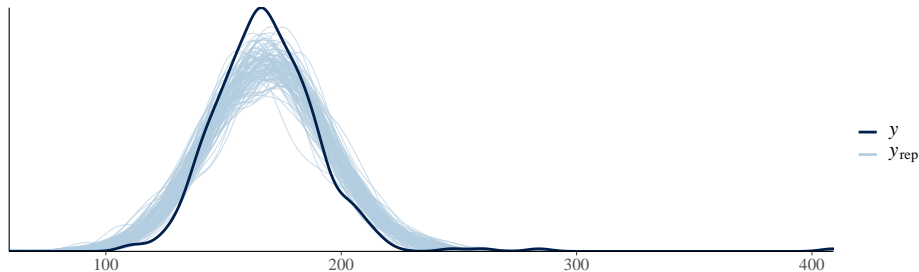
### Posterior predictive checks

Let's say we know that our model is working as expected, since we already used fake data to test the recovery of the parameters.

## Example 2: Investigating adaptation effects

### Posterior predictive checks

To do posterior predictive checks for our last example, using brms, we need to do:

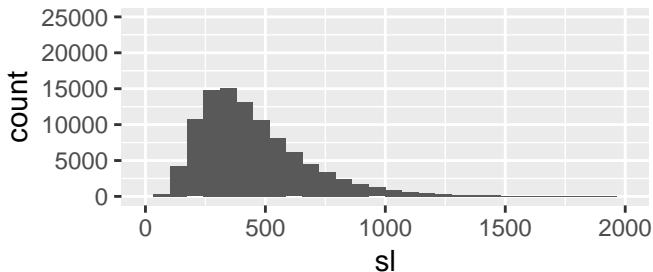


**Figure 9:** Posterior predictive check of model m2.

## Example 2: Investigating adaptation effects

### Using the log-normal likelihood

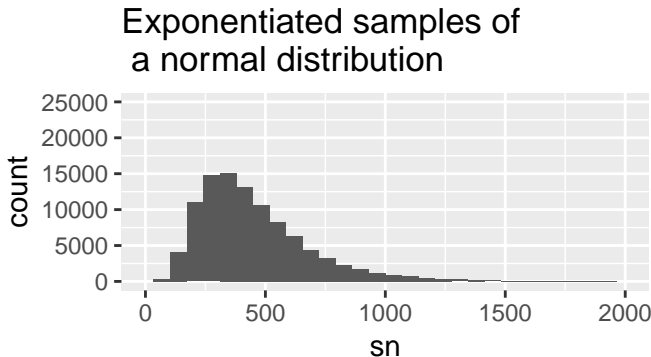
Log-normal distribution



**Figure 10:** The log-normal distribution.

## Example 2: Investigating adaptation effects

### Using the log-normal likelihood



**Figure 11:** Exponentiated samples from a log-normal distribution.

## Example 2: Investigating adaptation effects

### Re-fit the model assuming a log-normal likelihood

If we assume that RTs are log-normally distributed, we'll need to change our model:

$$Y_i \sim \text{LogNormal}(\alpha + \text{presses}_i \cdot \beta, \sigma) \quad (10)$$

where  $i = 1 \dots N$

But now the scale of our priors needs to change! They are no longer in milliseconds.

$$\begin{aligned} \alpha &\sim \text{Normal}(0, 10) \\ \beta &\sim \text{Normal}(0, 1) \\ \sigma &\sim \text{Normal}(0, 2) \text{ truncated so that } \sigma > 0 \end{aligned} \quad (11)$$

## Example 2: Investigating adaptation effects

Re-fit the model assuming a log-normal likelihood

```
priors_log <- c(set_prior("normal(0, 10)",  
                        class = "Intercept"),  
               set_prior("normal(0, 1)",  
                        class = "b",  
                        coef="presses"),  
               set_prior("normal(0, 2)",  
                        class = "sigma"))  
  
m2_logn<-brm(y~1+presses,noreading_data,  
             prior = priors_log,  
             iter = 2000, chains = 4,family = lognormal(),  
             control = list(adapt_delta = 0.99,  
                             max_treedepth=15))  
  
## Compiling the C++ model  
## Start sampling
```

## Example 2: Investigating adaptation effects

### Summarizing the posterior and inference

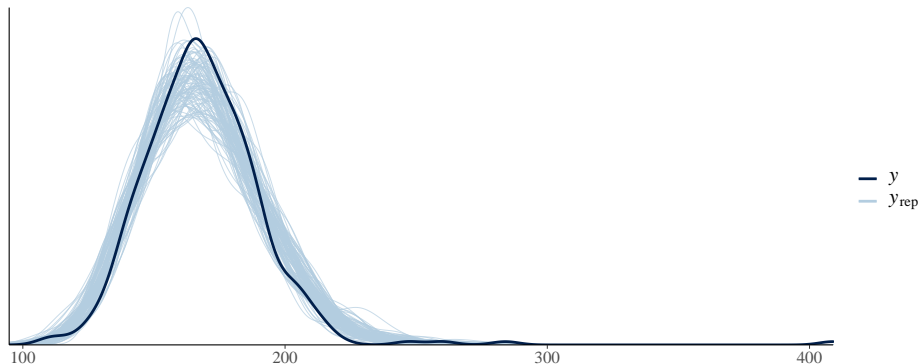
Next, we turn to the question of what we can report as our results, and what we can conclude from the data.

- We can summarize the posterior and do inference as discussed in Example 1.
- If we want to talk about the effect estimated by the model, we summarize the posterior of  $\beta$  in the following way:
- $\hat{\beta} = 0.08$ , 95% CrI =  $[0.063, 0.097]$ ,  $P(\beta > 0) \approx 1$

## Example 2: Investigating adaptation effects

### Posterior predictive checks and distribution of summary statistics

We can now verify whether our predicted datasets look more similar to the real dataset. See Figure 10.



**Figure 12:** Posterior predictive check.



## Example 2: Investigating adaptation effects

### Posterior predictive checks and distribution of summary statistics

```
m2_logn<-brm(y~1+presses,noreading_data,  
             prior = priors_log,  
             iter = 2000, chains = 4,  
             family = lognormal(),  
             control = list(adapt_delta = 0.99,  
                             max_treedepth=15))
```

# General workflow

This is the general workflow that we suggest for a Bayesian model.

- 1 Define the full probability model:
  - a. Decide on the likelihood.
  - b. Decide on the priors.
  - c. Write the brms or Stan model.
- 2 Do prior predictive checks to determine if priors make sense.
- 3 Check model using fake data simulations:
  - a. Simulate data with known values for the parameters.
  - b. Fit the model and do MCMC diagnostics.
  - c. Verify that it recovers the parameters from simulated data.
- 4 Fit the model with real data and do MCMC diagnostics.
- 5 Evaluate the model's fit (e.g., posterior predictive checks, distribution of summary statistics). This may send you back to 1.
- 6 Inference/prediction/decisions.
- 7 Conduct model comparison if there's an alternative model (to be discussed later).