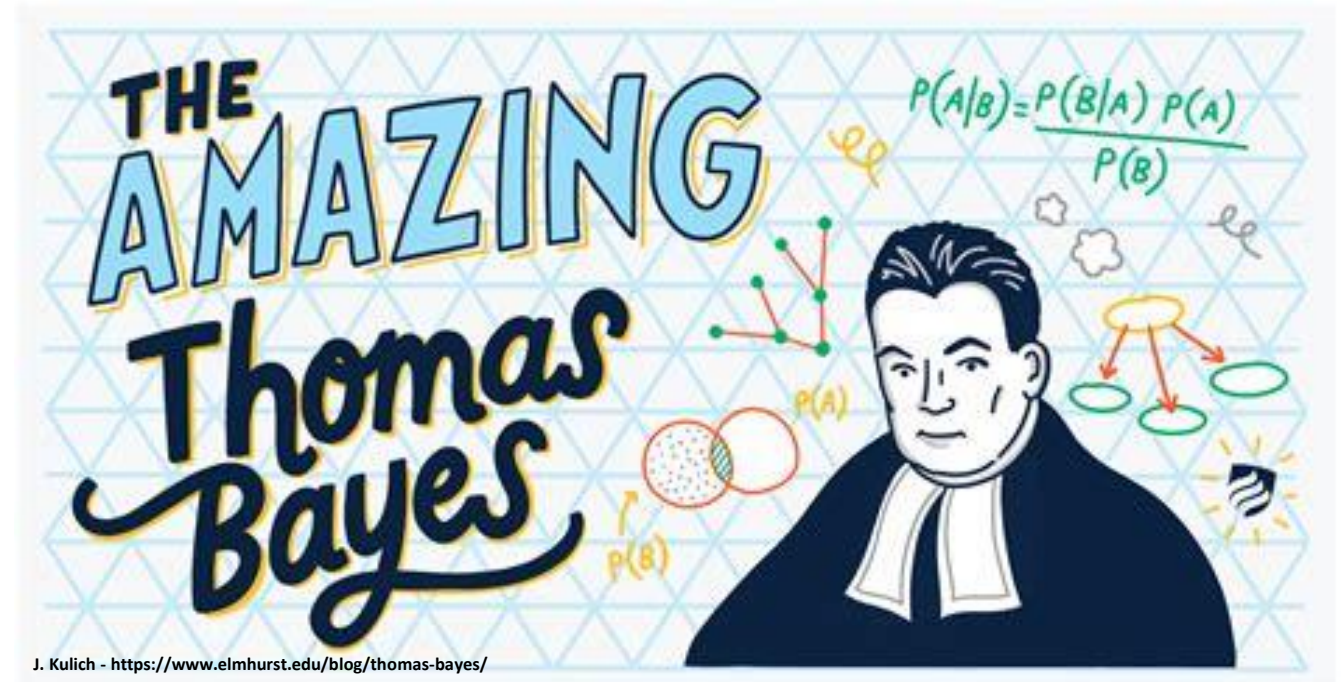


Bayesian stats with Stan and brms ()

Béatrice Capolla & Florent Déry



An immersion into british philosophy

David Hume: a great skeptic

Criticizes Descartes' rationalism:
our ideas cannot be confirmed by immediate
perceptions

1711 ———— 1776

Hume's fork:

**Demonstrative
statement**

—
irrefutable
(a priori
knowledge),
 $2 + 2 = 4$

**Probable
statement**

—
refutable ,
needs
empirical proof



Georgios Magkakis.

- A Treatise of Human Nature: Being an Attempt to introduce the experimental Method of Reasoning into Moral Subjects (1739)
- Enquiries concerning Human Understanding (1748)

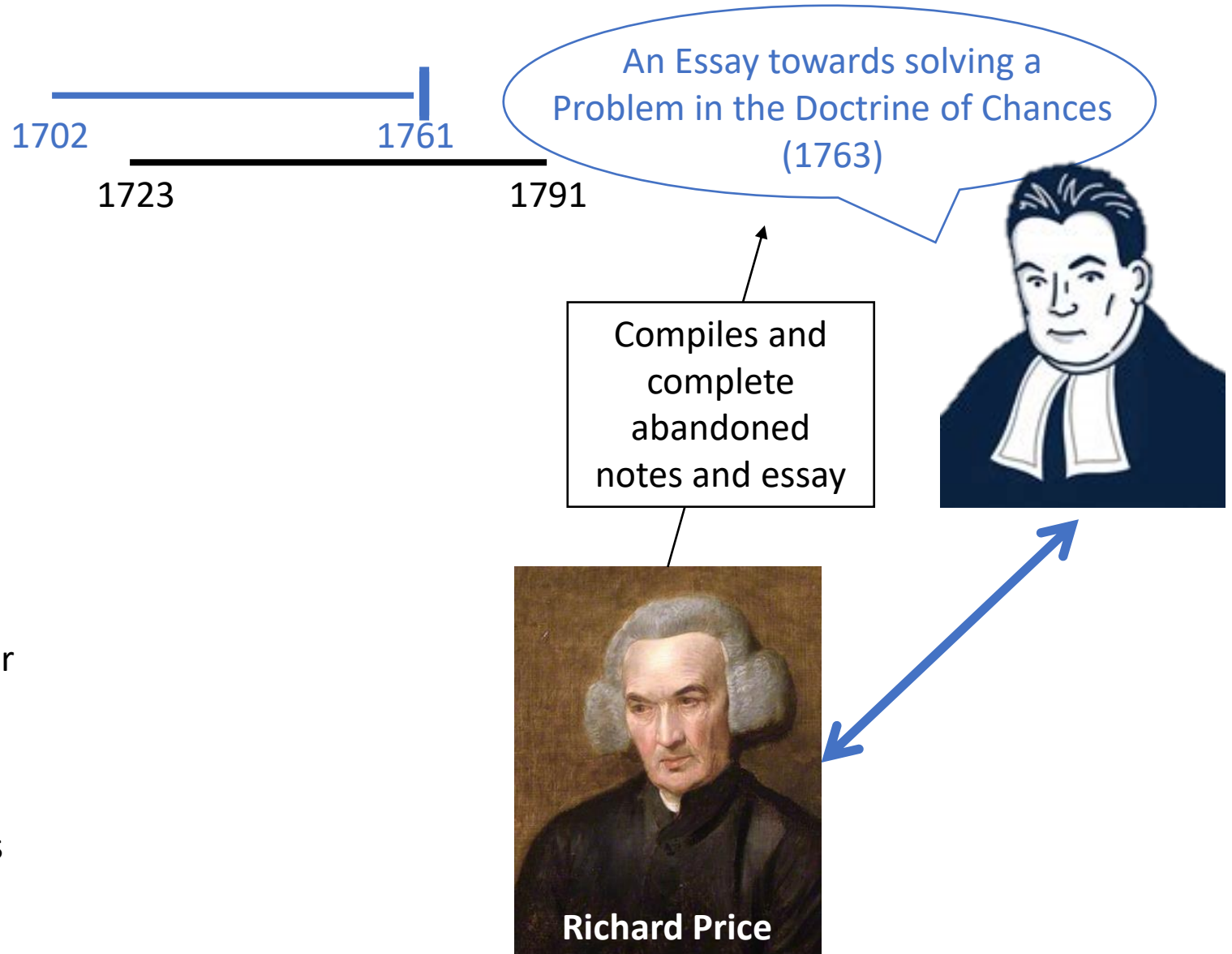
Concludes that
belief is at the
center of
rationalism, not
reason

Concludes that reports
of miracles change
nothing regarding our
understanding of
human existence

→ Clergy opposes

Thomas Bayes' theorem might be seen as a rebuttal of David Hume's work

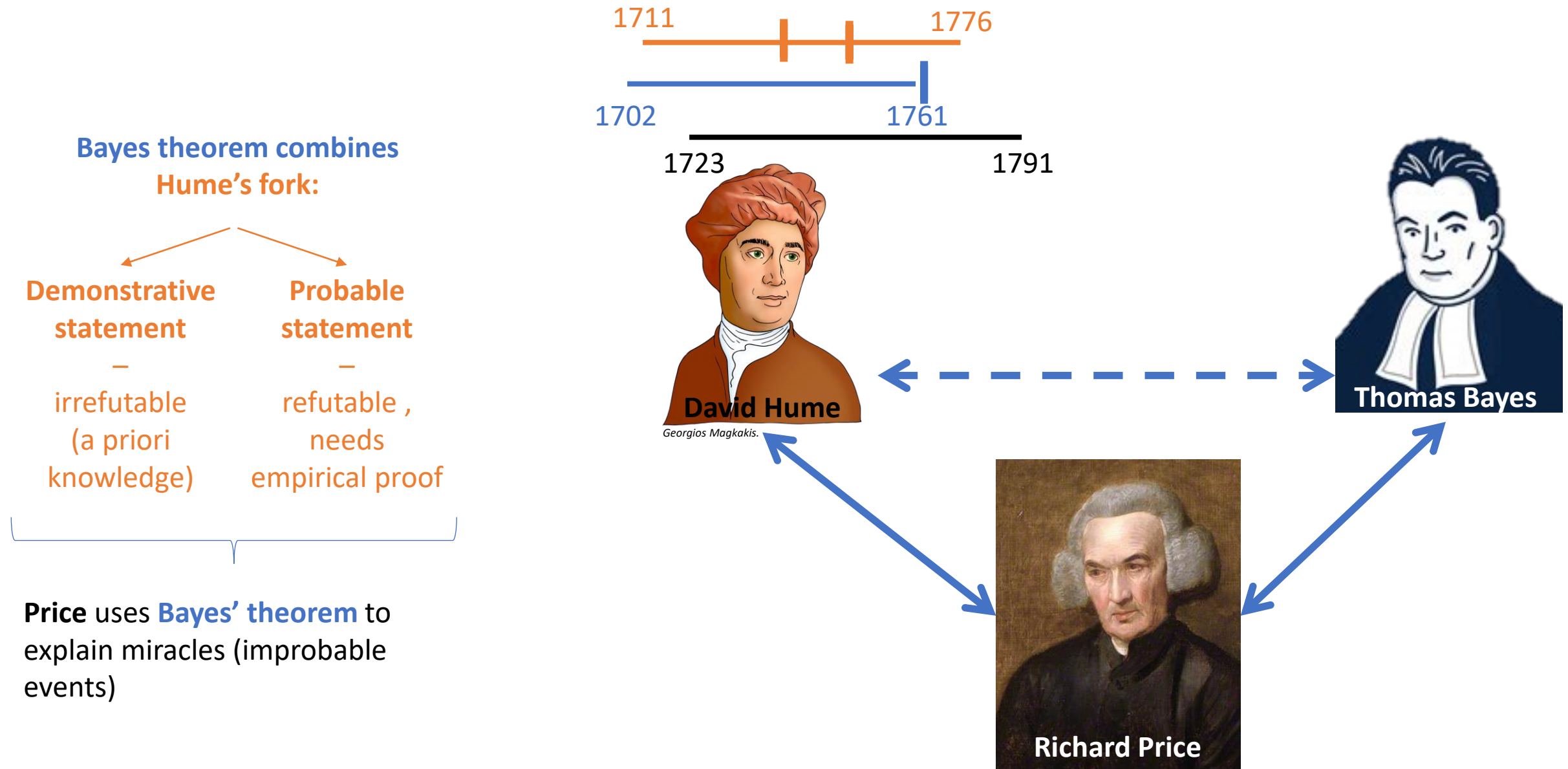
Thomas Bayes :
Highly skilled amateur
mathematician &
clergyman.
Started working on his
theorem shortly after
Hume's conclusion on
miracles



Richard Price :
Clergyman and pioneer
insurance statistician.

Completed and
published Bayes works
after his death

Thomas Bayes' theorem might be seen as a rebuttal of David Hume's work



Did you say probable(-ility)?

- **Conditional probability:** If y and x are two random variables

conditional probability = $p(y|x)$;

the probability of y for a given value of x (or **the probability of y knowing x**).

- **Joint probability:** probability of obtaining both a certain value of x and a certain value of y , noted $p(x, y)$, can be calculated in two ways:

$$p(x, y) = p(x)p(y|x) = p(y)p(x|y)$$

Diagram illustrating the calculation of joint probability $p(x, y)$ as the product of marginal and conditional probabilities in two ways:

- From $p(x)p(y|x)$:
 - $p(x)$: probability of getting x *
 - $p(y|x)$: probability of getting y knowing x
- From $p(y)p(x|y)$:
 - $p(y)$: probability of getting y *
 - $p(x|y)$: probability of getting x knowing y

- **Marginal probability:** probability of a variable y , $p(y)$, is its probability if we ignore the value of the other variables. If we do not know $p(y)$ directly, but we know $p(y, x)$ for each possible value of another variable x , then $p(y)$ corresponds to the sum of the joint probabilities of x and y for each value of x .

$$p(y) = \sum_x p(y, x) = \sum_x p(y|x)p(x)$$

Bayes' theorem



**Prior
distribution**

**Likelihood
function**

**Frequentist
approach !**

**Posterior
distribution**

$$p(x|y) = \frac{p(x)p(y|x)}{p(y)}$$

**Probability
of data**

Marginal probability

constant often omitted as it does not change the outcomes

Bayes' theorem



What we know
A Priori

The model

Parameters
estimates

$$p(x|Data) = \frac{p(x)p(Data|x)}{p(Data)}$$

Probability
of data

Bayes' theorem



Joint probability can become Bayes theorem by splitting the two right parts by $p(y)$;

$$p(x, y) = p(x)p(y|x) = p(y)p(x|y)$$

becomes

$$p(x|y) = \frac{p(x)p(y|x)}{p(y)}$$

Posterior distribution

Prior distribution

Likelihood function (model)

Probability of data

We can then calculate the probability distribution of x conditional of y if we know:

- the probability distribution of y conditional of x , and;
- the marginal probability distribution of x .

As for the denominator $p(y)$, this can be obtained by taking the sum (or the integral) of $p(x)p(y|x)$ over the set of possible values of x .

So, how do we solve Bayes' theorem?

Bayes theorem: can be hand solved for simple cases, but not for most of modern day use. --> Need a powerful computer.

Instead of resolving complex integrals, it is estimated by simulating many, many, many posterior distributions

Widely used algorithms:

WinBUGS, BUGS, JAGS, STAN, Nimble



Frequentist

Parameter is a known value for a given population. The sample serves to estimate this parameter.

Sample is one of many sample possibles from a population.

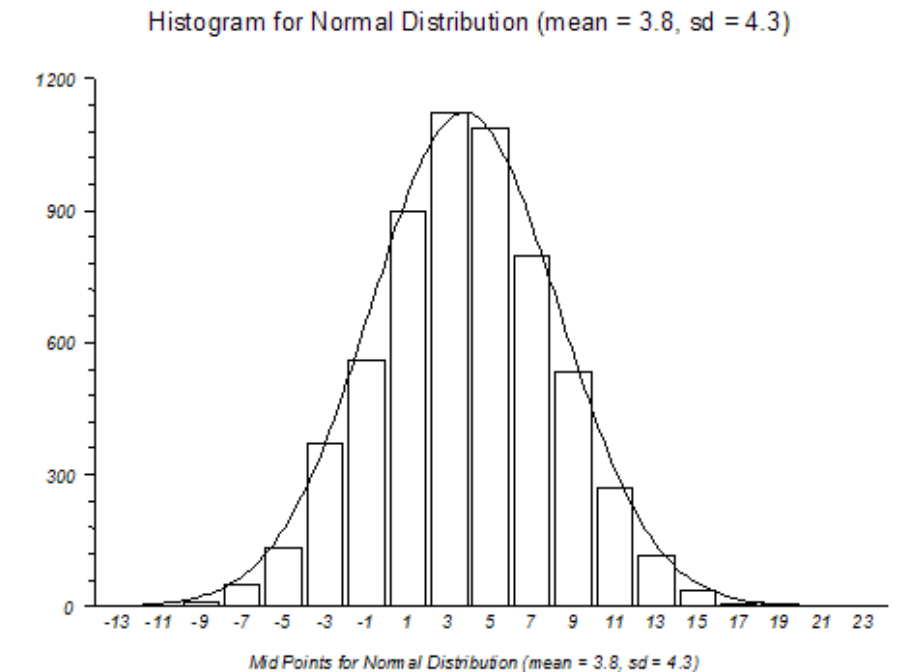
95% confidence intervals: 95% of possible samples of x would produce an interval containing the said value.

Bayesian

Parameter is a random value drawn from a distribution

Sample represents truth (what we actually know)

95% confidence intervals: Probability to observe the parameter within the interval's boundary. → credible interval.



Stan... who?

INSTALLATION DOCUMENTATION COMMUNITY ABOUT US YOUR SUPPORT SEARCH



RStan

the R interface to Stan

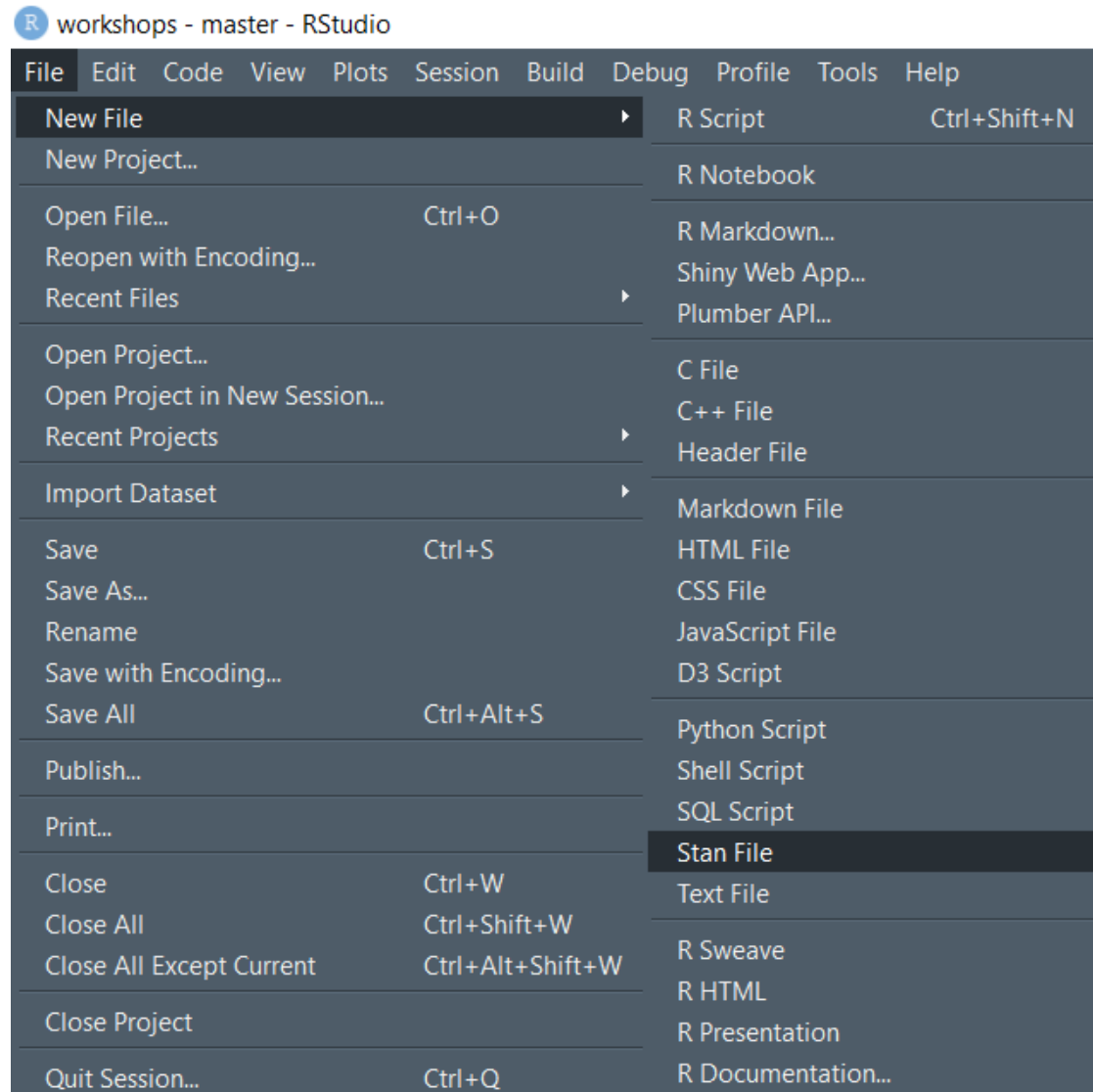
Download and Get Started

Instructions for downloading, installing, and getting started with RStan on all platforms.

- [RStan Quick Start Guide \(GitHub\)](#)

<https://mc-stan.org/users/interfaces/rstan>

To open a new Stan file



Empty Stan file

```
1 // ←
2 // This Stan program defines a simple model, with a
3 // vector of values 'y' modeled as normally distributed
4 // with mean 'mu' and standard deviation 'sigma'.
5 //
6 // Learn more about model development with Stan at:
7 //
8 //   http://mc-stan.org/users/interfaces/rstan.html
9 //   https://github.com/stan-dev/rstan/wiki/RStan-Getting-Started
10 //
11
12 // The input data is a vector 'y' of length 'N'.
13 data {
14   int<lower=0> N;
15   vector[N] y;
16 }
17
18 // The parameters accepted by the model. Our model
19 // accepts two parameters 'mu' and 'sigma'.
20 parameters {
21   real mu;
22   real<lower=0> sigma;
23 }
24
25 // The model to be estimated. We model the output
26 // 'y' to be normally distributed with mean 'mu'
27 // and standard deviation 'sigma'.
28 model {
29   y ~ normal(mu, sigma);
30 }
31
32 |
```

Use two forward slashes instead of a pound (#) when you want to write text and not code

Where to click if you have trouble coding in Stan

Define your data

Define your parameters

Define your model

Our data : BTdata.txt

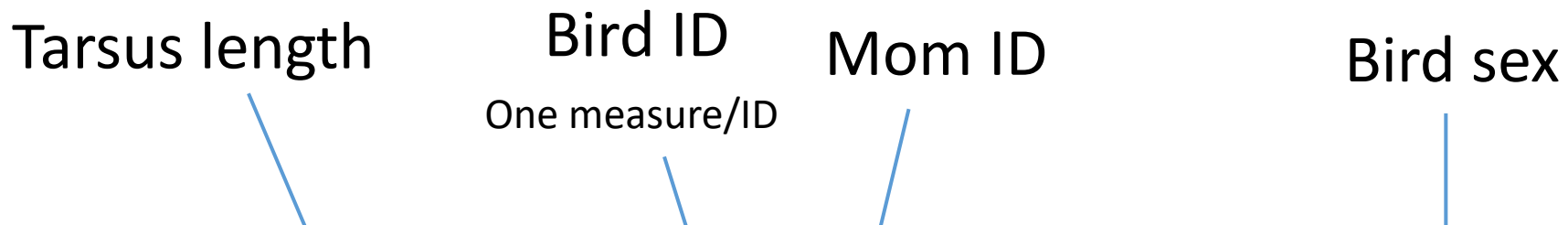
Tarsus length

Bird ID

One measure/ID

Mom ID

Bird sex




	tarsus	back	animal	dam	fosternest	hatchdate	sex
1	-1.89229718	1.146421150	R187142	R187557	F2102	-0.68740208	Fem
2	1.13610981	-0.759652092	R187154	R187559	F1902	-0.68740208	Male
3	0.98468946	0.144937259	R187341	R187568	A602	-0.42798142	Male
4	0.37900806	0.255584701	R046169	R187518	A1302	-1.46566405	Male
5	-0.07525299	-0.300699223	R046161	R187528	A2602	-1.46566405	Fem
6	-1.13519543	1.557721860	R187409	R187945	C2302	0.35028055	Fem
7	-1.13519543	-0.426824377	R187507	Fem3	C1902	-0.42798142	Male
8	1.89321156	-1.340762577	R187028	R187030	C1302	-0.94682274	Fem
9	-0.37809369	0.067481541	R046128	R187517	C602	-1.98450537	Fem

Stan hates factor... so make it integer!

```
15 # Convert factor to integer ####
16 library(dplyr)
17
18 BTdata$sex_no <- as.integer((BTdata$sex))
19
20 BTdata$dam_no <- as.integer((BTdata$dam))
21 dam_ID <- distinct(BTdata, dam_no, dam) %>%
22   arrange(dam_no)
23
```

Creates a DF
indicating which
dam_no corresponds
to which dam ID



	tarsus	back	animal	dam	fosternest	hatchdate	sex	sex_no	dam_no
1	-1.89229718	1.146421150	R187142	R187557	F2102	-0.68740208	Fem	1	56
2	1.13610981	-0.759652092	R187154	R187559	F1902	-0.68740208	Male	2	57
3	0.98468946	0.144937259	R187341	R187568	A602	-0.42798142	Male	2	61
4	0.37900806	0.255584701	R046169	R187518	A1302	-1.46566405	Male	2	38
5	-0.07525299	-0.300699223	R046161	R187528	A2602	-1.46566405	Fem	1	43
6	-1.13519543	1.557721860	R187409	R187945	C2302	0.35028055	Fem	1	94
7	-1.13519543	-0.426824377	R187507	Fem3	C1902	-0.42798142	Male	2	3
8	1.89321156	-1.340762577	R187028	R187030	C1302	-0.94682274	Fem	1	23
9	-0.37809369	0.067481541	R046128	R187517	C602	-1.98450537	Fem	1	37

Our model

$$\mu_{ij} = \beta_{0j} + \beta_i x_{ij} + \varepsilon_{0ij}$$

Tarsus length

Intercept

Sex

Error

$i = \text{observation number}$

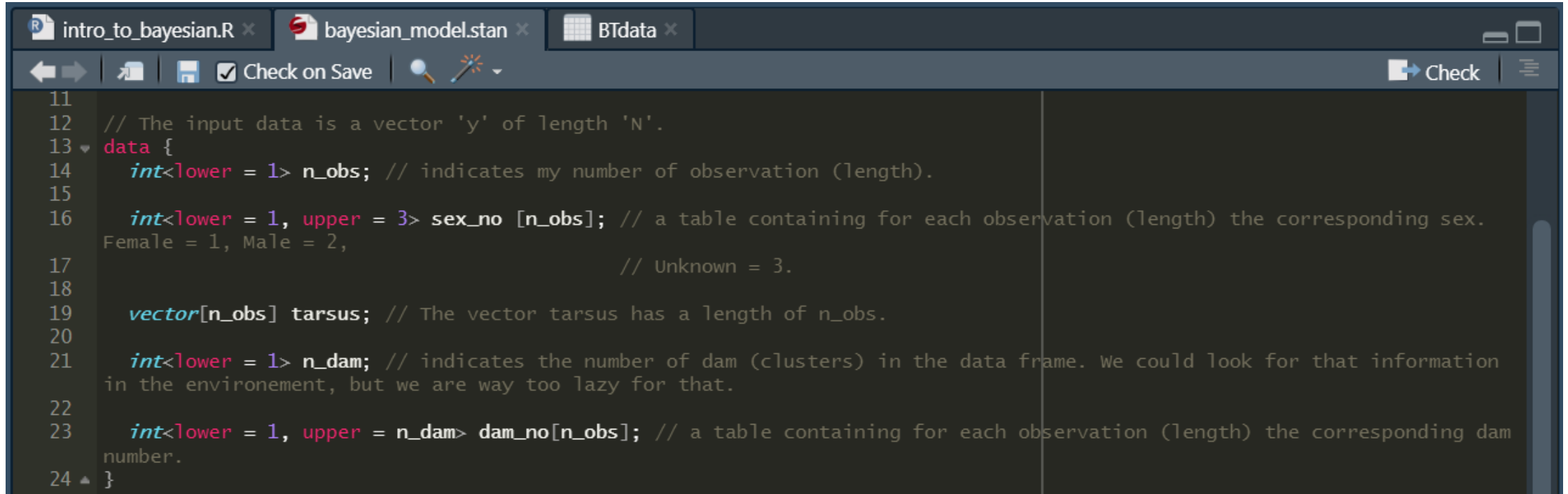
$j = \text{group (mother ID)}$

$\beta_{0j} = \beta_0 + \mu_{0j}$

$\mu_{0j} \sim N(0, \sigma_{\mu_0}^2)$

$\varepsilon_{0ij} \sim N(0, \sigma_{e_0}^2)$

Coding data in Stan

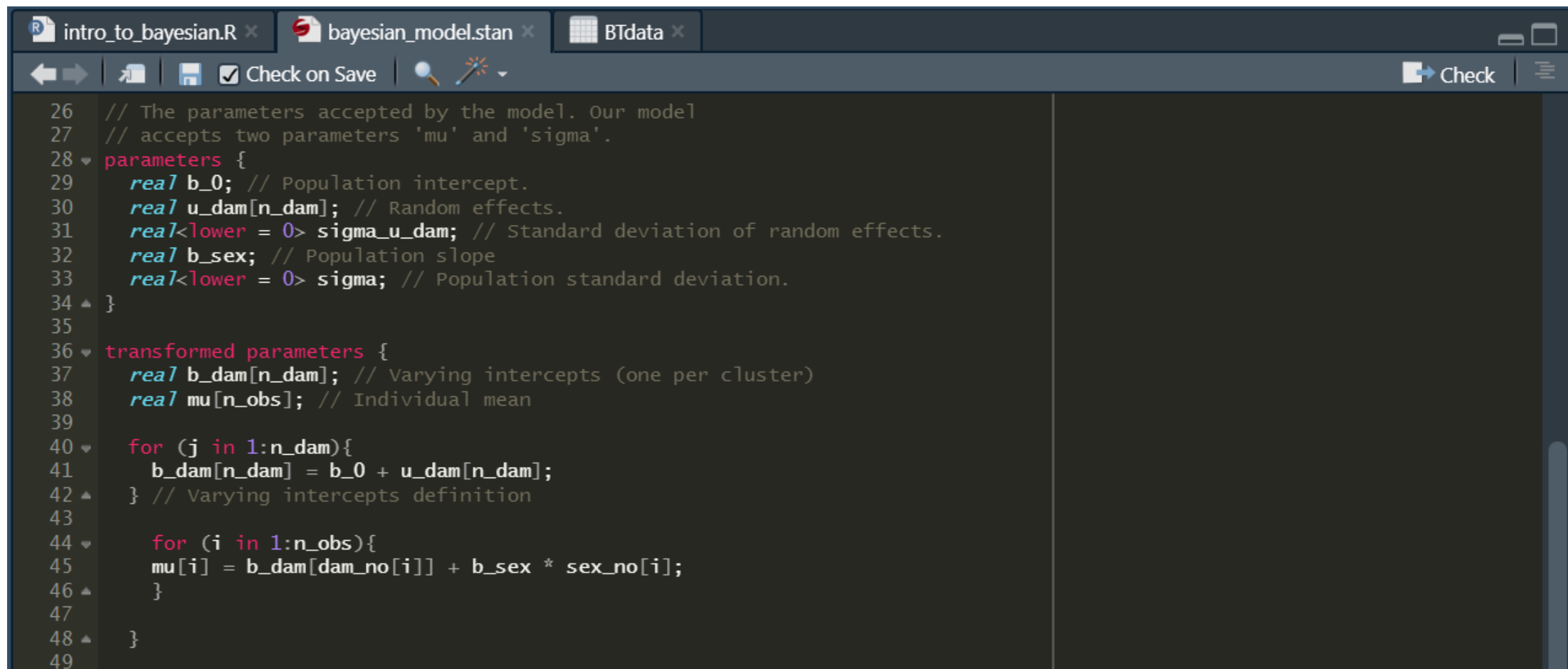


```
11
12 // The input data is a vector 'y' of length 'N'.
13 data {
14   int<lower = 1> n_obs; // indicates my number of observation (length).
15
16   int<lower = 1, upper = 3> sex_no [n_obs]; // a table containing for each observation (length) the corresponding sex.
17   // Female = 1, Male = 2,
18   // Unknown = 3.
19
20   vector[n_obs] tarsus; // The vector tarsus has a length of n_obs.
21
22   int<lower = 1> n_dam; // indicates the number of dam (clusters) in the data frame. We could look for that information
23   // in the environment, but we are way too lazy for that.
24   int<lower = 1, upper = n_dam> dam_no[n_obs]; // a table containing for each observation (length) the corresponding dam
25   // number.
26 }
```

Use chevrons to
indicate boundaries!

Use brackets to
indicate the length of
your data!

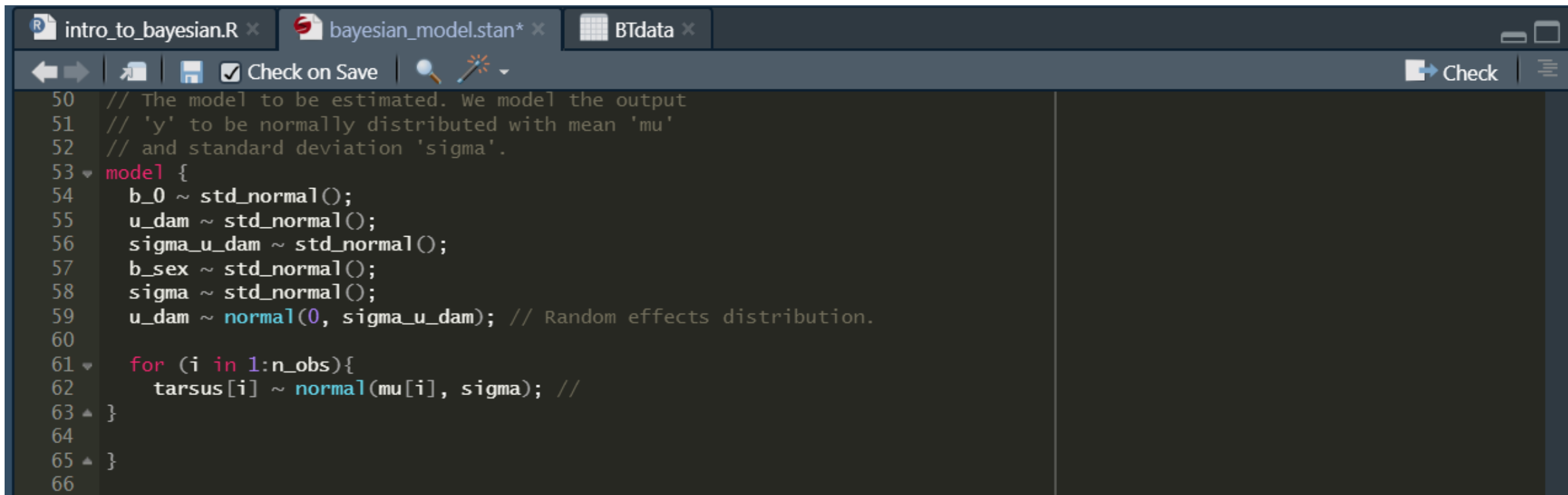
Coding parameters in Stan



The screenshot shows an RStudio editor window with three tabs: 'intro_to_bayesian.R', 'bayesian_model.stan', and 'BTdata'. The 'bayesian_model.stan' tab is active, displaying Stan code. The code defines parameters, transformed parameters, and a likelihood function. The parameters section includes a population intercept, random effects, and standard deviations. The transformed parameters section defines varying intercepts and the individual mean. The likelihood function is defined in the 'model' block.

```
26 // The parameters accepted by the model. Our model
27 // accepts two parameters 'mu' and 'sigma'.
28 parameters {
29   real b_0; // Population intercept.
30   real u_dam[n_dam]; // Random effects.
31   real<lower = 0> sigma_u_dam; // Standard deviation of random effects.
32   real b_sex; // Population slope
33   real<lower = 0> sigma; // Population standard deviation.
34 }
35
36 transformed parameters {
37   real b_dam[n_dam]; // Varying intercepts (one per cluster)
38   real mu[n_obs]; // Individual mean
39
40   for (j in 1:n_dam){
41     b_dam[j] = b_0 + u_dam[j];
42   } // Varying intercepts definition
43
44   for (i in 1:n_obs){
45     mu[i] = b_dam[dam_no[i]] + b_sex * sex_no[i];
46   }
47
48 }
49
```

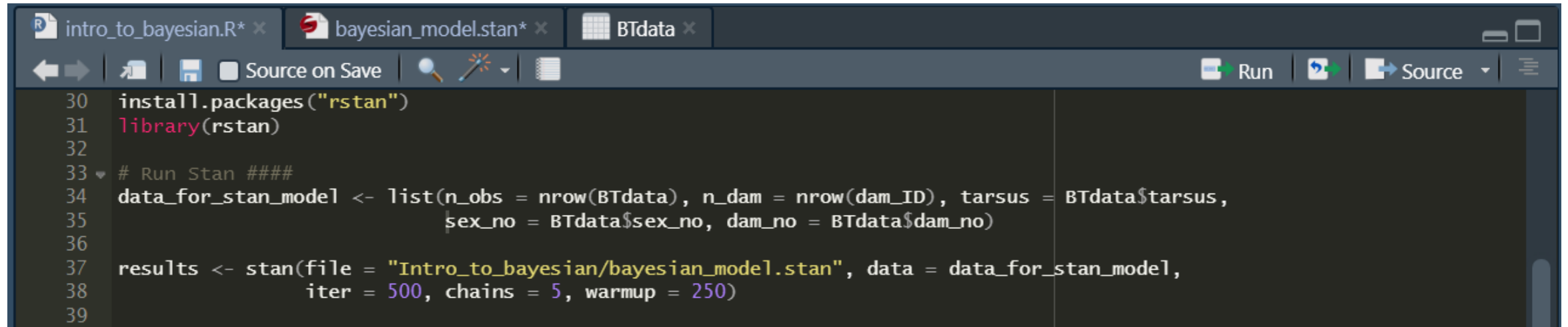
Coding models in Stan



The screenshot shows the RStudio interface with three tabs: 'intro_to_bayesian.R', 'bayesian_model.stan*', and 'BTdata'. The 'bayesian_model.stan' tab is active, displaying a Stan model. The code defines a model with parameters `b_0`, `u_dam`, `sigma_u_dam`, `b_sex`, and `sigma`, all drawn from standard normal distributions. A random effects distribution is defined for `u_dam` as `normal(0, sigma_u_dam)`. The data is modeled as `tarsus[i] ~ normal(mu[i], sigma)` for each observation `i` from 1 to `n_obs`. The interface includes a toolbar with icons for navigation, saving, and checking, and a 'Check' button on the right.

```
50 // The model to be estimated. We model the output
51 // 'y' to be normally distributed with mean 'mu'
52 // and standard deviation 'sigma'.
53 model {
54   b_0 ~ std_normal();
55   u_dam ~ std_normal();
56   sigma_u_dam ~ std_normal();
57   b_sex ~ std_normal();
58   sigma ~ std_normal();
59   u_dam ~ normal(0, sigma_u_dam); // Random effects distribution.
60
61   for (i in 1:n_obs){
62     tarsus[i] ~ normal(mu[i], sigma); //
63   }
64
65 }
66
```

Run your Stan model in R



The screenshot shows an RStudio interface with three tabs: 'intro_to_bayesian.R*', 'bayesian_model.stan*', and 'BTdata'. The 'intro_to_bayesian.R*' tab is active, displaying the following R code:

```
30 install.packages("rstan")
31 library(rstan)
32
33 # Run Stan ####
34 data_for_stan_model <- list(n_obs = nrow(BTdata), n_dam = nrow(dam_ID), tarsus = BTdata$tarsus,
35                             sex_no = BTdata$sex_no, dam_no = BTdata$dam_no)
36
37 results <- stan(file = "Intro_to_bayesian/bayesian_model.stan", data = data_for_stan_model,
38                 iter = 500, chains = 5, warmup = 250)
39
```

The code installs the 'rstan' package, loads it, and then runs a Stan model. The data for the model is extracted from the 'BTdata' object. The model file is 'bayesian_model.stan' located in the 'Intro_to_bayesian' directory. The model is run with 500 iterations, 5 chains, and 250 warmup iterations.