

A Beginner's Guide to Stat-JR's TREE interface version 1.0.3

Programming and Documentation by

William J. Browne, Christopher M.J. Charlton, Danius T.
Michaelides*, Richard M.A. Parker, Bruce Cameron, Camille
Szmaragd, Huanjia Yang*, Zhengzheng Zhang, Harvey
Goldstein, Kelvyn Jones, George Leckie and Luc Moreau*

Centre for Multilevel Modelling,

University of Bristol.

*Electronics and Computer Science,

University of Southampton.

August 2015

A Beginner's Guide to Stat-JR's TREE interface version 1.0.3

© 2015. William J. Browne, Christopher M.J. Charlton, Danius T. Michaelides, Richard M.A. Parker, Bruce Cameron, Camille Szmaragd, Huanjia Yang, Zhengzheng Zhang, Harvey Goldstein, Kelvyn Jones, George Leckie and Luc Moreau.

No part of this document may be reproduced or transmitted in any form or by any means, electronic or mechanical, including photocopying, for any purpose other than the owner's personal use, without the prior written permission of one of the copyright holders.

ISBN: To be confirmed

Printed in the United Kingdom

1.	About Stat-JR.....	2
1.1	Stat-JR: software for scaling statistical heights.....	2
1.2	About the Beginner’s guide	4
2	Installing and Starting Stat-JR	4
2.1	Installing Stat-JR.....	4
2.2	The use of third party software and licenses.....	4
2.3	Starting up <i>TREE</i>	5
3	Quick-start guide.....	6
3.1	Stage 1: Selecting a template & dataset.....	7
3.2	Stage 2: Providing template-specific input	9
3.3	Stage 3: Outputting the files to run the desired execution	11
3.4	Stage 4: Running the execution	12
3.5	Stage 5: The results	13
4	A detailed guide with worked examples.....	14
4.1	The structure and layout of the <i>TREE</i> interface.....	14
4.2	Application 1: Analysis of the tutorial dataset using the eStat engine.....	24
4.2.1	Summarising the dataset and graphs	24
4.2.2	Single-level Regression.....	29
4.2.3	Multiple chains	35
4.2.4	Adding gender to the model	36
4.2.5	Including school effects.....	38
4.2.6	Caterpillar plot	40
4.3	Interoperability – a brief introduction	42
4.3.1	So why are we offering Interoperability?	42
4.3.2	Regression in eStat revisited.....	42
4.3.3	Interoperability with WinBUGS.....	44
4.3.4	Interoperability with OpenBUGS	47
4.3.5	Interoperability with JAGS	49
4.3.6	Interoperability with MLwiN	51
4.3.7	Interoperability with R	53
4.3.8	Interoperability with AML.....	58
4.4	Application 2: Analysis of the Bangladeshi Fertility Survey dataset	59
4.4.1	The Bangladeshi Fertility Survey dataset	59
4.4.2	Modelling the data using logistic regression	60

4.4.3	Multilevel modelling of the data.....	63
4.4.4	Comparison between software packages.....	67
4.4.5	Orthogonal parameterisation.	71
4.4.6	Predictions from the model	73
4.5	Miscellaneous other topics e.g. Data Input/Export	75
5	References	76
6	Appendix: List of Third Party Software that are used by Stat-JR	78

Acknowledgements

The Stat-JR software is very much a team effort and is the result of work funded initially under three ESRC grants: the LEMMA 2 and LEMMA 3 programme nodes (Grant: RES-576-25-0003 & Grant: RES-576-25-0032) as part of the National Centre for Research Methods programme, and the e-STAT node (Grant: RES-149-25-1084) as part of the Digital Social Research programme. The work has continued with the ESRC grant ES/K007246/1.

We are therefore grateful to the ESRC for financial support to allow us to produce this software.

All nodes have many staff that, for brevity, we have not included in the list on the cover. We acknowledge therefore the contributions of:

Fiona Steele, Rebecca Pillinger, Paul Clarke, Mark Lyons-Amos, Liz Washbrook, Sophie Pollard, Robert French, Nikki Hicks, Mary Takahama and Hilary Browne from the LEMMA nodes at the Centre for Multilevel Modelling.

David De Roure, Tao Guan, Alex Fraser, Toni Price, Mac McDonald, Ian Plewis, Mark Tranmer, Pierre Walthery, Paul Lambert, Emma Housley, Kristina Lupton and Antonina Timofejeva from the e-STAT node.

A final acknowledgement to Jon Rasbash who was instrumental in the concept and initial work of this project. We miss you and hope that the finished product is worthy of your initials.

WJB August 2015.

1. About Stat-JR

1.1 Stat-JR: software for scaling statistical heights.

The use of statistical modelling by researchers in all disciplines is growing in prominence. There is an increase in the availability and complexity of data sources, and an increase in the sophistication of statistical methods that can be used. For the novice practitioner of statistical modelling it can seem like you are stuck at the bottom of a mountain, and current statistical software allows you to progress slowly up certain specific paths depending on the software used. Our aim in the Stat-JR package is to assist practitioners in making their initial steps up the mountain, but also to cater for more advanced practitioners who have already journeyed high up the path, but want to assist their novice colleagues in making their ascent as well.

One issue with complex statistical modelling is that using the latest techniques can involve having to learn new pieces of software. This is a little like taking a particular path up a mountain with one piece of software, spotting a nearby area of interest on the mountainside (e.g. a different type of statistical model), and then having to descend again and take another path, with another piece of software, all the way up again to eventually get there, when ideally you'd just jump across! In Stat-JR we aim to circumvent this problem via our interoperability features so that the same user interface can sit on top of several software packages thus removing the need to learn multiple packages. To aid understanding, the interface will allow the curious user to look at the syntax files for each package to learn directly how each package fits their specific problem.

To complete the picture, the final group of users to be targeted by Stat-JR are the statistical algorithm writers. These individuals are experts at creating new algorithms for fitting new models, or better algorithms for existing models, and can be viewed as sitting high on the peaks with limited links to the applied researchers who might benefit from their expertise. Stat-JR will build links by incorporating tools to allow this group to connect their algorithmic code to the interface through template-writing, and hence allow it to be exposed to practitioners. They can also share their code with other algorithm developers, and compare their algorithms with other algorithms for the same problem. A template is a pre-specified form that has to be completed for each task: some run models, others plot graphs, or provide summary statistics; we supply a number of commonly used templates and advanced users can use their own – see the Advanced User's Guide. It is the use of templates that allows a building block, modular approach to analysis and model specification.

At the outset it is worth stressing that there a number of other features of the software that should persuade you to adopt it, in addition to interoperability. The first is flexibility – it is possible to fit a very large and growing number of different types of model. Second, we have paid particular attention to speed of estimation and therefore in comparison tests, we have found that the package compares well with alternatives. Third it is possible to embed the software's templates inside an e-book which is exceedingly helpful for training and learning, and also for replication. Fourth, it provides a very powerful, yet easy to use environment for accessing state-of-the-art Markov Chain Monte Carlo procedures for calculating model estimates and functions of model estimates, via its eStat engine. The eStat engine is a newly-developed estimation engine with the advantage of being transparent in that all the algebra, and even the program code, is available for inspection.

While this is a beginner's guide – it is a beginner's guide to the software. We presume that you have a good understanding of statistical models which can be gained from for example the LEMMA online course (<http://www.bristol.ac.uk/cmm/learning/online-course/index.html>) . It also pre-supposes familiarity with MCMC estimation and Bayesian modelling – the early chapters of Browne (2012) available at <http://www.bristol.ac.uk/cmm/media/software/mlwin/downloads/manuals/2-33/mcmc-web.pdf> provide a practical introduction to this material.

Many of the ideas within the Stat-JR system were the brainchild of Jon Rasbash (hence the “JR” in Stat-JR). Sadly, Jon died suddenly just as we began developing the system, and so we dedicate this software to his memory. We hope that you enjoy using Stat-JR and are inspired to become part of the Stat-JR community: either through the creation of your own templates that can be shared with others, or simply by providing feedback on existing templates.

Happy Modelling,

The Stat-JR team.

1.2 About the Beginner's guide

We have written three initial guides to go with the software: this Beginner's Guide will cover how to start up and run the software, with a particular focus on the *TREE (Template Reading and Execution Environment)* interface. It will provide some simple examples and is designed for the researcher who wishes to be able to use the software package without worrying too much about how the mathematics behind the modelling works. As such, it does not go into detail on how users can contribute to extending the software themselves: that is covered in the second, Advanced User's, guide, designed for those who want to understand in greater detail how the system works. There is also a third, E-book User's, guide which deals with the software's alternative *DEEP (Documents with Embedded Execution and Provenance)* E-book interface.

As well as these three Guides, we also publish support, such as answers to frequently asked questions, on our website (<http://www.bristol.ac.uk/cmm/software/statjr>), where you can also find our forum in which users can discuss the software.

In this Beginner's Guide we first describe how to install Stat-JR, and then provide a 'Quick-start' guide as a quick visual overview, with brief notes, of the basics of how to work with Stat-JR via *TREE*. There then follows more detailed sections which provide further explanation, together with point-and-click examples for you to work through.

We look at an example application taken from education research, fitting a Normal response model for a continuous outcome. Here our aim is more to illustrate how to use the software than primarily how to do the best analysis of the dataset in question, and we will demonstrate the interoperability features with the other software packages that link to Stat-JR as well. We will then look at a second example from demography that illustrates binomial response models for a discrete outcome.

2 Installing and Starting Stat-JR

2.1 Installing Stat-JR

Stat-JR has a dedicated website (<http://www.bristol.ac.uk/cmm/software/statjr>) from which you can request a copy of the software, and which contains instructions for installation.

2.2 The use of third party software and licenses

Stat-JR is written primarily in the Python package but also makes use of many other third party software packages. We are grateful to the developers of these programs for allowing us to use their products within our package. When you have installed Stat-JR you will find a directory entitled licences in which you can find subdirectories for each package detailing the licensing agreement for each. The list of software packages that we are using can be found in the Appendix to this document.

2.3 Starting up *TREE*

Stat-JR's interface is viewed and operated via a web browser, but it is started by running an executable file.

To start Stat-JR select the *Stat-JR TREE* link from the *Centre for Multilevel Modelling* suite on the start up menu. This action opens a command prompt window in the background to which commands are printed out. This window is useful for viewing what the system is doing: for example, on the machine on which we have run *TREE*, you can see commands like the following:

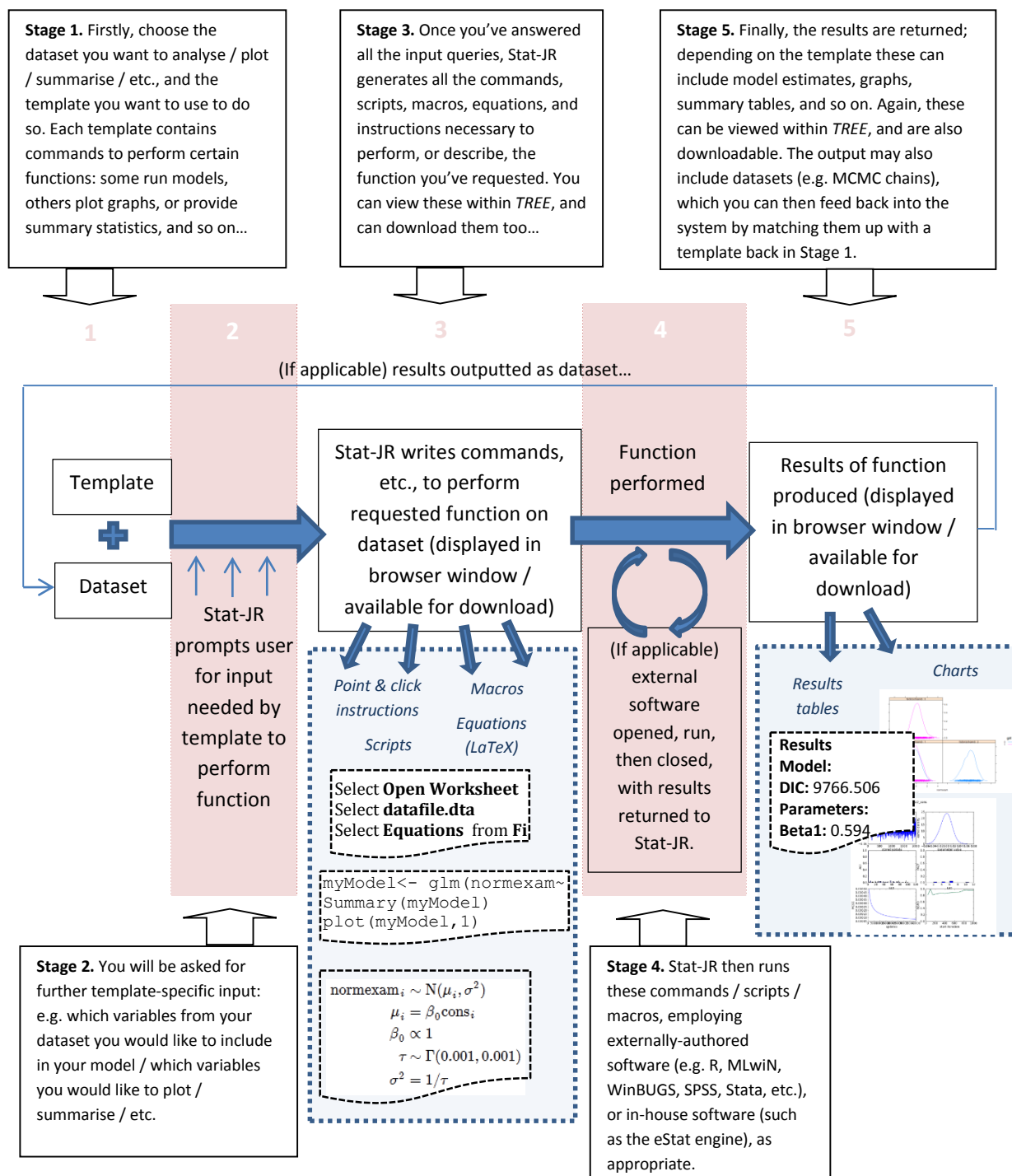
```
WARNING:root:Failed to load package GenStat_model (GenStat not found)
WARNING:root:Failed to load package Minitab_model (Minitab not found)
WARNING:root:Failed to load package Minitab_script (Minitab not found)
WARNING:root:Failed to load package SABRE (Sabre not found)
http://0.0.0.0:55534/
```

The most important command when starting up is the final line (the precise five-digit number written out towards the end of the line will likely differ, though). This only appears when the program has successfully performed all its initial set-up routines. This may take a while, particularly the first time you use the program. You should then be able to view the start page of *TREE* in your browser; if you can't, then try refreshing the browser window, or typing **localhost:55534** (in this example) into the address bar. The lines such as `WARNING:root:Failed to load package GenStat model (GenStat not found)` are not necessarily problematic but are warning you that the Genstat statistical package has not be found and loaded on your particular machine.

3 Quick-start guide

This section provides an overview ‘quick-start’ guide to using Stat-JR, via the *TREE* application; for more detailed instructions, together with worked point-and-click examples, see later sections. We’re assuming you’ve installed Stat-JR, and can see the opening page of the *TREE* application in your browser (see Section 2).

When operating Stat-JR through *TREE*, you generally proceed through the following five stages:



Below we briefly highlight the main features, with screenshots, of each of these five stages.

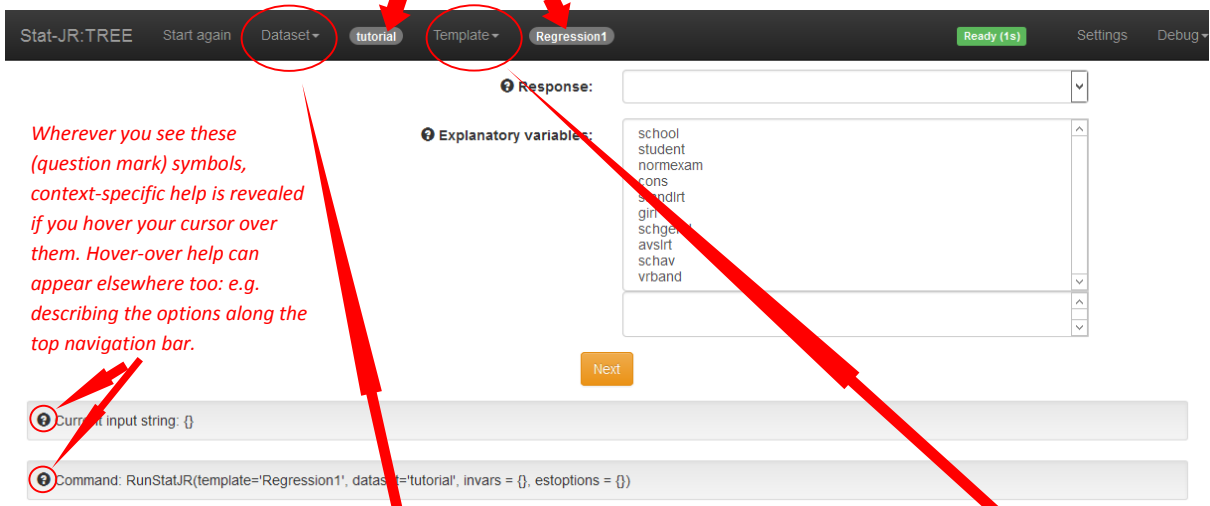
3.1 Stage 1: Selecting a template & dataset

- On opening Stat-JR, the page below, containing introductory information, will be displayed in a web browser. To proceed to choosing a template and dataset, click on the **Begin** button.

The screenshot shows the Stat-JR 1.0.3 welcome page. At the top left, the text 'Stat-JR: TREE' is displayed. At the top right, there are two menu items: 'Settings' and 'Debug'. A red arrow points from the 'Settings' menu to a text annotation: 'This link, available on all subsequent pages of TREE too, allows you to change settings such as paths to data and template folders, paths to interoperating software, and optimisation settings for generated code.' Another red arrow points from the 'Debug' menu to a text annotation: 'If you have modified any files in the templates, datasets or packages folders, then you can reload their contents into the current session via the Debug menu here, which is also available on all subsequent pages of TREE.' The main heading is 'Welcome to Stat-JR 1.0.3'. Below it, a paragraph of text mentions the developers and a link to 'webpages' which is circled in red. A red arrow points from this link to a text annotation: 'A link to the Stat-JR webpages which contain further support, including frequently asked questions & a user forum'. At the bottom left, there is a blue 'Begin' button circled in red. A red arrow points from this button to a text annotation: 'Click here to progress to the next screen where you can choose a template & dataset, and can start specifying your inputs...'

- Having pressed **Begin**, the page below will be displayed. Note that here, and on other screens, wherever you see the question mark symbols, **context-specific help** is revealed if you hover your cursor over them. Hover-over help can appear elsewhere too: e.g. describing the options along the top navigation bar.
- Here you can specify the **template** and **dataset** you want to use, and then begin to specify your **inputs**.
- Selecting **Dataset > Choose** or **Template > Choose** from the top bar will reveal lists of available datasets and templates. For each, find the one you want from the list, and then press **Use**.
- Note, when choosing a template, you can use the **cloud terms to help your search**: the blue tags describe functional aspects of the templates, whilst the red terms describe which engines / packages the templates support (you can combine search terms by clicking on more than one, and cancel your selections by pressing **[reset]**).

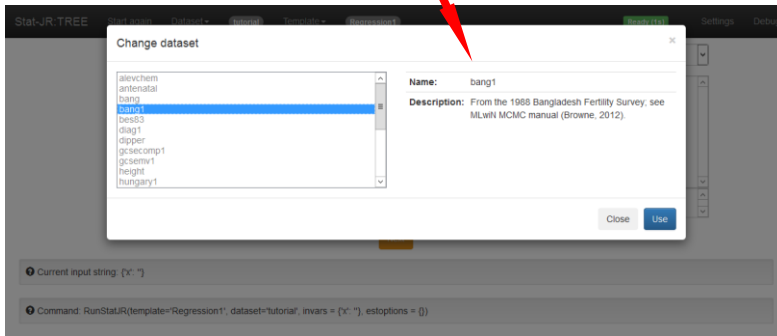
Here you can see which dataset and template are currently selected. Hovering your cursor over these names will reveal a description of each (if available).



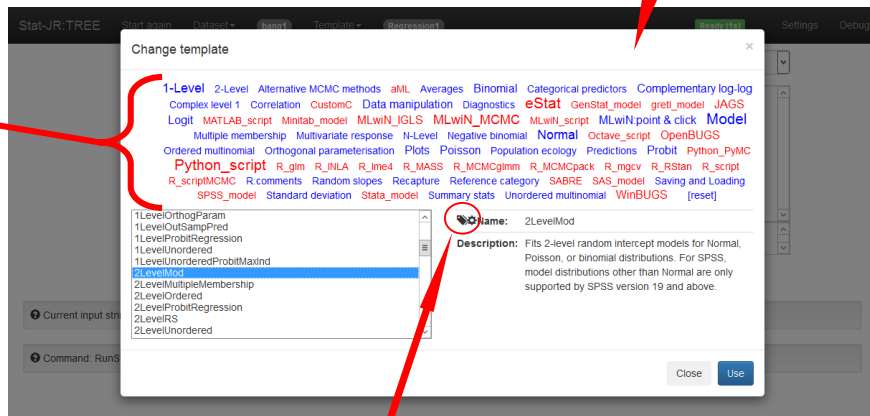
Wherever you see these (question mark) symbols, context-specific help is revealed if you hover your cursor over them. Hover-over help can appear elsewhere too: e.g. describing the options along the top navigation bar.

Clicking on the down arrow symbol just to the right of the **Dataset** heading in the top bar will bring up a menu. Select **Choose** to bring up the window, below, allowing you to nominate a dataset other than that currently selected...

...and likewise for the **Template**...



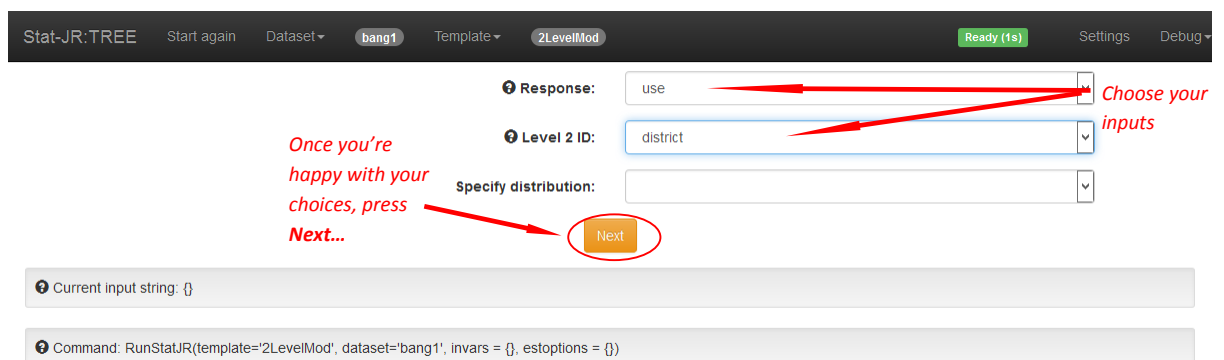
You can select one or more of these terms to help you find relevant templates; the blue tags describe the functional aspects of a template, whilst the red terms describe the engines / packages supported by a template. To unselect terms, press [reset]



Clicking the 'label' symbol brings up a list of tags, whilst clicking the 'cog' symbol brings up a list of supported engines / packages.

3.2 Stage 2: Providing template-specific input

- Once your desired **Dataset** and **Template** is selected, you can start answering the input questions back on the main page. These are required by Stat-JR to allow the template to perform the appropriate executions with your dataset; these inputs vary between templates, and also within templates too, depending on your earlier choices as you progress through the screens.
- For multi-choice lists you can de-select variables by simply clicking on their name in the list of selected items.
- Press **Next** each time you've completed the input questions on the current page.
- Then, if applicable, more inputs will be revealed, and those you have already selected will be greyed-out. However, you can still change each input via the **remove** button which you'll see next to each one. Alternatively, to re-specify *all* your inputs, press **Start again** (in the top bar).
- When asked for the **Name of the output results**, this will be the name given to any outputted dataset which results from running the template (see Stage 5).



Stat-JR: TREE Start again Dataset **bang1** Template **2LevelMod** Ready (1s) Settings Debug

If, at any point, you want to re-specify all your inputs, then press **Start again**

For multi-select lists, you can de-select variables by clicking on their name here

You can remove specific inputs via these buttons here

As you progress through the screens, you can see your choices reflected in the input string and the RunStatJR command, at the bottom; a record of your inputs is also kept under **Template > Set inputs** (via the black bar at the top), allowing you to automatically populate the inputs boxes with your previous choices (see later section); the RunStatJR command, on the other hand, can be used to call Stat-JR via a command line

Again, once you're happy with your inputs, press **Next**

Response: use remove

Level 2 ID: district remove

Specify distribution: Binomial remove

Denominator: cons

Specify link function: logit

Explanatory variables: woman, district, use, lc, urban, educ, hindu, d_illit, d_pray, cons, age

treat cons as categorical

treat age as categorical

Store level 2 residuals? Yes No

Next

Current input string: {'y': 'use', 'L2ID': 'district', 'D': 'Binomial'}

Command: RunStatJR(template='2LevelMod', dataset='bang1', invars = {'y': 'use', 'L2ID': 'district', 'D': 'Binomial'}, estoptions = {})

Stat-JR: TREE Start again Dataset **bang1** Template **2LevelMod** Ready (1s) Settings Debug

Denominator: cons remove

Specify link function: logit remove

Explanatory variables: cons.age remove

Store level 2 residuals? Yes remove

Choose estimation engine: eStat remove

Number of chains: 4 remove

Random Seed: 1 remove

Length of burnin: 1000 remove

Number of iterations: 2500 remove

Thinning: 1 remove

Use default algorithm settings: Yes remove

Generate prediction dataset: No remove

Use default starting values: Yes remove

Name of output results: my_output

Next

Current input string: {'D': 'Binomial', 'storeresid': 'Yes', 'nchains': '4', 'link': 'logit', 'defaultalg': 'Yes', 'iterations': '2500', 'seed': '1', 'defaultsv': 'Yes', 'Engine': 'eStat', 'L2ID': 'district', 'burnin': '1000', 'n': 'cons', 'thinning': '1', 'y': 'use', 'x': 'cons.age', 'makepred': 'No'}

Command: RunStatJR(template='2LevelMod', dataset='bang1', invars = {'L2ID': 'district', 'D': 'Binomial', 'storeresid': 'Yes', 'n': 'cons', 'link': 'logit', 'y': 'use', 'x': 'cons.age'}, estoptions = {'Engine': 'eStat', 'burnin': '1000', 'defaultsv': 'Yes', 'thinning': '1', 'nchains': '4', 'defaultalg': 'Yes', 'iterations': '2500', 'seed': '1', 'makepred': 'No'})

(We've skipped a screen or two where we were asked about this input – some have default values, and we've changed a few...)

We've now completed all the inputs, and so we press **Next** for the final time...

This is the name given to any outputted dataset (e.g. MCMC chains produced by the model run)

3.3 Stage 3: Outputting the files to run the desired execution

- Once you're pressed **Next** after the final input, Stat-JR returns a number of initial outputs which you can view in the output pane at the bottom of the window.
- Note that Stat-JR hasn't done everything you want it to do yet: it's just producing preliminary files telling you what it's going to do, and how it's going to do it.
- To select particular content to view in the output pane, use the drop-down menu just above it.
- The **Popout** button, just above the output pane, allows you to view its contents in a new browser tab.
- Pressing **Run** performs the executions described by the scripts, etc, returned in the output pane.

*Via the **Edit** button, you can directly edit scripts and macros, e.g. to change model specification, plot characteristics, etc...*

Generate prediction dataset: No [remove](#)

Use default starting values: Yes [remove](#)

Name of output results: my_output [remove](#) *Press **Run** to perform the executions...*

Run

Current input string: {'D': 'Binomial', 'storeresid': 'Yes', 'nchains': '4', 'link': 'logit', 'defaultalg': 'Yes', 'iterations': '2500', 'outdata': 'my_output', 'seed': '1', 'defaultsv': 'Yes', 'Engine': 'eStat', 'L2ID': 'district', 'burnin': '1000', 'n': 'cons', 'thinning': '1', 'y': 'use', 'x': 'cons.age', 'makepred': 'No'}

Command: RunStatJR(template='2LevelMod', dataset='bang1', invars = {'L2ID': 'district', 'D': 'Binomial', 'storeresid': 'Yes', 'n': 'cons', 'link': 'logit', 'y': 'use', 'x': 'cons.age'}, estoptions = {'Engine': 'eStat', 'burnin': '1000', 'defaultsv': 'Yes', 'thinning': '1', 'nchains': '4', 'defaultalg': 'Yes', 'iterations': '2500', 'outdata': 'my_output', 'seed': '1', 'makepred': 'No'})

Edit equation.tex **Popout** *Click here to view the contents of the output pane, below, in a new browser tab...*

You can choose what to view in the output panel (here we've chosen to view the equation for the model we've specified), via this selection box

$$\begin{aligned}
 use_i &\sim \text{Binomial}(cons_i, \pi_i) \\
 \text{logit}(\pi_i) &= \beta_0 cons_i + \beta_1 age_i + u_{\text{district}[i]} \\
 u_{\text{district}[i]} &\sim N(0, \sigma_u^2) \\
 \beta_0 &\propto 1 \\
 \beta_1 &\propto 1 \\
 \tau_u &\sim \Gamma(0.001, 0.001) \\
 \sigma_u^2 &= 1/\tau_u
 \end{aligned}$$

3.4 Stage 4: Running the execution

- Once you're pressed **Run**, the executions specified by you are performed.
- Depending on your choices, this may take anything from a second or two (e.g. to produce a simple plot, fit a model using a non-iterative method of estimation, produce summary data, etc.), to many minutes (e.g. to run MCMC chains for a large number of iterations).
- If appropriate (e.g. if the template supports inter-operability, and if you have chosen to employ it when prompted), externally-authored software packages (e.g. R, MLwiN, WinBUGS, SPSS, etc.) are opened, run, then closed, and the results are returned to Stat-JR.
- Whilst the execution runs, you may see a lot of activity in the black command window, which may help you keep a track of progress.

The screenshot shows the Stat-JR interface with the following settings:

- Length of burnin: 1000 remove
- Number of iterations: 2500 remove
- Thinning: 1 remove
- Use default algorithm settings: Yes remove
- Use default prediction dataset: No remove
- Use default starting values: Yes remove
- Name of output results: my_output remove
- Extra iterations: (empty)

Buttons: Download, Add to ebook

Current input string: ('D': 'Binomial', 'storerresid': 'Yes', 'nchains': '4', 'link': 'logit', 'defaultalg': 'Yes', 'iterations': '2500', 'outdata': 'my_output', 'seed': '1', 'defaultsv': 'Yes', 'Engine': 'eStat', 'L2ID': 'district', 'burnin': '1000', 'n': 'cons', 'thinning': '1', 'y': 'use', 'x': 'cons.age', 'makepred': 'No')

Command: RunStatJR(template=2LevelMod, dataset=bang1, invars = ['L2ID': 'district', 'D': 'Binomial', 'storerresid': 'Yes', 'n': 'cons', 'link': 'logit', 'y': 'use', 'x': 'cons.age'], estoptions = ('Engine': 'eStat', 'burnin': '1000', 'defaultsv': 'Yes', 'thinning': '1', 'nchains': '4', 'defaultalg': 'Yes', 'iterations': '2500', 'outdata': 'my_output', 'seed': '1', 'makepred': 'No'))

equation.tex Popout

$$use_i \sim \text{Binomial}(cons_i, \pi_i)$$

Whilst it performs these executions, the progress gauge indicates that Stat-JR is still working...

You may see a lot of activity in the command window as the execution is performed...

3.5 Stage 5: The results

- Once the executions have run, the progress gauge, in the top-right corner, will change from “Working” to “Ready”, and the drop-down list, just above the output pane, will now be populated with more results.
- Depending on the template, a range of buttons / boxes appear above the output pane allowing you to e.g. **Download** the results, **Add to ebook**, and run chains for **Extra iterations**.
- If applicable, an outputted dataset now appears in the list of datasets (see **Dataset > Choose**, via the top bar).

The screenshot shows the Stat-JR: TREE interface. At the top, the status is 'Ready (71s)'. Below the status bar, there are several settings sections: 'Use default algorithm settings: Yes remove', 'Generate prediction dataset: No remove', 'Use default starting values: Yes remove', and 'Name of output results: my_output remove'. There are buttons for 'Extra Iterations', 'Download', and 'Add to ebook'. A 'More' button is also visible. Below these buttons, there are two informational boxes: one for the 'Current input string' and another for the 'Command'. At the bottom, there is a dropdown menu showing 'ModelResults' and a 'Popout' button. A table of results is displayed below the dropdown menu.

The outputted dataset (which we earlier chose to call 'my_output') will now appear in the list of datasets (see Dataset > Choose) allowing us to investigate it further by matching it up with another template ...

Stat-JR indicates it has finished running these executions, by being "Ready" again...

Here you can add, to an eBook, the inputs you have just entered, the details of the template and dataset you have just chosen, and the outputs you would like to be displayed...

You can Download results, and run for Extra iterations ...

The results (e.g. plots, model estimates, etc.) are added to the list of outputs; here we've chosen to display a summary table of results...

parameter	mean	sd	ESS	variable
sigma2_u	0.269406474316	0.089387391158	850	
beta_0	-0.539822657478	0.0856270762827	507	cons
beta_1	0.00853908304127	0.00552257842808	2421	age
tau_u	4.14221258894	1.45517648317	692	
u_0	-0.4563146855	0.205573010541	2043	district
u_1	-0.0482519589551	0.348344553015	2358	district
u_2	0.299541719112	0.485614002114	2318	district

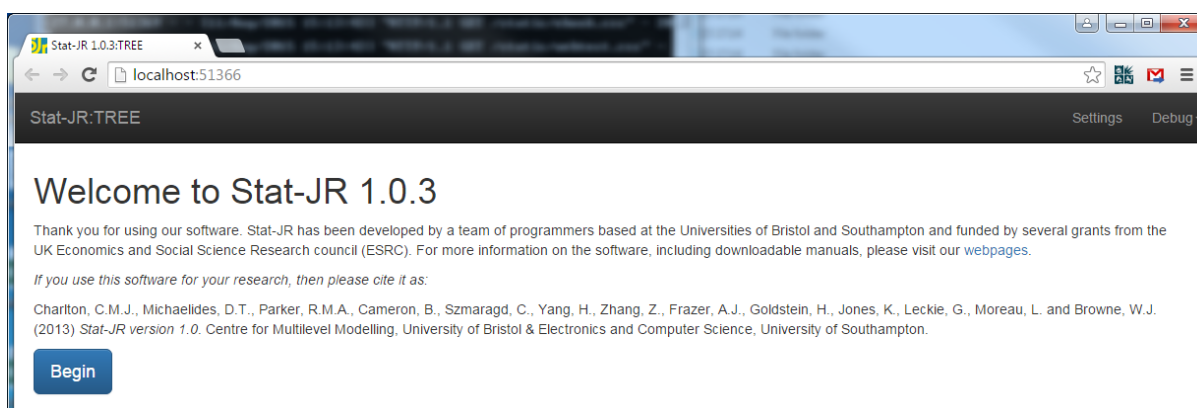
This ends the quick start guide. In the next chapter we describe the operation of TREE in more detail, and work through examples.

4 A detailed guide with worked examples

4.1 The structure and layout of the *TREE* interface

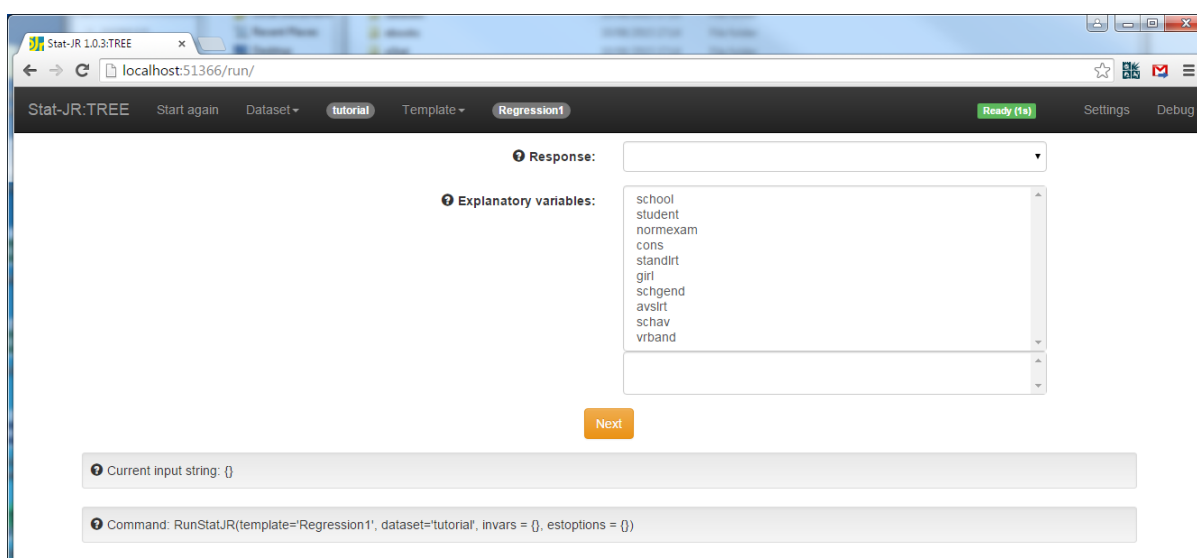
Stat-JR can be thought of as a system that manages the use of a set of templates written either by the developers, and supplied with the software, or by users themselves. Each template will perform a specific function: for example, fitting a specific family of models, summarising a dataset, or plotting a graph. The Stat-JR system therefore allows the user to select and use specific templates with their datasets, and to capture and display the outputs that result.

Returning to our start-up of the software, when the line `http://0.0.0.0:50215/` appears, and after refreshing the web browser, the browser window should appear as follows:



This is the start screen for the *TREE* interface to Stat-JR, and contains information on funders, authors, and a link to the Stat-JR website which contains further guidance, such as answers to frequently asked questions, and a user forum.

Pressing **Begin** returns the following screen:



At the top you'll see a black **title bar**. From left to right, this contains:

- a link (**Stat-JR:TREE**) back to the opening page;

- an option (**Start again**) to clear all inputs the user has chosen for the current template;
- a **Dataset** menu allowing the user to choose, drop (from temporary memory cache), return summary statistics for the current dataset, view (the entire dataset; see below), return a list of datasets, upload / download (see Section 4.5) datasets. For example, selecting **Dataset > Choose** returns a scrollable list of all the datasets that the system is aware of: i.e. those which appear in the *datasets* subdirectory of this installation of Stat-JR. This pane can be used to change the selected dataset via the **Use** button; for inputting your own data set you can use the **Upload** button.
- the name of the currently-selected dataset (in the grey box) – if you hover your cursor over this name, it returns a textual description of the dataset;
- a **Template** menu allowing the user to choose, list (described below), upload individual templates not already uploaded in the current session or set the template inputs as a list. If you select **Template > Choose**, a box appears which contains a scrollable list of all the templates that the system is aware of: i.e. those which appear in the *templates* subdirectory of this installation of Stat-JR. This can be used to change the selected template via the **Use** button. As we anticipate there being many templates, each template has defined ‘*tags*’ which are terms to describe what it does. These appear as blue phrases in the ‘cloud’ above the list of templates, whereas the estimation engines supported by each template appear in the cloud in red. When you select a template, its name and description appear to the right of the list. Clicking on the symbol that looks like a baggage label returns the tags for that template, whereas clicking on the ‘cog’ symbol returns a list of engines that particular template supports;
- the name of the currently-selected template (in the grey box) – again, if you hover your cursor over this name, it also returns a description of the template;
- a progress gauge indicating whether Stat-JR is “Ready”, “Working”, “Idle” or whether it has encountered an “Error”;
- a link to a page containing options to change a variety of **Settings** (discussed further below);
- a **Debug** button; this produces a drop-down list from which one can choose to reload the templates, datasets or packages, allowing users upload changes to files they make outside the *TREE* interface, without having to start-up *Stat-JR* again. For example, a user could paste a new dataset into the datasets directory, or modify a template in the templates directory, and reload them so that they appear in their lists in the browser window.

The **Settings** window, accessible via the black **title bar**, displays a number of settings that the program uses with each possible software package: some of these are relatively straightforward, such as where the executables for each package are found, and some are relatively advanced, such as for the eStat engine, optimisation, starting values and standalone code options.

We will now look at The **View dataset** window:

Select **Dataset > Choose** from the menu in the black **title bar**.

Scroll down the dataset list, towards the bottom, and click on **rats**; its name and description will appear to the right of the list.

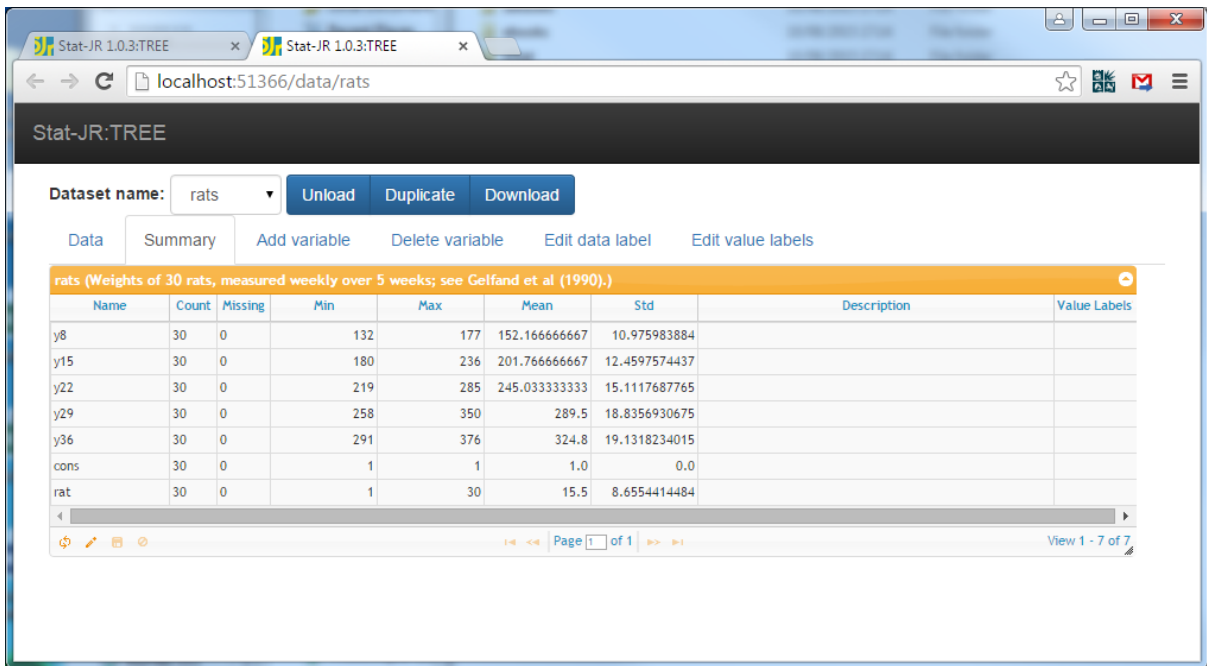
Click on the **Use** button, and the name of the current dataset (in the grey box in the black title bar at the top) should have changed accordingly.

Select **Dataset > View**; this will open a new tab in your browser: if you click on this you will be able to see the dataset we have just selected, as follows:

	y8	y15	y22	y29	y36	cons	rat
1	151	199	246	283	320	1	1
2	145	199	249	293	354	1	2
3	147	214	263	312	328	1	3
4	155	200	237	272	297	1	4
5	135	188	230	280	323	1	5
6	159	210	252	298	331	1	6
7	141	189	231	275	305	1	7
8	159	201	248	297	338	1	8
9	177	236	285	350	376	1	9
10	134	182	220	260	296	1	10
11	160	208	261	313	352	1	11
12	143	188	220	273	314	1	12
13	154	200	244	289	325	1	13
14	171	221	270	326	358	1	14
15	163	216	242	281	312	1	15
16	160	207	248	288	324	1	16
17	142	187	234	280	316	1	17
18	156	203	243	283	317	1	18
19	157	212	259	307	336	1	19
20	152	203	246	286	321	1	20
21	154	205	253	298	334	1	21

The **rats** dataset is a small, longitudinal animal growth dataset which contains the weights of 30 laboratory rats on 5 weekly occasions from 8 days of age (see Gelfand et al (1990) for more details). The five measurements are labelled **y8**, **y15**, **y22**, **y29** and **y36**, respectively, and the dataset also contains a constant column – a set of 1's, named **cons**, and a rat identifier column, **rat**. Initially, we are going to perform a regression analysis of the initial weight (**y8**) on the final weight (**y36**), including an intercept (**cons**). The boxes above the data allow the user to quickly add a new variable or delete an existing variable from the dataset. We can also view a summary of the dataset:

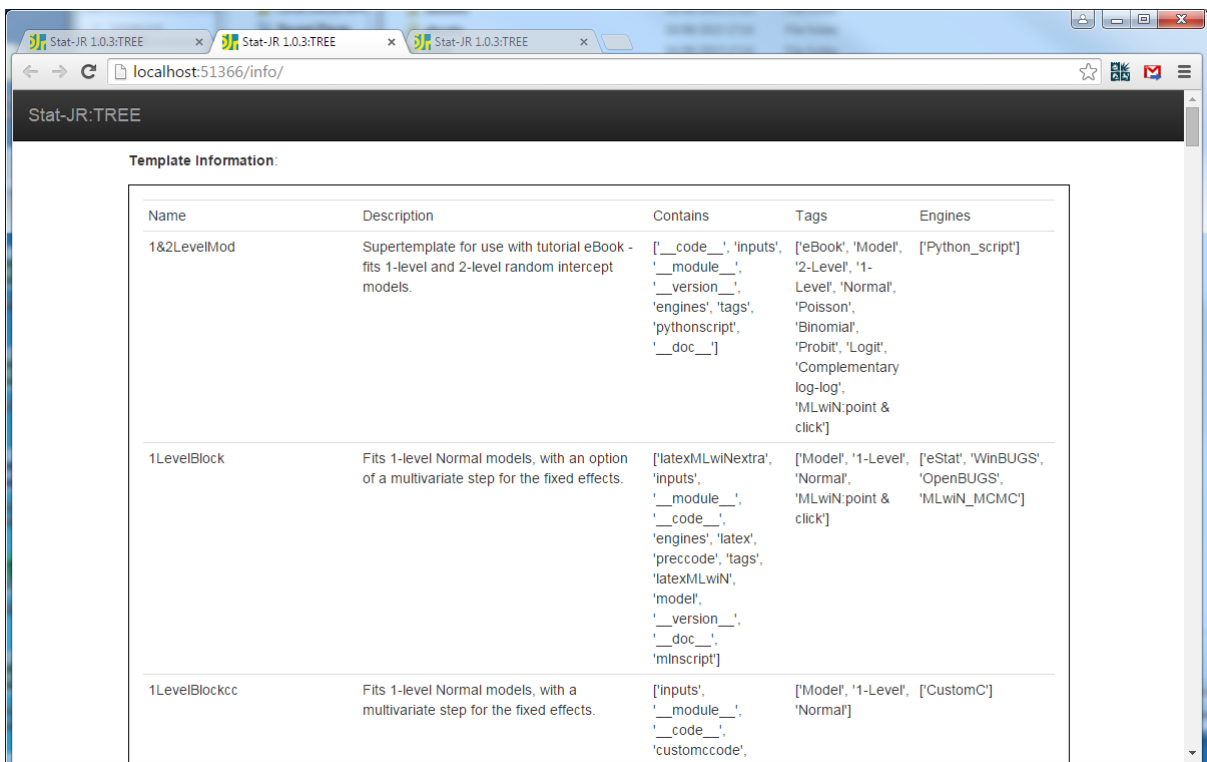
To view a summary of the dataset, click on the **Summary** tab in the tabs above the data and the screen will look as follows:



Here we get a very short summary of the dataset, giving, for each variable, the minimum value, maximum value, mean and standard deviation. If the dataset has had descriptions added or has categorical variables then they will appear in the last two columns. More extensive summaries are available by using specific templates to summarise datasets, as we will see later.

Let's now look at the **Template** menu:

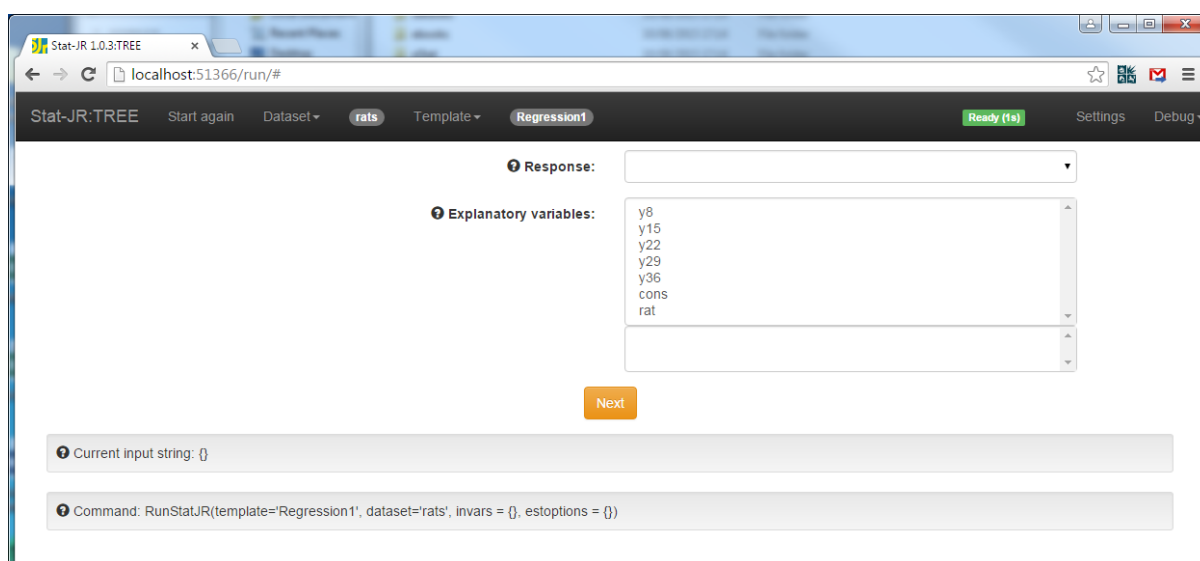
Back on the main page, if you click on **Template > List** the following screen will appear in a new tab:



This rather busy screen (we don't reproduce it all here, due to its length) contains, in the two columns on the left, a tabular list of all the templates that are available with a short description of what each template does. The other three columns are of more interest to advanced users, but contain a list of functions in the template code, tags that identify the template type, and the engines that are supported by the template.

We will next demonstrate running a template, using the default **Regression1** template that fits a 1-level Normal response regression model, this is appropriate as the response, the weights of the rats, is a continuous measure

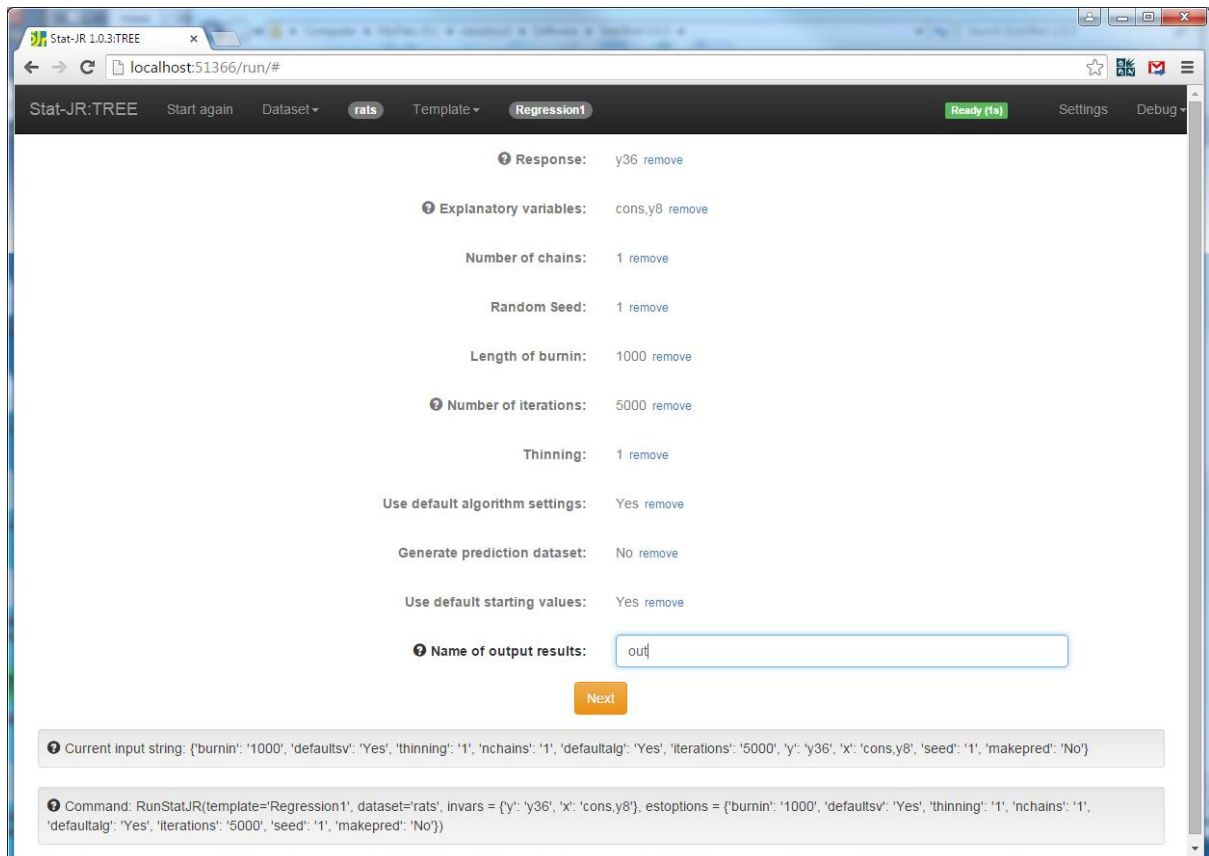
Return to the main menu screen, which should look as follows:



In the middle of the screen you can see the inputs required for this template (these are template-specific, and may change when you use a different template). Since some inputs are conditional (i.e. are only required when earlier inputs take specific values), the opportunity to specify inputs proceeds through sequential steps. Here we see the two initially-required inputs for the **Regression1** template are the **Response** variable and **Explanatory variables**. Since this template only allows for 1 response variable to be specified, a pull-down list is displayed, but since it allows for several explanatory variables to be specified, a multiple selection list is displayed for that input. In the case of the latter, variables are selected by clicking on their name in the left-hand list; to de-select them, click on their name in the right-hand list.

The **Start again** link (in the top black bar) will clear any inputs the user has already selected and return you to the first template input screen (i.e. the current screen, in this case), whilst the **Next** button will allow the user to move on and specify further inputs once those on the current screen have all been chosen.

Use the input controls and the **Next** button(s) to fill in the screen as follows:



Note that an option to **remove** appears next to each input previously submitted; this will remove the current input, but keep the other inputs you have specified (as far as it can; if they are conditional on the input you have removed, then they will be, out of necessity, removed too).

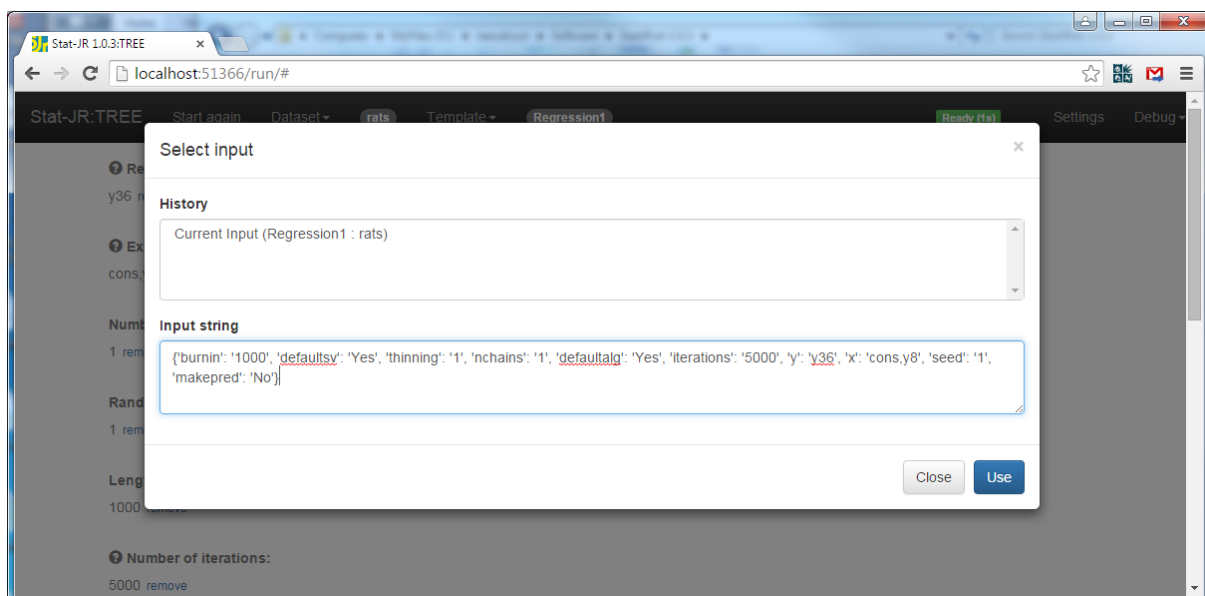
So, here we are performing a regression of the initial weight (y_8) on the final weight (y_{36}), including an intercept (*cons*).

The other inputs refer to the Monte Carlo Markov chain (MCMC) estimation procedures in Stat-JR. MCMC estimation methods are simulation-based, and so require certain parameters to be set. The methods involve taking a series of random (dependent) draws from the posterior distribution of the model parameters in order to summarise each parameter. The inputs required here are as follows:

- *the number of chains*: this is the number of starting points from which we will take random draws;
- *random seed*: the value from which random numbers are initially drawn. This allows repeatability, as a run using the same starting values and random seed will give the same answers. When multiple chains are used this seed is generally multiplied by the chain number to give a unique seed for each chain;
- *length of the burnin*: the initial length of the chain (i.e. the number of iterations at the start) which are excluded from the parameter summaries (the rationale for this is explained a little further in the example, below, with the *tutorial* dataset);
- *number of iterations*: the length of chain following the burnin, from which the parameter summaries are drawn;
- *thinning*: this determines how often the values are stored: i.e. store every n th iteration.

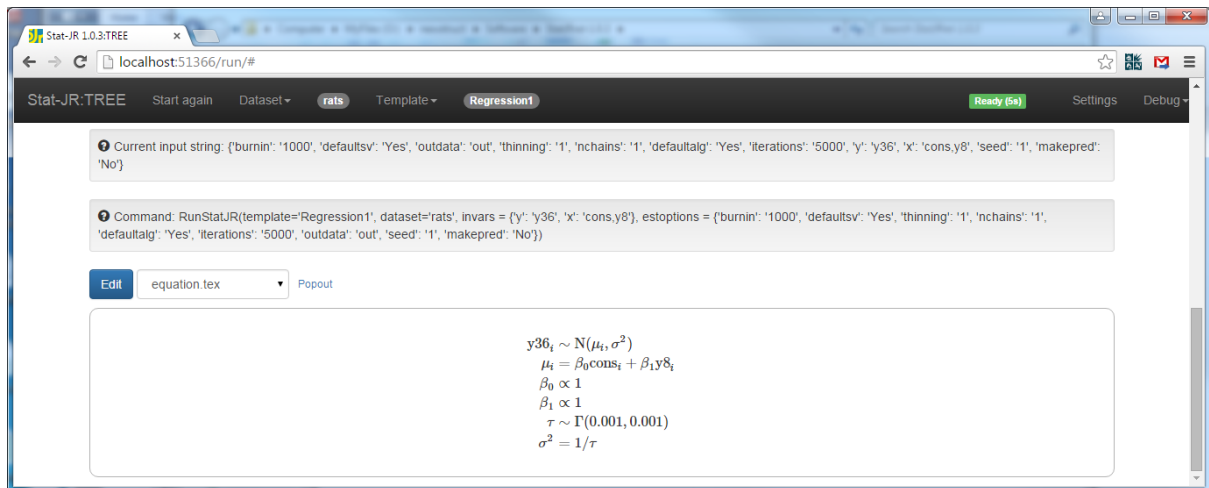
By answering 'Yes' to the question *use default algorithm settings*, we have used defaults for other settings for which we will therefore not be prompted to complete. By answering 'No' to *generate prediction dataset* we have chosen not to generate a dataset of predictions from our model. By answering 'Yes' to *use default starting values* we have chosen not to start the chain at values of our choosing, instead accepting Stat-JR's defaults. We will discuss MCMC estimation in slightly more detail in the applications in the next section. The final input we're asked for is the *name of output results*: this is the name (here we've chosen *out*) given to any dataset of parameter chains that results from running the template.

You will notice, towards the bottom of the window, a box with a rather long text string labelled **Current input string** above it and another labelled as **Command** below it. The input string allows the user to specify all the inputs directly and this is done via the **Set Inputs** option in the **Template** pull down list, without having to point-and-click through the list as we have done. These have to be formatted in a certain way, however, as illustrated by the current (*Input String*) text string which Stat-JR has written for us as a result of our inputs.

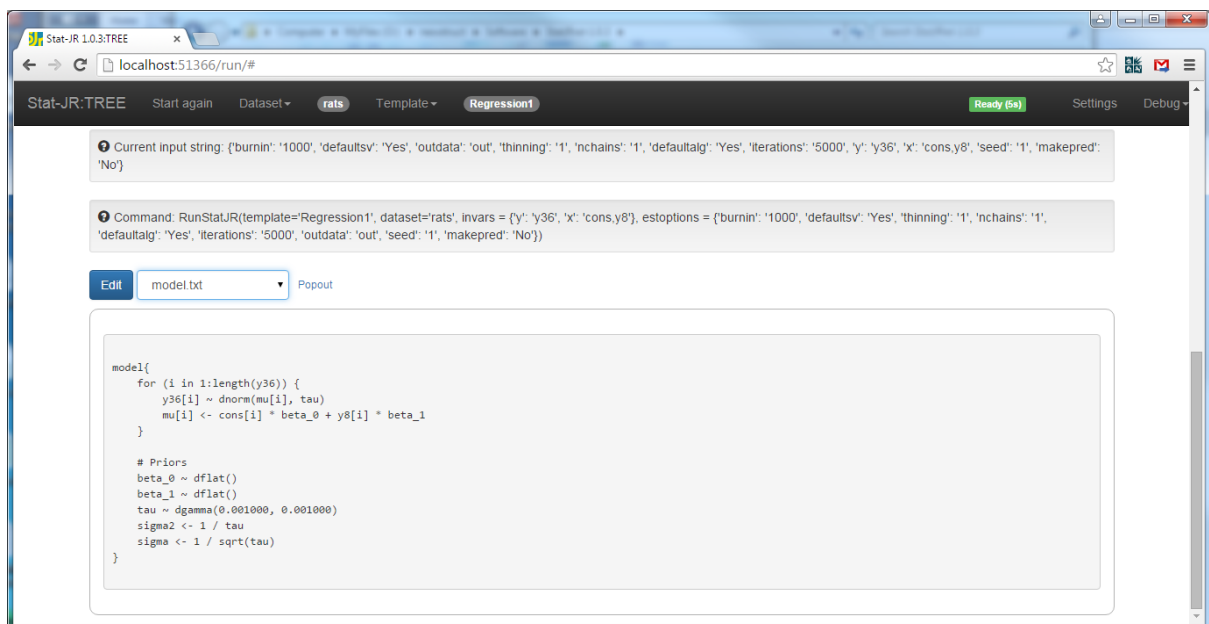


This (i.e. the string between, and including, the curly brackets: in this example `{'burnin': '1000'... 'makepred': 'No'}`) can be copied and pasted into the box labelled **input string** in the above window, and the **Use** button pressed (following any edits you would like to make to the input values), in order to specify inputs directly. Alternatively one can use the history feature to revert to older inputs. Returning to the main window, the second text string (labelled *Command*) can be used by the command driven version of Stat-JR to perform the same operations, although we will not discuss this further here.

Clicking on the **Next** button will now pre-process the template inputs, and will result in the following new pane at the bottom of the window:



The object currently specified in the pull-down list (*equation.tex* is selected by default here) appears in the pane below it. These objects are any outputs constructed by Stat-JR before and during the execution of the template, so here we see a nice mathematical description of the model. If we now select the object *model.txt* from the list we see a description of the regression model that we wish to fit in the language that is used by the eStat engine:



At this point we haven't actually run the template, and so the objects that can be selected from the pull down list are those present pre-model run, and include computer code to actually fit the model.

Click the **Run** button to run the template.

Once the progress gauge, towards the right of the black title bar, has changed from "Working" (blue) to "Ready" (green), select *ModelResults* from the pull-down list.

The screen will then look as follows:

Stat-JR 1.0.3: TREE

localhost:51366/run/#

Stat-JR: TREE Start again Dataset **rats** Template **Regression1** Ready (1s) Settings Debug

Extra Iterations: More

Download Add to ebook

Current input string: {'burnin': '1000', 'defaults': 'Yes', 'outdata': 'out', 'thinning': '1', 'nchains': '1', 'defaultalg': 'Yes', 'iterations': '5000', 'y': 'y36', 'x': 'cons,y8', 'seed': '1', 'makepred': 'No'}

Command: RunStatJR(template='Regression1', dataset='rats', invars = {'y': 'y36', 'x': 'cons,y8'}, estoptions = {'burnin': '1000', 'defaults': 'Yes', 'thinning': '1', 'nchains': '1', 'defaultalg': 'Yes', 'iterations': '5000', 'outdata': 'out', 'seed': '1', 'makepred': 'No'})

ModelResults Popout

Results

Parameters:

parameter	mean	sd	ESS	variable
tau	0.00418080675558	0.00110904973838	4327	
beta_0	169.180410959	31.4018966849	23	cons
beta_1	1.02242050055	0.205619681782	23	y8
sigma2	257.218314014	73.2560156192	4429	
sigma	15.8879481878	2.18893042306	4397	
deviance	250.689939715	2.28300617933	403	

Model:

Statistic	Value
Dbar	250.689939715
D(thetabar)	248.013138421
pD	2.67680129348
DIC	253.366741008

Here we see parameter estimates, along with standard deviations (SDs), as a measure of precision for each parameter. We will explain these further in the next section. At the top of the screen shot above (which is in fact the middle of the full window, vertically-speaking) we now have a few additional buttons. The **Extra Iterations** box, along with the **More** button, will allow us to run for longer (i.e. for a number of iterations additional to those we have already run for). The **Download** button will produce a zipped file that contains a folder with files for many of the objects contained in the two pull-down lists whilst the **Add to ebook** button can be used if one wants to construct an ebook to be used with the *DEEP* eBook interface into Stat-JR.

You'll recall that we earlier named the output results *'out'*, so if we choose this from the pull-down list just above the output pane, we'll be able to view it, as follows:

Stat-JR 1.0.3: TREE
localhost:51366/run/#
Stat-JR: TREE Start again Dataset **rats** Template Regression1 Ready (1s) Settings Debug

Command: RunStatJR(template='Regression1', dataset='rats', invars = {'y': 'y36', 'x': 'cons,y8'}, estoptions = {'burnin': '1000', 'defaults': 'Yes', 'thinning': '1', 'nchains': '1', 'defaultalg': 'Yes', 'iterations': '5000', 'outdata': 'out', 'seed': '1', 'makepred': 'No'})

out Popout

iteration	chain	tau	beta_0	beta_1	sigma2	sigma	deviance
1	1	0.00288518002951	205.489839661	0.819819491159	346.598822178	18.6171647191	253.491151791
2	2	0.00355026627789	196.364972102	0.823041334967	281.669013456	16.7829977494	250.465223393
3	3	0.00337715317864	201.360440503	0.805661906673	296.107386044	17.2077710946	249.839227094
4	4	0.00488282814749	205.747671086	0.784315702633	204.799343699	14.3108121258	249.587746728
5	5	0.00299592578601	203.383693904	0.793137713746	333.786639398	18.2698286656	250.783047997
6	6	0.00394771670837	202.148374143	0.815630644772	253.310982999	15.9157463852	249.301144304
7	7	0.00243484081172	196.10488089	0.804675385154	410.704467901	20.2658448603	255.768571423
8	8	0.00439076387497	199.260547643	0.83049110559	227.750803385	15.0914148901	248.959684971
9	9	0.00651709766023	198.173447585	0.792648488211	153.442537175	12.3871924654	259.508753637
10	10	0.00284766420737	206.488667221	0.772200262327	351.164999515	18.7393969891	251.38549422
11	11	0.00410410577617	209.647455254	0.794806635121	243.658437316	15.6095623678	253.280948695
12	12	0.00290512711805	208.12525737	0.802042431984	344.219016712	18.5531403464	253.474956979
13	13	0.00298889148998	201.023619949	0.822192714798	334.572199544	18.2913148665	250.74923858
14	14	0.0039487089206	194.258616129	0.853870178626	253.247332257	15.9137466442	248.832562639
15	15	0.00398014644491	195.373724985	0.87587350772	251.247036721	15.850740102	250.410945995
16	16	0.00274330185388	195.617064614	0.862741807536	364.524231479	19.0925176831	251.607779298
17	17	0.00467992957874	194.845981125	0.864004465259	213.678428954	14.6177436342	249.056906502
18	18	0.00431451833524	190.916814268	0.865492729324	231.775582417	15.2241775613	249.229469469
19	19	0.0037582194272	200.328243082	0.805435982072	266.083452382	16.3120646266	249.670869303
20	20	0.0052505796049	200.325388889	0.80382956858	190.455164048	13.8005494111	250.543489206
21	21	0.00394930422779	199.926836813	0.846218456775	253.209158455	15.912547202	250.632325746
22	22	0.00502986456319	198.312899419	0.840389451779	198.812510245	14.1000890155	249.493896804
23	23	0.00498151618344	201.277482336	0.81888328435	200.742096016	14.168348387	249.51587256
24	24	0.00432901193958	200.951152711	0.834241583974	230.999593893	15.1986707936	250.074462576
25	25	0.00542238675379	198.163614546	0.810410168159	177.860407651	13.3364315936	252.17758698
26	26	0.00455976756075	202.871623973	0.803061550114	219.309424587	14.8090993848	249.165280566
27	27	0.00510966170594	204.421761084	0.820922602561	195.707672553	13.9895558383	252.621119088

View 1 - 30 of 5,000

Here we see columns containing the chains of values for each parameter in the model. As well as being able to view this file here, it is also a dataset (stored in temporary memory) and so will appear in the dataset list (at least for the duration of our current session using the software) accessible via the **Dataset** menu in the top title bar (emboldened to indicate that it has been created in this run of the software). This means that we can string templates together, as we can select *out* as a dataset and perform operations on it using another template.

This ends our whistle-stop tour of many of the windows in Stat-JR. We will next look at a practical application.

4.2 Application 1: Analysis of the tutorial dataset using the eStat engine

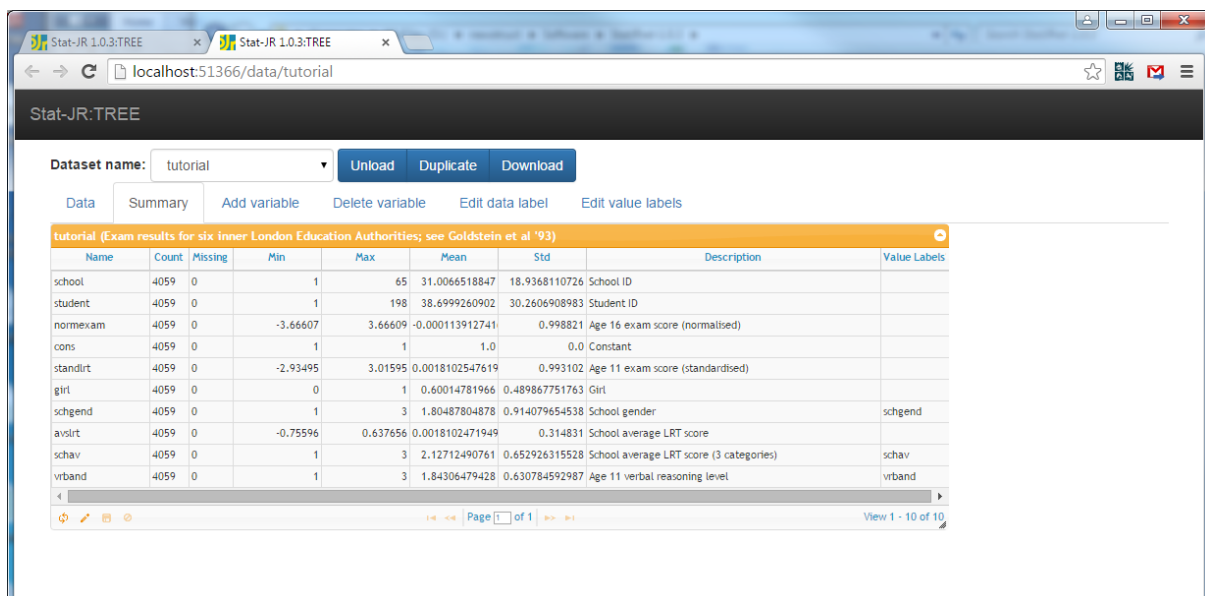
4.2.1 Summarising the dataset and graphs

In this section we will look at performing some analysis of an example dataset from education. The dataset in question is known as the **tutorial** dataset, and is used as an example in the MLwiN software manuals (see, for example, Browne 2012). In fact, much of the material here owes a lot to Browne (2012), which employs similar analysis but using MLwiN.

Let us start by looking at the tutorial dataset.

Select **tutorial** via **Dataset > Choose** (see the title bar), then click **Use**.

If you then select **Dataset > View**, and click on the Summary tab the following should appear in a new tab in the browser window containing summary information, as follows:



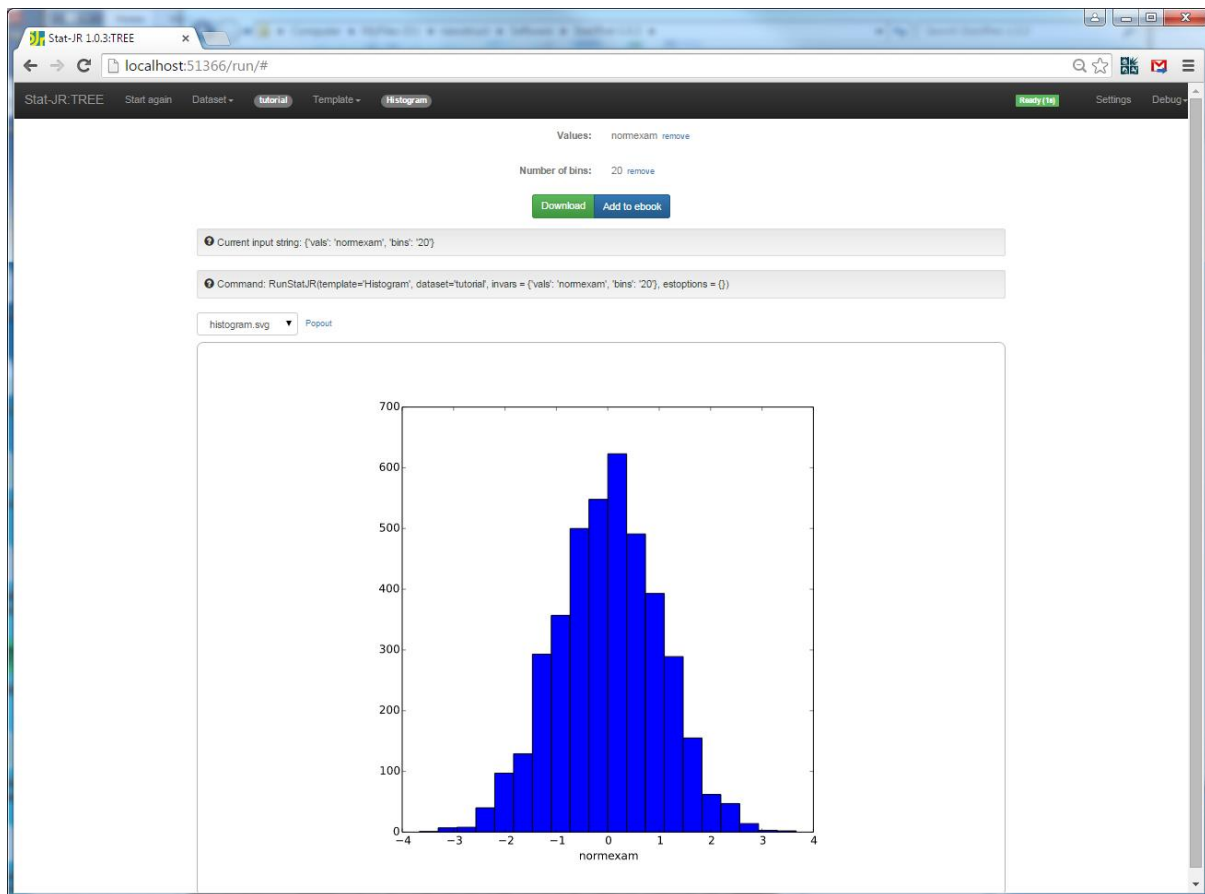
Name	Count	Missing	Min	Max	Mean	Std	Description	Value Labels
school	4059	0	1	65	31.0066518847	18.9368110726	School ID	
student	4059	0	1	198	38.6999260902	30.2606908983	Student ID	
normexam	4059	0	-3.66607	3.66609	-0.000113912741	0.998821	Age 16 exam score (normalised)	
cons	4059	0	1	1	1.0	0.0	Constant	
standlrt	4059	0	-2.93495	3.01595	0.0018102547619	0.993102	Age 11 exam score (standardised)	
girl	4059	0	0	1	0.60014781966	0.489867751763	Girl	
schgend	4059	0	1	3	1.80487804878	0.914079654538	School gender	schgend
avslrt	4059	0	-0.75596	0.637656	0.0018102471949	0.314831	School average LRT score	
schav	4059	0	1	3	2.12712490761	0.652926315528	School average LRT score (3 categories)	schav
vrband	4059	0	1	3	1.84306479428	0.630784592987	Age 11 verbal reasoning level	vrband

The **tutorial** dataset contains data on exam scores of 4059 secondary school children from 65 schools at age 16. These exam scores have been normalised to have a mean of zero and a standard deviation of one and are named *normexam*. There are several predictor variables, including a (standardised) reading test (*standlrt*) taken at age 11, the pupils gender (*girl*), and the school's gender (*schgend*) which takes values 1 for mixed, 2 for boys and 3 for girls. Each variable is described in the Description column and if you hover over any of the variables that say "Yes" in the value labels column, the category labels will be displayed.

We can explore the dataset in more detail, prior to fitting any models, by using the many data manipulation templates available in Stat-JR. We will first look at some plots of the data:

Select **Template > Choose** and then select **Histogram** from the template list that appears and click **Use**.

Fill in the inputs as shown below and click **Next** and then **Run** and select *histogram.svg* from the output list.



Here you will see, in the output pane, a histogram plot that shows that the response variable we will model, *normexam*, appears Normally-distributed.

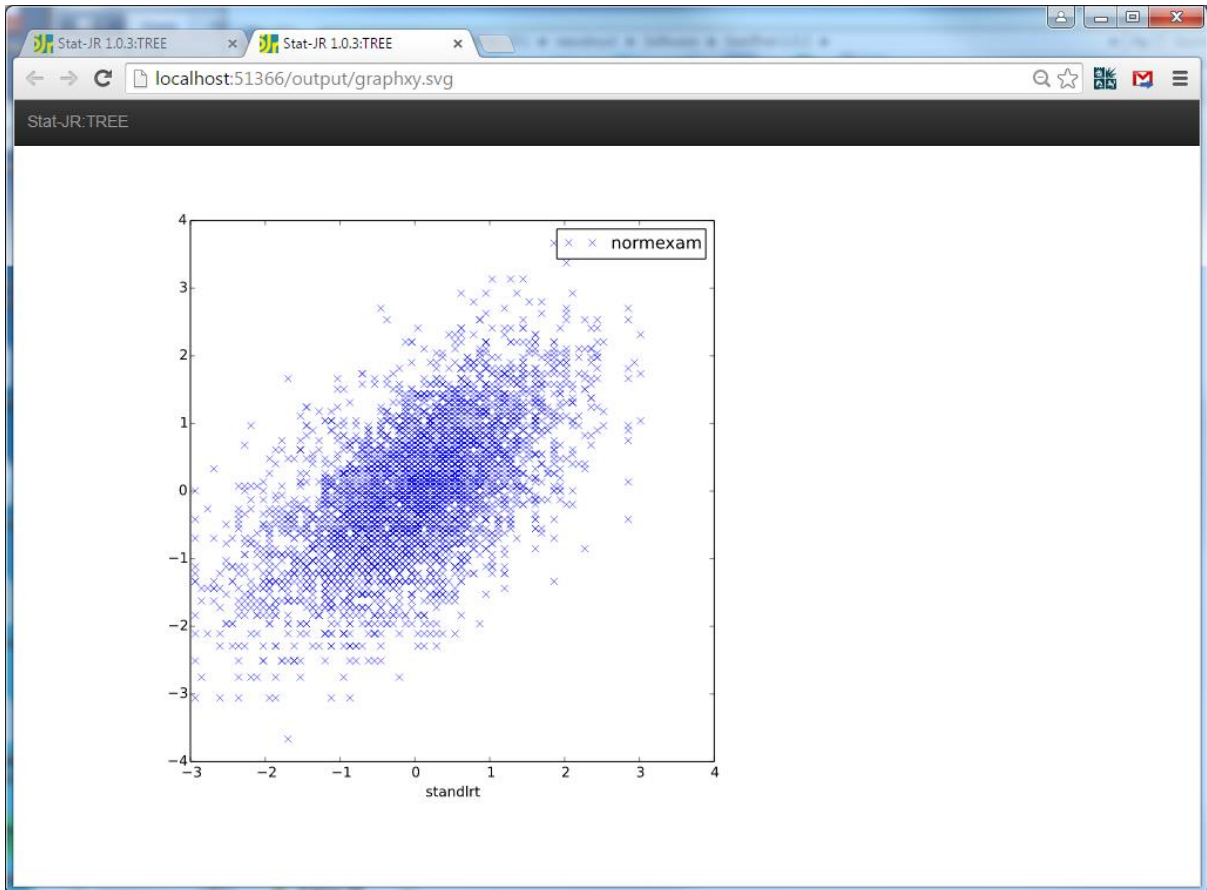
Select **Template > Choose** and this time select **XYPlot** from the template list, then click **Use**.

Fill in the inputs as shown below and click **Next** and then **Run** and select *graphxy.svg* from the list.



Here we see that there appears to be a positive relationship between *normexam* and *standlrt*, with pupils that have higher intake scores performing better, on average, at age 16.

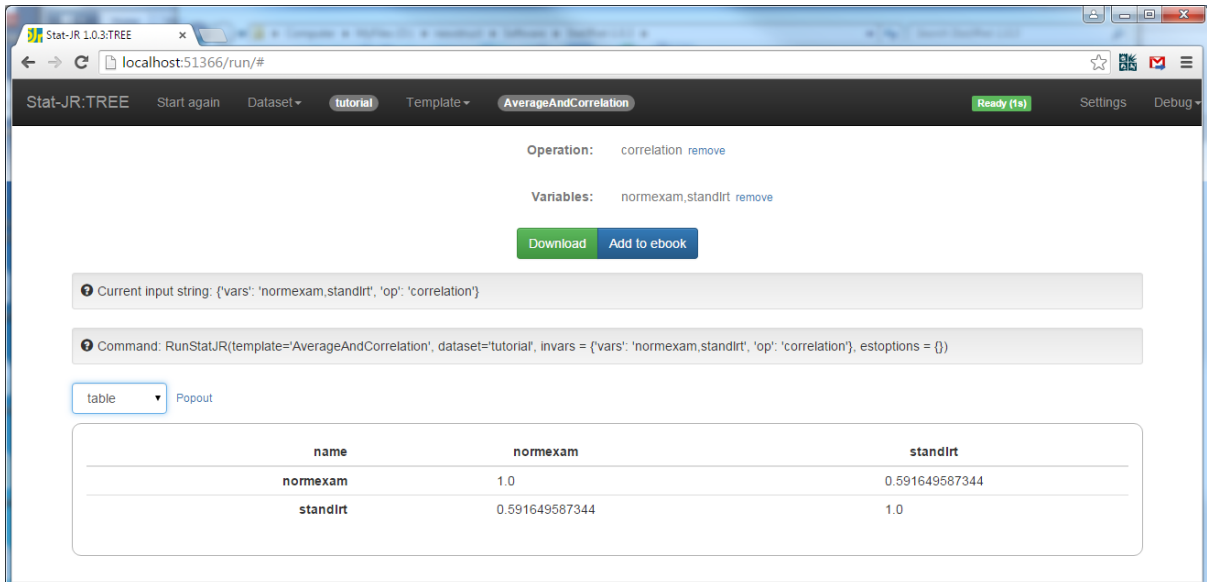
We can display the graph in a separate tab in the browser window by clicking on the **Popout** button next to the pull down list:



For the sake of brevity, for the remainder of this documentation we will assume you now know how to change template/dataset, and also how to display output in separate tabs, so we'll refrain from repeating this information in detail again.

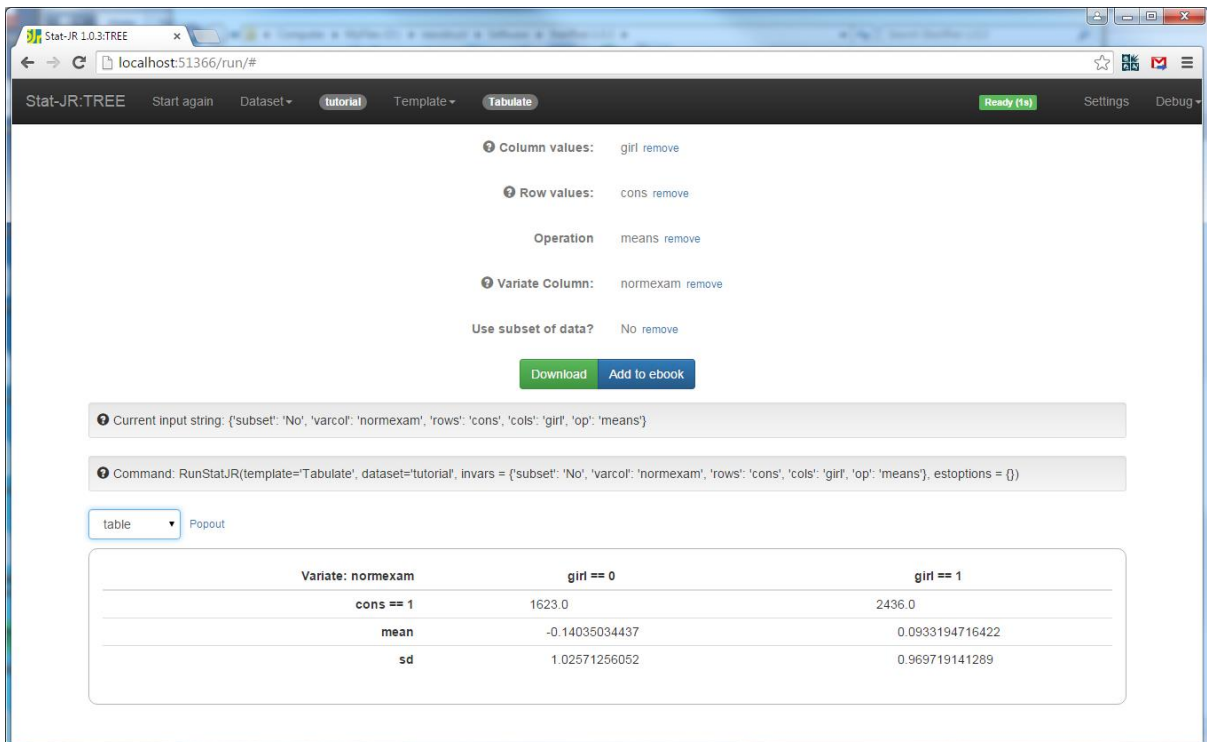
Next, we might like to examine how correlated the two variables, *normexam* and *standlrt*, actually are:

Select **AverageandCorrelation** as the template, and complete the inputs as follows before clicking on **Next** and **Run** and selecting **table** from the outputs:



Here we see that the correlation is 0.592, so fairly strong and positive. We might also like to look at how exam score varies by gender:

Select **Tabulate** as the template, and complete the inputs as follows, before clicking on **Next** and **Run** and selecting **table** from the output list:

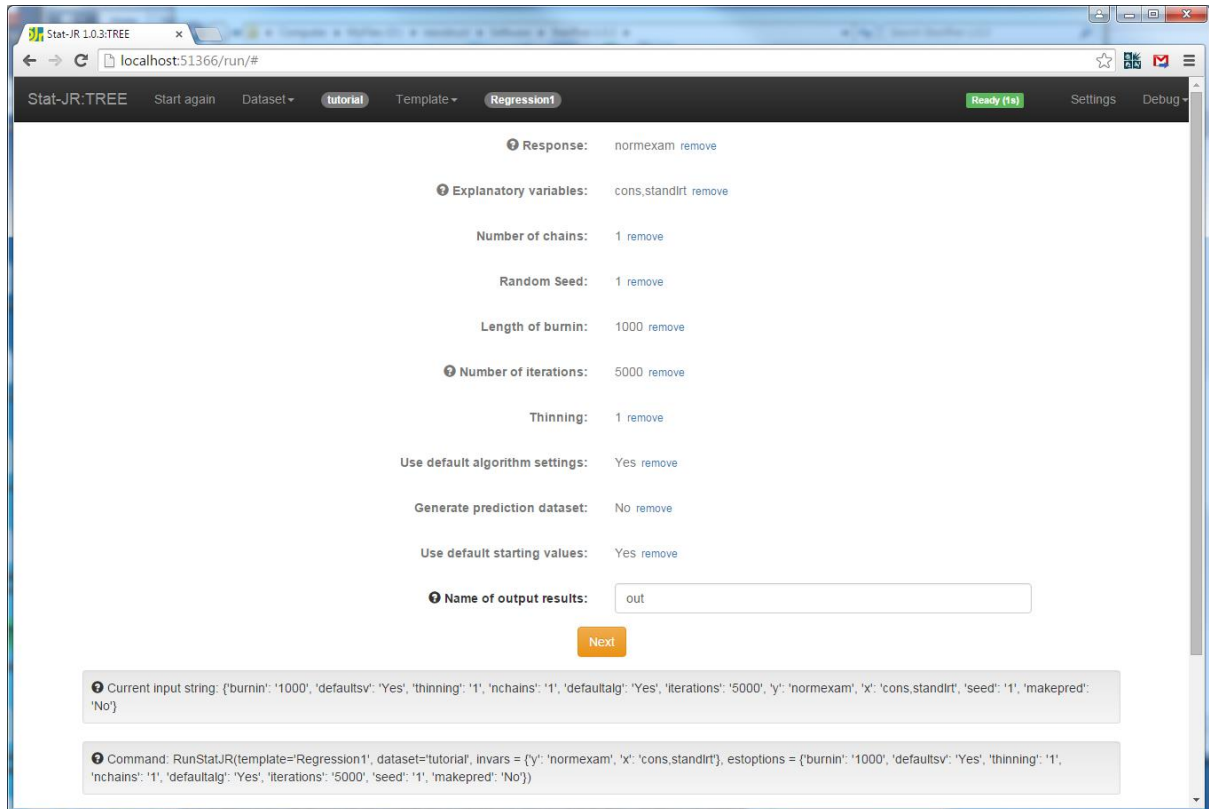


We have to enter variables for column values and row values, and so here we have specified column values as *girl* (taking value 1 for girls and 0 for boys) and row values as *cons* (which is a constant), and then we get 2 columns in the output labelled 0 and 1 for boys and girls, respectively. Looking at the means, it appears that girls do slightly better than boys, and looking at the standard deviations (sds) they are slightly less variable than boys in their scores. Let us now consider performing some statistical modelling on the dataset.

4.2.2 Single-level Regression

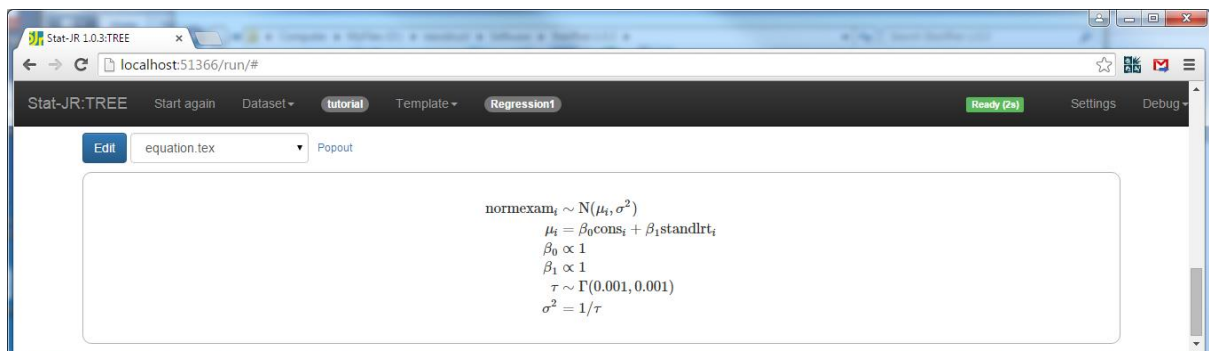
As in the last chapter, with the **rats** dataset, we will start by fitting a simple linear regression model to the **tutorial** dataset. Here we will regress *normexam* on *standlrt* by using a modelling template.

Select **Regression1** as the template and fill it in as follows:



Here we are fitting a linear regression, and so have *standlrt* as an explanatory variable, but also *cons* (which is a column of 1s) as we would like to include an intercept as well. For now we have set-up the MCMC estimation options as we did for the **rats** dataset, and we will overwrite the output file *out*.

Clicking on the **Next** button will populate a pull down list of objects created by Stat-JR at the bottom of the screen and by default we see the object *equation.tex* :



In the pane we find a mathematical representation of the chosen model. Note that the file is a LaTeX file that is being rendered in the browser by a piece of software called MathJaX (v2.3, 2013), so if you are a LaTeX-user you can copy this file straight into a document. If we instead choose *model.txt* from the list we see the following:

```

model{
  for (i in 1:length(normexam)) {
    normexam[i] ~ dnorm(mu[i], tau)
    mu[i] <- cons[i] * beta_0 + standlrt[i] * beta_1
  }

  # Priors
  beta_0 ~ dflat()
  beta_1 ~ dflat()
  tau ~ dgamma(0.001000, 0.001000)
  sigma2 <- 1 / tau
  sigma <- 1 / sqrt(tau)
}

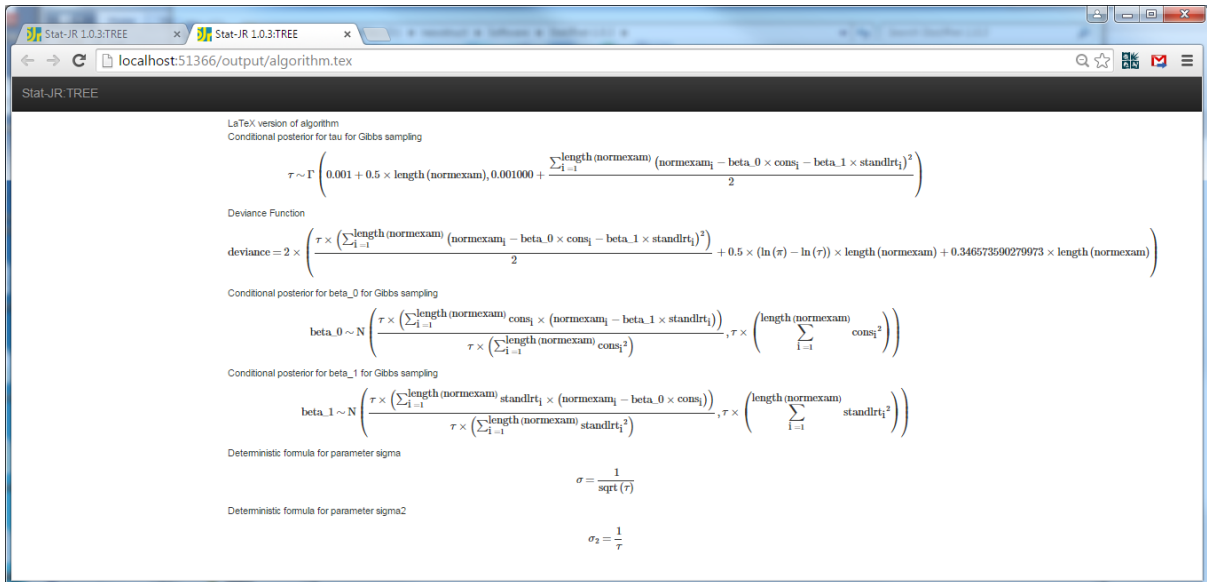
```

Here we see the text file that represents the model we wish to fit in the language that the algebra system used by the built-in eStat engine requires. The **Regression1** template only uses the eStat MCMC-based estimation engine, so as you can see in the mathematical formulae in *equation.tex* we are fitting a Bayesian version of a linear regression, and the last four lines of the output are prior distributions for the unknown parameters, β_0 , β_1 and the precision τ (where $\tau=1/\sigma^2$).

Whilst we will keep our description of Bayesian statistics and MCMC estimation to a minimum, and recommend Chapter 1 of Browne (2012) for more details, in brief we are interested in the joint posterior distribution of all unknown parameters given the data (and the prior distributions specified). In practice, in complex models, this distribution has many dimensions (in our simple regression we have 3 dimensions) and is hard to evaluate analytically. Instead, MCMC algorithms work by simulating random draws from a series of conditional posterior distributions (which can be evaluated). It can then be shown (by some mathematics) that after a period of time (required for the simulations to move from their possibly arbitrary starting point) that the draws will be a dependent sample from the joint posterior distribution of interest. It is common, therefore, to throw away the first n draws which are deemed a **burn-in** period.

For the simple linear regression, it is a mathematical exercise to show that the conditional posterior distributions have standard forms and are Normal (for the fixed effect) and Gamma (for the precision = 1 /variance). The eStat engine has a built in algebra system which takes the text file (*model.txt*) in the left-hand pane and returns the conditional posterior distributions; you can view these as follows:

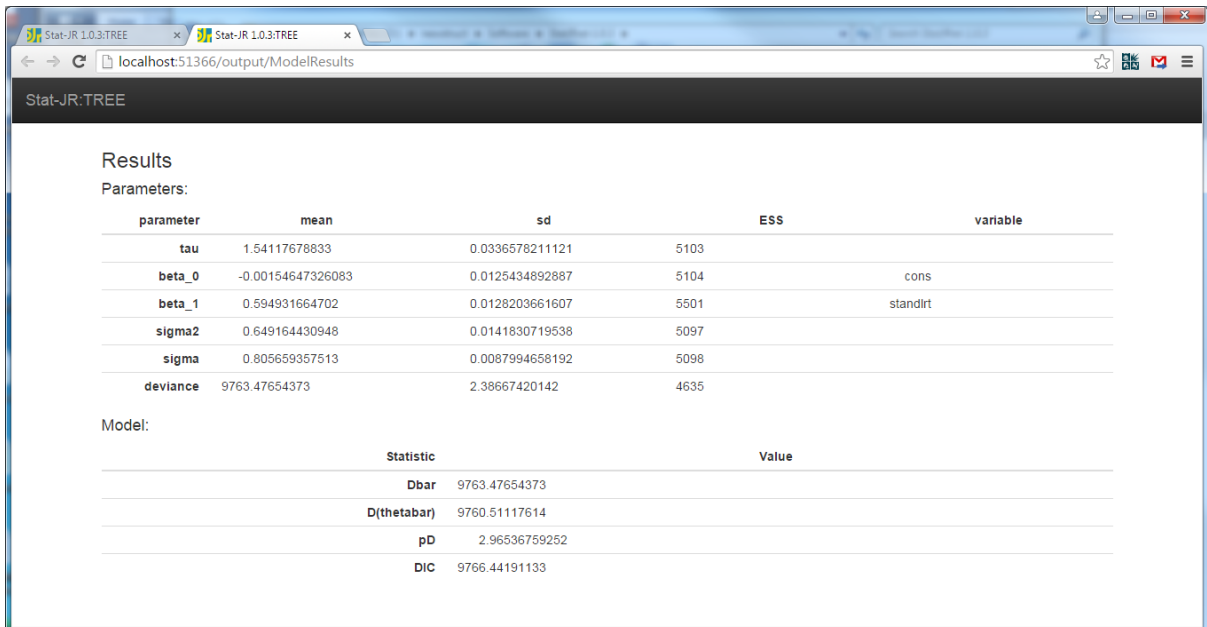
Select *algorithm.tex* from the list and click on the **Popout** button and the algebra steps will appear in a new tab as follows:



The eStat engine then takes these posterior distributions and wraps them up into computer code (C++) which it will compile and run for the model. By default this will be several pieces of code that are linked together by Stat-JR, although the **Settings** screen (accessible via a link towards the top of the main menu screen, as we saw earlier) has an option to output completely standalone code that can be taken away and run separately from the Stat-JR system; this is, however, a topic for more advanced users.

Returning to the tab, in the browser window, containing the model template, click on the **Run** button and wait for the model to run.

Then select *ModelResults* from the pull down list and pop it out into a separate tab.



Here the model results can be split into two parts:

The first part of the results (under the heading ‘*Parameters*’) contains the actual parameter estimates. Here, for each parameter, we get 3 numbers: a posterior mean estimate (*mean*), a posterior standard deviation (*sd*), and an effective sample size (*ESS*).

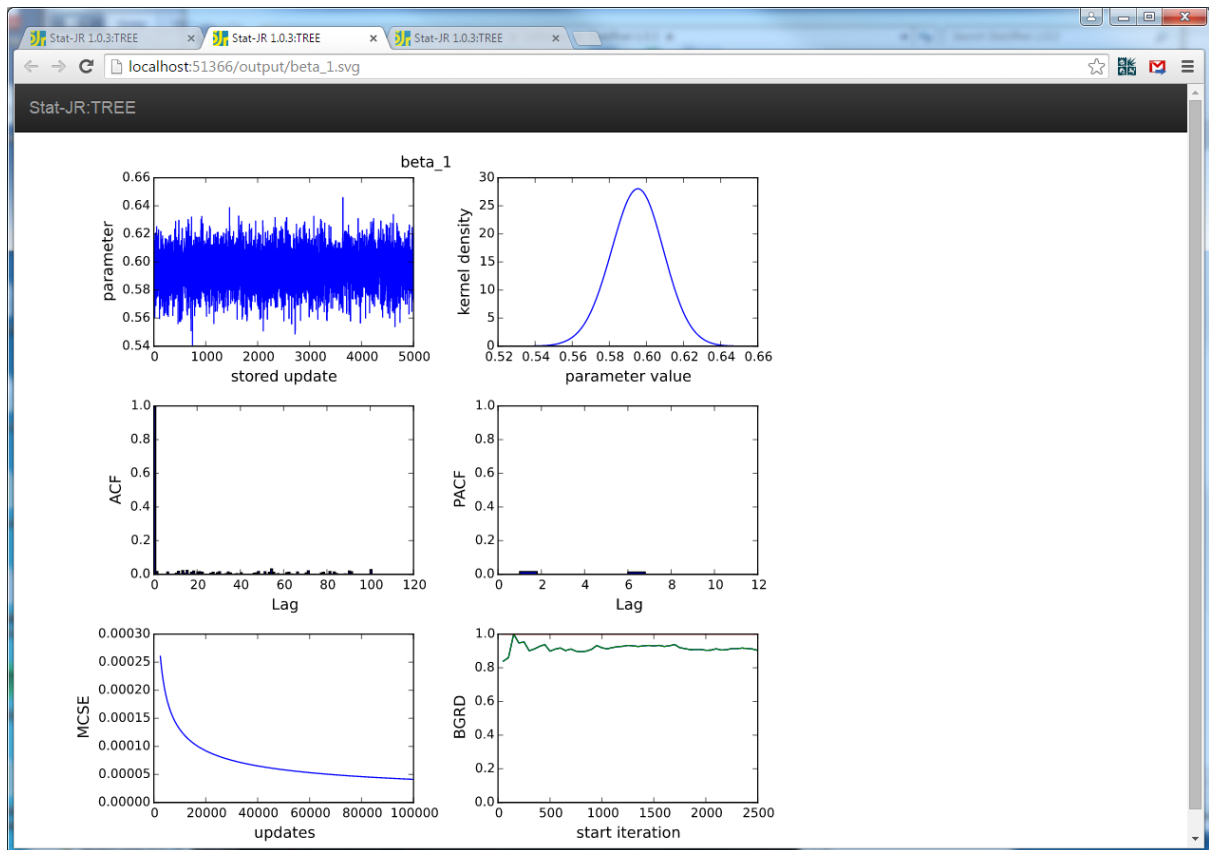
Here we see that *beta_0* has a mean estimate of approximately 0, which we would expect as both the response and predictor have been normalised, or standardised. The slope *beta_1* has mean 0.595 with standard deviation 0.013, and is highly significant, as its mean estimate is many times its standard deviation (a Bayesian equivalent of a standard error) The value 0.595 represents the average increase in the *normexam* score for a 1-point (1 sd, due to standardising) increase in *standlrt*. The residual variance, *sigma2*, has value 0.649 showing that, as the initial response variance was 1.0, *standlrt* has explained 35.1% of the variability.

The ESS is a diagnostic which reflects the simulation-based (stochastic) nature of the MCMC estimation procedure: we have based our results on the 5,000 iterations post burn-in, but we know that the method produces dependent samples, and so the ESS gives an equivalent number of independent samples for the parameters involved; in effect a measure of the information content of the chain In this case, all parameters have ESS of > 4000, and so the chains are almost independent.

The second part (under the heading ‘*Model*’) refers to the model fit for this particular model and the DIC diagnostic (Spiegelhalter et al. 2002). The DIC diagnostic is an information criterion which is a measure of how good a specific model is, consisting of a combination of how well the model fits the data (here defined by the model deviance) and how complex the model is (here defined by *pD*: the effective number of parameters). Basically the better fitting the model is, the better the model is, but it has to be penalised by how complex it is. The DIC statistic is defined as the deviance of the mean + 2*pD*. In this example the deviance at the mean (*D*(*thetabar*)) is 9760.5 and *pD* is ~3 (reflecting the three parameters of the model that are being estimated) and so we have a DIC value of 9766.4. This number is not particularly interesting in isolation but it is when we compare values for several models.

We can also get more information from the diagnostic plots that are available in the list of objects

Return to the model run tab in the browser window, and select *beta_1.svg* from the pull-down list above the output pane and pop it out into a separate tab.



This “sixway” plot gives several graphs that are constructed from the chain of 5,000 values produced for *beta_1*. The top-left graph shows the values plotted against iteration number, and is useful to confirm that the chain is ‘mixing well’, meaning that it visits most of the posterior distribution in few iterations. The top-right graph contains a kernel density plot which is like a smoothed histogram and represents the posterior distribution for this parameter. Here the shape is symmetric and looks like a Normal distribution which we expect given theory for fixed effects in a normal model.

The two graphs in the middle row are time series plots known as the autocorrelation (ACF) and partial autocorrelation (PACF) functions. The ACF indicates the level of correlation within the chain; this is calculated by moving the chain by a number of iterations (called the lag) and looking at the correlation between this shifted chain and the original. In this case, the autocorrelation is very small for all lags. The PACF picks up the degree of auto-regression in the chain. By definition a Markov chain should act like an autoregressive process of order 1, as the Markov definition is that the future state of the chain is independent of all the past states of the chain given the current value. If, for example, in reality the chain had additional dependence on the past 2 values, then we would see a significant PACF at lag 2. In this case all PACF values are negligible. All of this suggests that we have good mixing and it would be appropriate to proceed to the interpretation of the parameters.

The bottom-left plot is the estimated Monte Carlo standard error (MCSE) plot for the posterior estimate of the mean. As MCMC is a simulation-based approach this induces (Monte Carlo) uncertainty due to the random numbers it uses. This uncertainty reduces with more iterations, and is measured by the MCSE, and so this graph details how long the chain needs to be run to achieve a specific MCSE. The sixth (bottom-right) plot is a multiple chains diagnostic and doesn’t make much

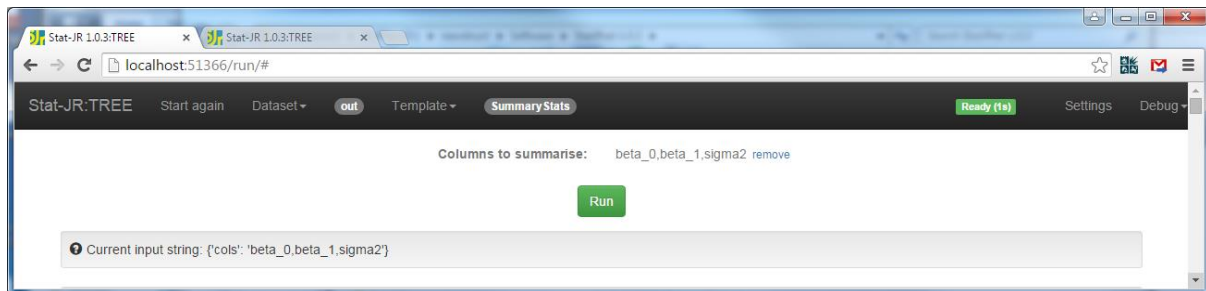
sense when we have run only one chain, and we will therefore consider running multiple chains in the next section.

We can also get some other diagnostics and summary statistics for the model as follows:

Click on the **Template** pull down list at the top of the screen and select **Choose** and **SummaryStats** as the template.

Next click on the **Dataset** pull down list and select **Choose** and **out** as the dataset.

Run the **SummaryStats** template and select the inputs as follows before clicking on **Run**:



Now select **table** from the drop-down list of outputs, and display it in a separate tab:

The screenshot shows the Stat-JR: TREE web interface displaying a summary statistics table. The browser address bar shows 'localhost:51366/output/table'. The table has columns for 'name', 'beta_0', 'beta_1', and 'sigma2'. The data is as follows:

name	beta_0	beta_1	sigma2
N	5000	5000	5000
mean	-0.00154647326083	0.594931664702	0.649164430948
sd	0.0125434892887	0.0128203661607	0.0141830719538
median	-0.00137256910389	0.595057913446	0.649016571249
min	-0.0423723866943	0.54096895575	0.604003596165
max	0.0457521057663	0.645871124458	0.70543401519
2.5%	-0.0260586560258	0.568970859499	0.622427120967
5%	-0.0220177345695	0.57359082263	0.626297012023
50%	-0.00137256910389	0.595057913446	0.649016571249
95%	0.0191401171679	0.615577564578	0.672825396309
97.5%	0.0231234059113	0.61974651492	0.676982988991
IQR	0.0168228609912	0.0168564529234	0.019379812426
ESS	5104	5501	5097
BD	240935	27	32

Here we see a more extensive summary of the three parameters of interest. This summary table includes various quantiles of the distribution which are calculated by sorting the chain and picking the values that lie x% into the sorted chain (where x is 2.5, 5, 50 etc.). These allow for accurate interval estimates that do not rely on a Normal distribution assumption. The inter-quartile range (IQR) is similarly calculated by picking the values that lie 25% and 75% through the sorted list and calculating the distance between them.

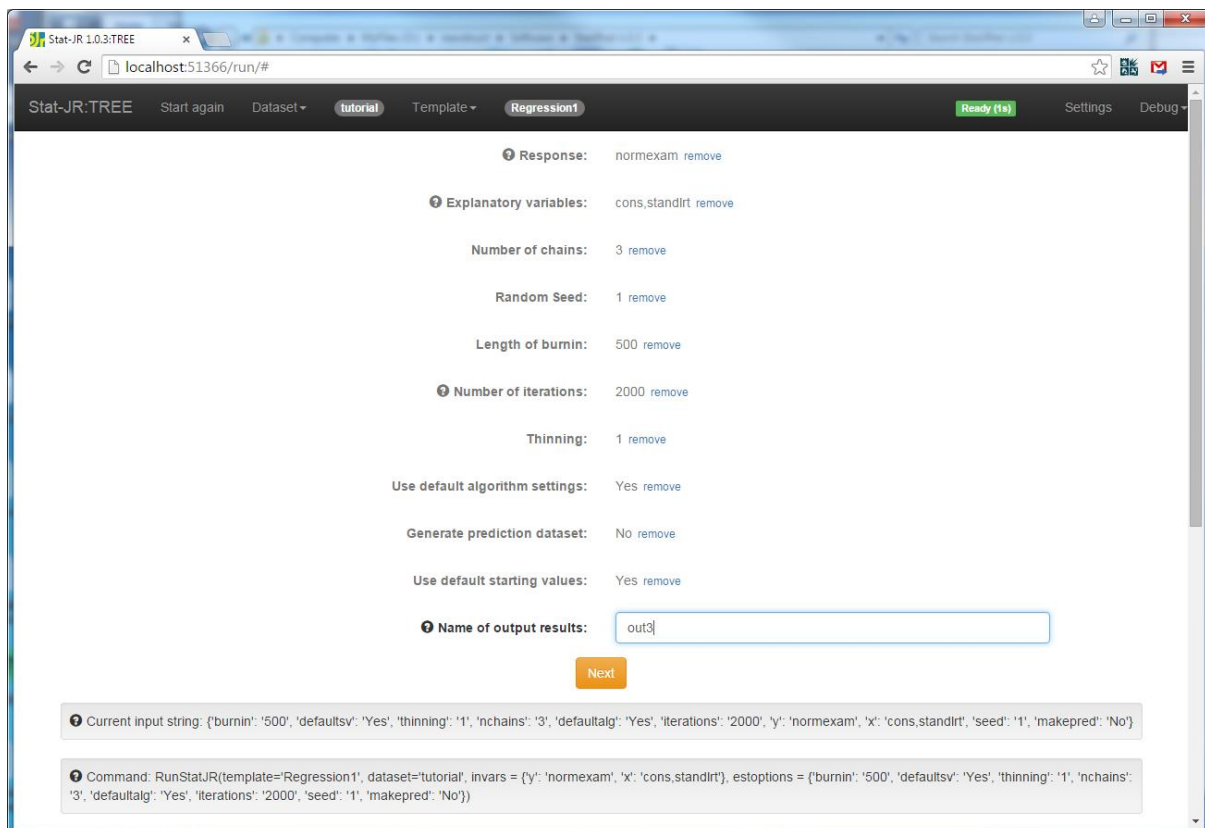
The final statistic is an MCMC diagnostic designed to suggest a length of chain to be run. The Brooks-Draper diagnostic is based on measuring the mean estimate to a particular accuracy (number of significant figures set to 2 by default). For example, it states that to quote σ_2 as 0.65 with some desired accuracy only requires 32 iterations. The anomaly here is β_0 , however, since the true value is 0 we have difficulty quoting such a value to 2 significant figures!

4.2.3 Multiple chains

MCMC methods are more complicated to deal with than classical methods as we have to specify many estimation parameters, including how long to run the MCMC chains for. The idea of running chains for a longer period is to counteract the fact that the chains are serially-correlated, and therefore are not independent samples from the distribution. Another issue that might cause problems is that the posterior distribution of interest may have several possible maxima (i.e. may be multimodal). This is not usually an issue in the models we cover in this book, but it is still a good idea to start off the estimation procedure from several places, or with several runs with different random number seeds, to confirm we get the same answers.

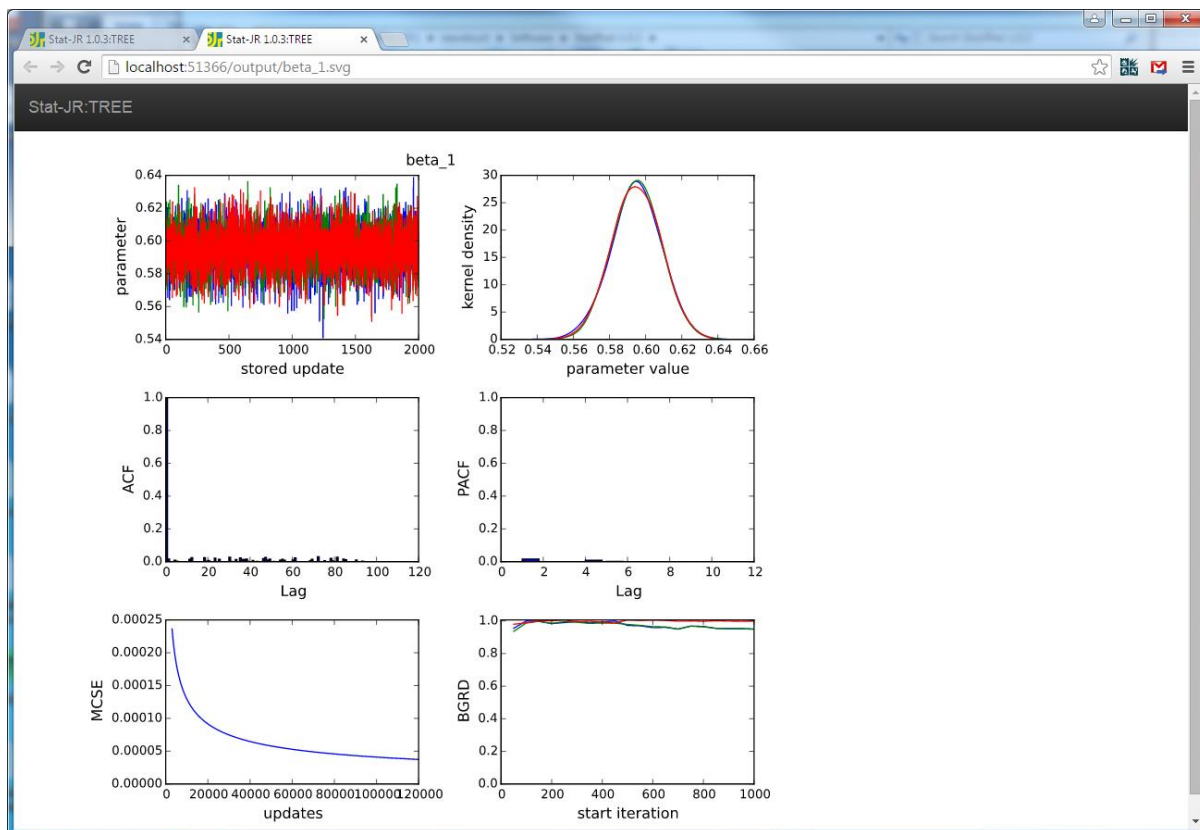
From the top bar change **Template** and **Dataset** using the respective pull down lists and **Choose** so you have **Regression1** as the template and **tutorial** as the dataset.

This time fill in the screen as follows:



Click on the **Next** and **Run** buttons.

When the model has run select *beta_1.svg* from the outputs list and pop it out to view it in a new tab.



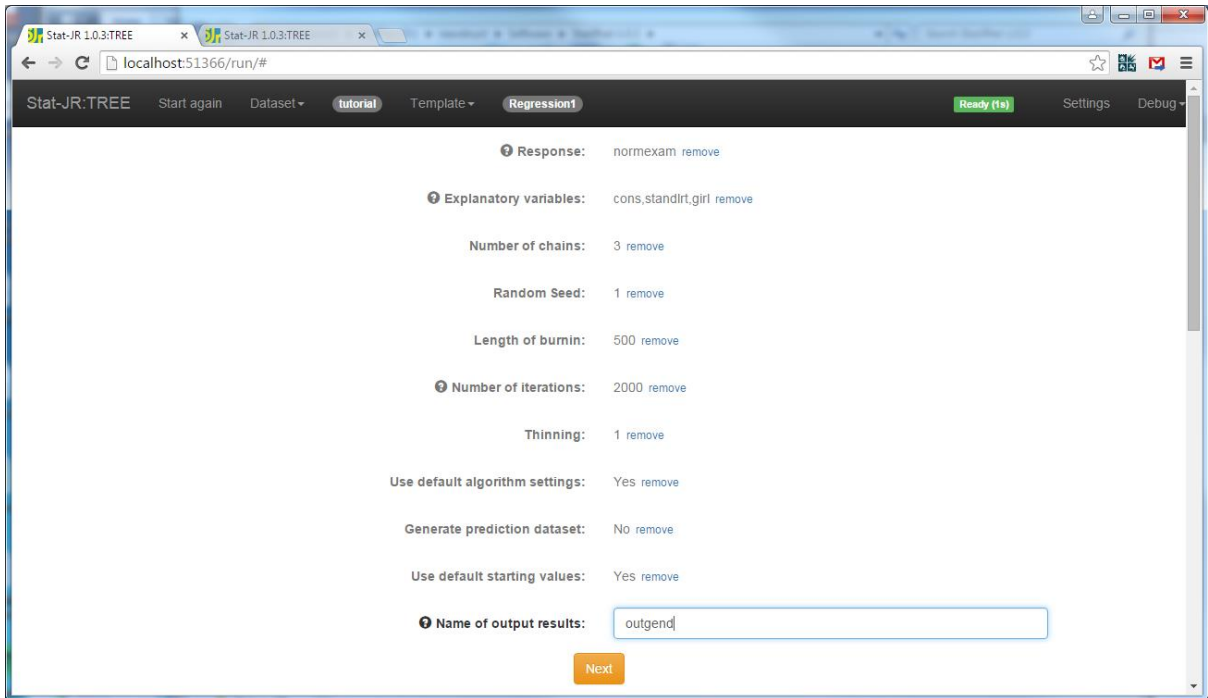
Here we see the three chains superimposed on each other in the top-left pane – note the chain looks primarily red simply because this chain (chain 3) has been plotted on top of the other two, and due to good mixing obscures them. Each chain has its own kernel plot in the top-right pane and this also suggests that, by the similarity of shape and position, the chains are mixing well.

We have previously described what all the graphs here mean in Section 4.2.2, apart from the Brooks-Gelman-Rubin diagnostic plot (BGRD; Brooks and Gelman, 1998) in the bottom-right corner. This plot looks at mixing across the chains: the green and blue lines measure variability between and within the chains, and the red is their ratio. For good convergence this red line should be close to 1.0, and here the values get close to 1.0 fairly quickly. We can have a lot of faith in the estimates of our model.

4.2.4 Adding gender to the model

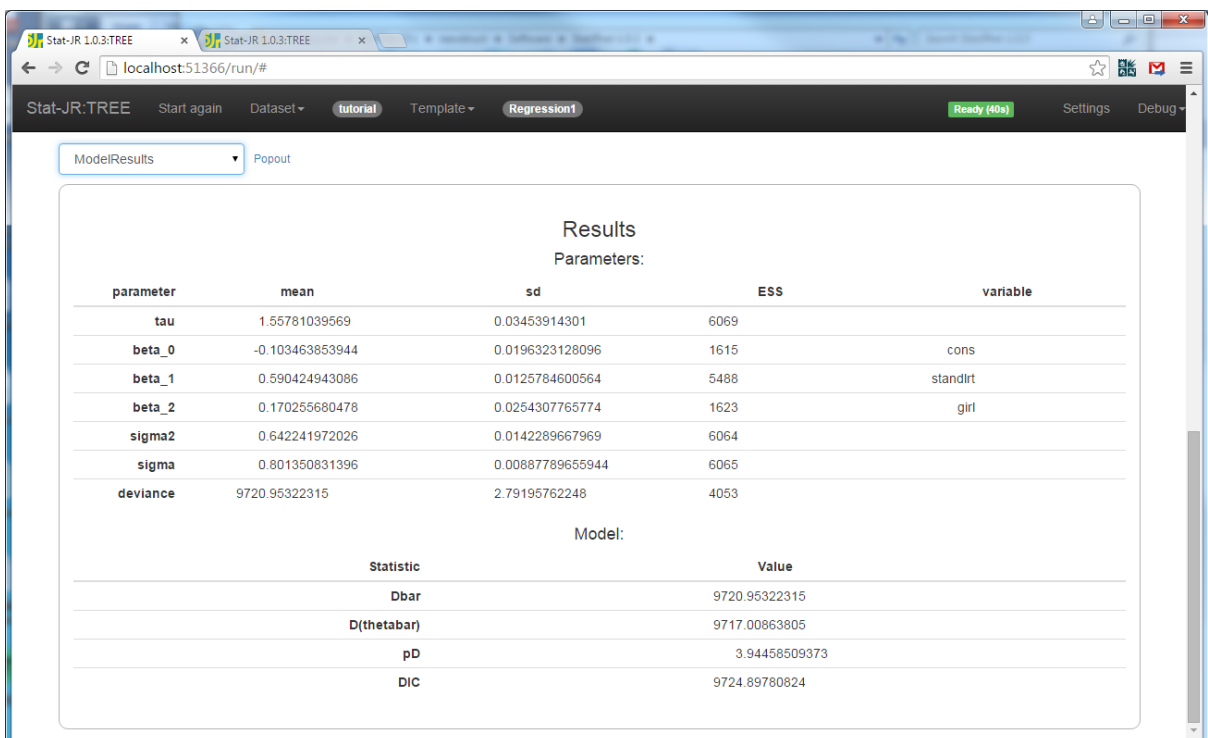
We have so far been more focused on understanding the MCMC methods but now we will return to modelling. We next wish to look at whether gender has an additional effect on *normexam* on top of that we have observed for intake score (*standlrt*).

To do this, click on the remove link next to explanatory variables in the browser window, and fill-in the template as follows:



Click on **Next** and then **Run** to run the model.

After the model finishes running select *ModelResults* from the drop-down list of outputs, and display in a new tab.

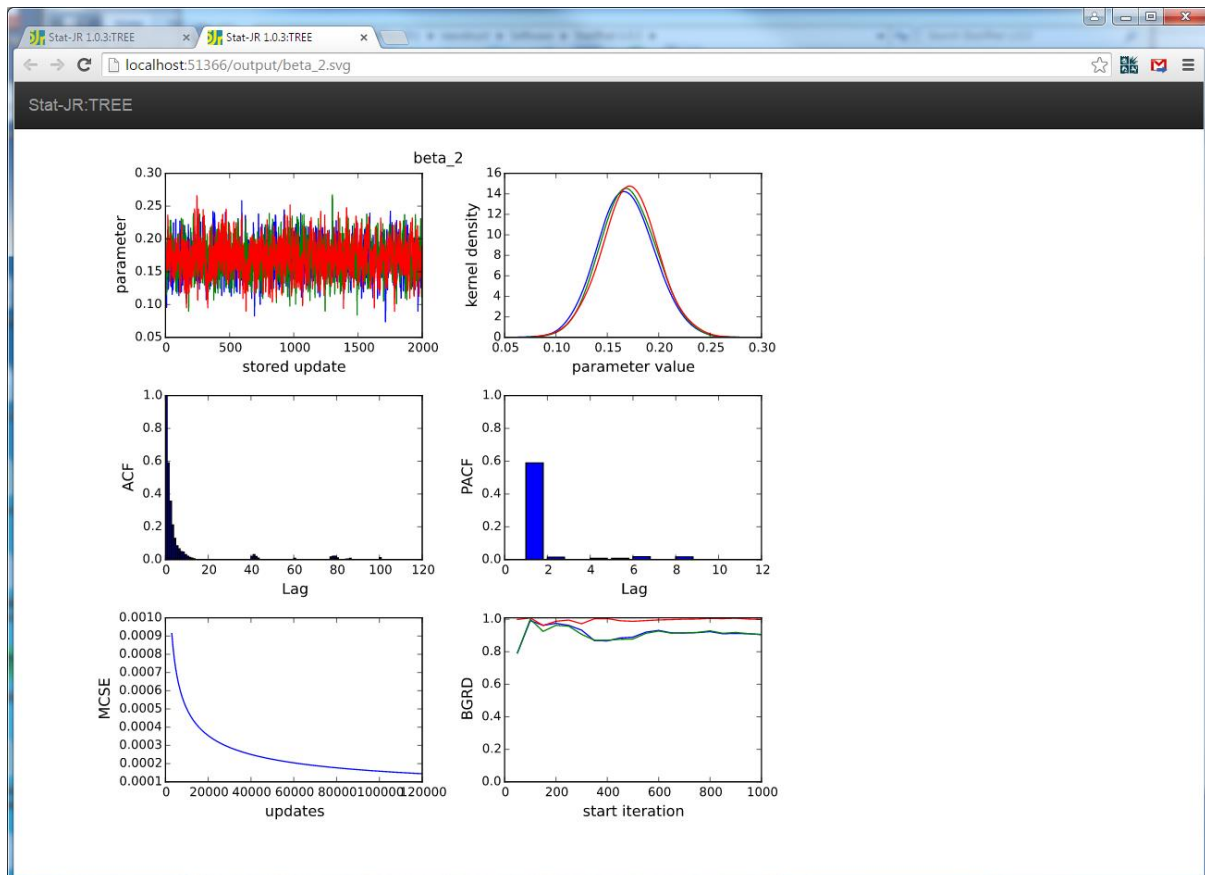


This new model has one additional fixed effect parameter (beta_2) associated with gender, and we see it has a positive effect (0.170) which appears highly-significant (at least twice its sd, which is 0.025). Note that in our earlier tabulation we saw that the difference in gender means was 0.093- (-

0.140) = 0.233 and so the effect here is somewhat smaller, probably due to correlation between gender and intake score.

Looking at the DIC diagnostic to assess whether this model is better we see this has dropped from 9766.4 to 9724.9, which is a big drop, and so the model with gender is indeed much better.

Finally we see that the ESS for two of the parameters is lower (beta_0 and beta_2), at around 1600, so the model doesn't mix quite as well; however, these ESS are still large enough not to require further iterations. Here is the graph for **beta_2.svg**, displayed in a new tab:



We see reasonable mixing, and can clearly see the significance of the effect as well (as the kernel density plot in the top-right corner indicates that 0 is nowhere near the posterior distribution). From a modelling perspective we have thus far ignored the fact that our data is a two-stage sample and that we should account for the clustering of the pupils within secondary schools. To do this we need to fit a 2-level model, and use a different template.

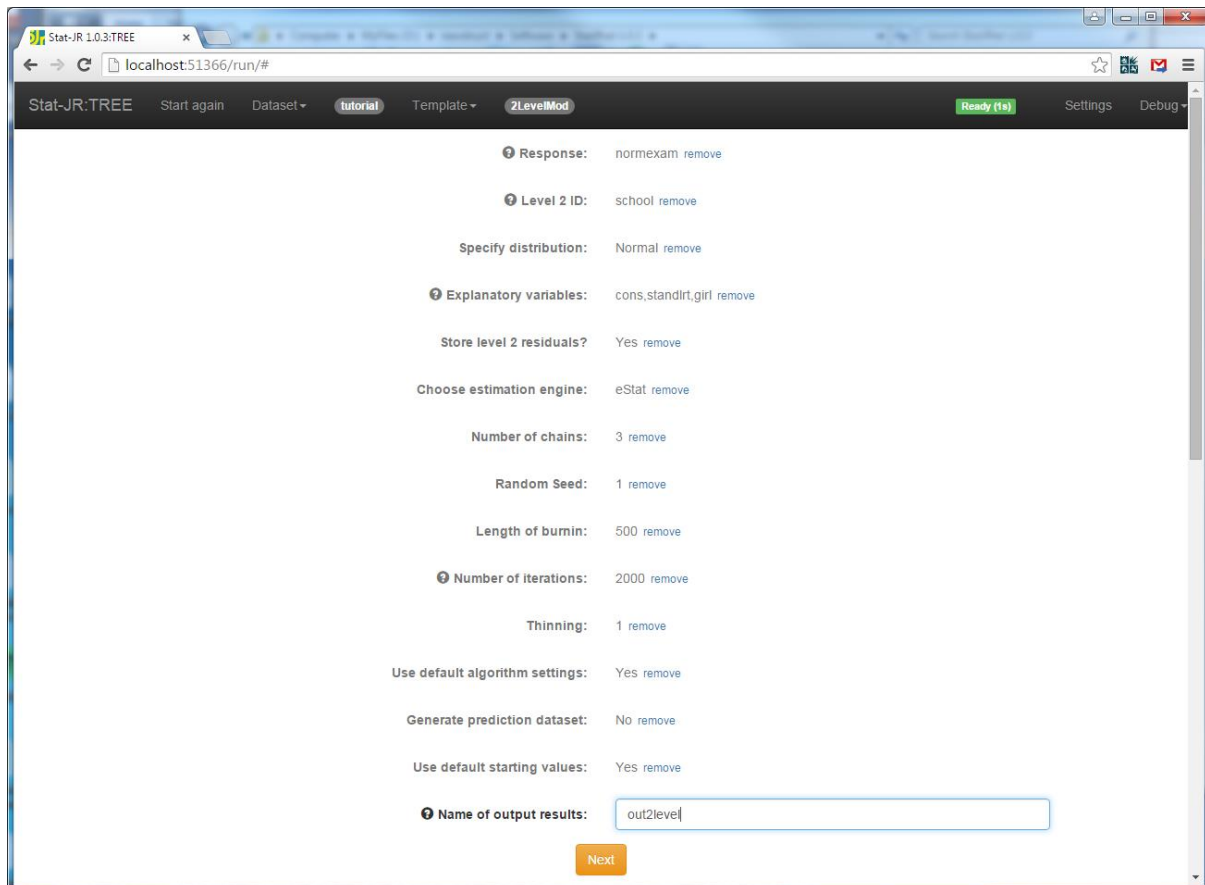
4.2.5 Including school effects

Stat-JR contains many different model-fitting templates some of which can fit whole families of models and some of which can fit just one or two specific models. We have thus far looked at the rather restrictive **Regression1** template that only fits single level normal response models. To include school effects we will now look at the **2LevelMod** template, which not only includes a set of random

effects but also supports different response types and estimation engines, features that we will look at later.

On the **Template** pull down list at the top of the screen select **Choose** and select **2LevelMod** as the template and stick with **tutorial** for the dataset.

Set-up the inputs as shown below:



The screenshot shows the Stat-JR 1.0.3:TREE web interface. The browser address bar shows localhost:51366/run/#. The interface has a dark header with 'Stat-JR: TREE', 'Start again', 'Dataset' (tutorial), 'Template' (2LevelMod), 'Ready (1s)', 'Settings', and 'Debug'. The main content area lists various configuration options for a 2-level model:

- Response: normexam remove
- Level 2 ID: school remove
- Specify distribution: Normal remove
- Explanatory variables: cons, standlrt, girl remove
- Store level 2 residuals?: Yes remove
- Choose estimation engine: eStat remove
- Number of chains: 3 remove
- Random Seed: 1 remove
- Length of burnin: 500 remove
- Number of iterations: 2000 remove
- Thinning: 1 remove
- Use default algorithm settings: Yes remove
- Generate prediction dataset: No remove
- Use default starting values: Yes remove
- Name of output results: out2level

An orange 'Next' button is located at the bottom center of the configuration area.

Press **Next** and then **Run** to fit the model. Note that running will take a while as we are storing all 65 school effects and so for each one the software needs to construct diagnostic plots.

When the model finishes select **ModelResults**, from the output list, and show the results in a separate tab.

parameter	mean	sd	ESS	variable
sigma2_u	0.0927580841793	0.019214800522	3418	
tau	1.77808634602	0.0398110038171	6072	
deviance	9184.86189301	11.9608582571	5978	
beta_0	-0.0909418181226	0.0429473833425	319	cons
beta_1	0.559532031983	0.012593774994	4951	standlrt
beta_2	0.170213502116	0.0329991198223	775	girl
u_0	0.398604325785	0.0921121960575	2286	school
u_1	0.430788328464	0.105398279624	2899	school
u_2	0.518891434178	0.104348689937	2873	school
u_3	0.0376716646194	0.0893074505889	2328	school
u_4	0.241875184368	0.121985255709	3779	school
u_5	0.469549038534	0.0907376260424	2064	school
u_6	0.30512934547	0.0871035400229	1998	school
u_7	-0.0997709439726	0.0825388991247	1862	school
u_8	-0.11362155163	0.12136834076	3965	school
u_9	-0.311431588694	0.106620310371	3132	school
u_10	0.266481255227	0.100348401637	2510	school
u_11	-0.0558801364388	0.108899186418	3065	school
u_12	-0.155371453148	0.096835512947	2894	school
u_13	-0.161928176094	0.0650231057667	948	school

Here if you scroll down we see that the DIC value for the two-level model is 9245, compared with 9725 for the simpler model, showing that it is important to account for the two levels in the data. If you scroll down to the beta fixed effect parameters, as shown in the table below, you will find that their mean estimates have changed little.

Parameter	Single level Mean(sd)	Single level ESS	2level Mean(sd)	2level ESS
beta_0	-0.103 (0.0196)	1615	-0.091 (0.0429)	319
beta_1	0.590 (0.0126)	5488	0.560 (0.0126)	4951
beta_2	0.170 (0.0254)	1623	0.170 (0.0330)	775

The standard deviations for *beta_0* and *beta_2* have increased due to taking account of the clustering, and the ESS values have reduced due to correlation in estimating the fixed effects and level 2 residuals.

4.2.6 Caterpillar plot

The random effects in the 2-level model are also interesting to look at, and one graph that is often used is a caterpillar plot. This can be produced in Stat-JR using a template specifically designed for producing this plot. This template requires the user to select all the 'u's to be displayed in the plot, which can be time-consuming if there are many of them:

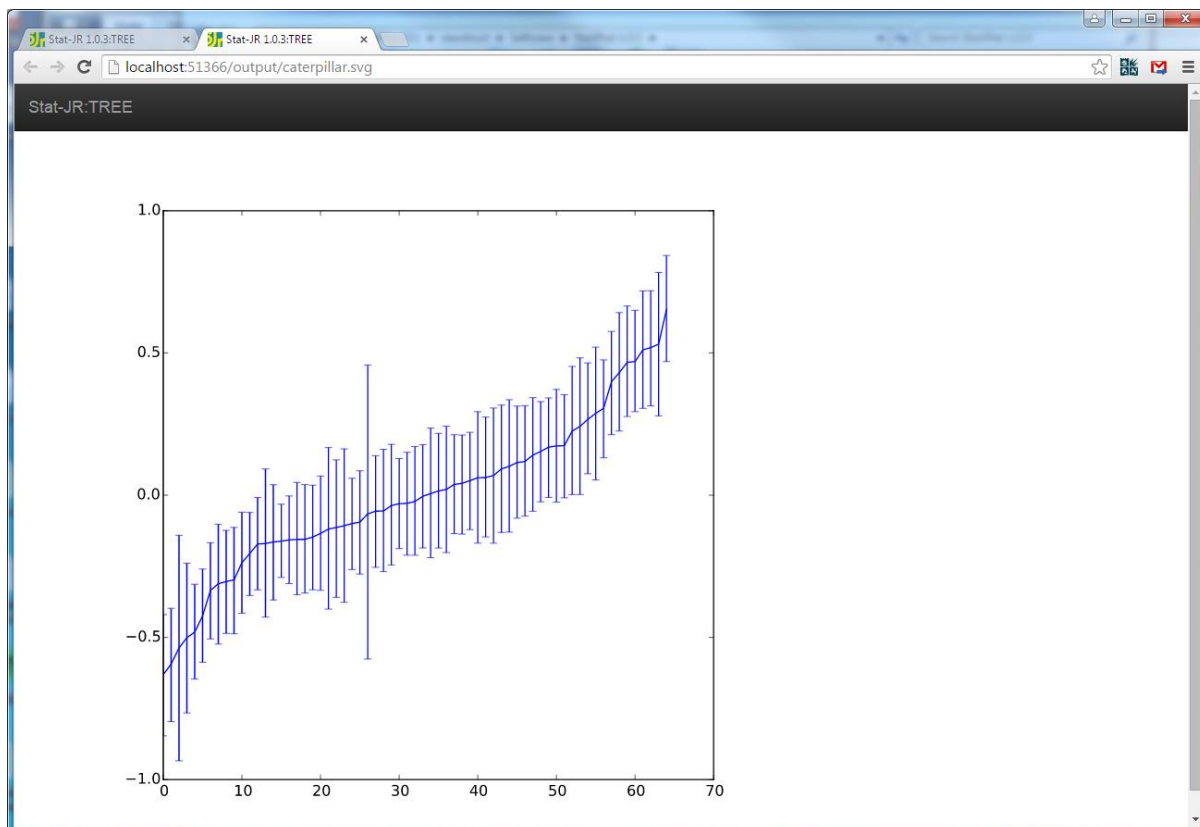
From the top bar we need to select **Choose for Template and Dataset**.

Choose **CaterpillarPlot95** as the template and **out2level** as the dataset.

You need now to select all the u 's from $u0$ to $u64$ which is best done by clicking on $u0$ and holding down the mouse and scrolling down to multiselect all the u 's together

Once all are selected press the **Next** and **Run** buttons.

Select **caterpillar.svg** in the pull down list and view in a new tab as follows:



This graph shows the schools in order of ascending mean whilst the bars give a 95% confidence interval around each mean. The school in the middle with the wide confidence interval (i.e. very large bars) has only 2 pupils and so there is much greater uncertainty in the estimate.

In this chapter we have explored fitting three models to the tutorial dataset. This has illustrated how the Stat-JR system works, how to interpret the output from MCMC and eStat, and how to compare models via the DIC diagnostic. There are better models that can be fitted to the dataset: for example, we could look at treating the effect of intake score (*standlrt*) as random, and fit a random slopes model using the template **2LevelRS**; in the future we may add material on this subject to this manual, but for now we leave this as an exercise for the reader. Next we turn to the interoperability features of Stat-JR.

4.3 Interoperability – a brief introduction

In this section we look at interoperability with other software packages. In order to run this section you will need to have installed the other packages and told Stat-JR where they are. For more details look at the Stat-JR website (www.bristol.ac.uk/cmm/software/statjr/).

4.3.1 So why are we offering Interoperability?

There are many motivations that could be given for the benefits of having an interoperability interface. First and foremost it opens up functionality in other software packages through a common interface.

One important feature that the template, **Regression1AML**, which we cover at the end of this chapter, shows is that not all model templates need to use the built-in eStat engine. It would be perfectly reasonable for a user to construct a template that fitted a specific family of models in the WinBUGS software and then novice users would have access to a user-friendly interface to such models without having to understand the subtleties of writing WinBUGS code; it can thus play an important role introducing packages, such as WinBUGS, to new users. This follows earlier work: for example the MLwiN-WinBUGS interface that we developed 10 years ago.

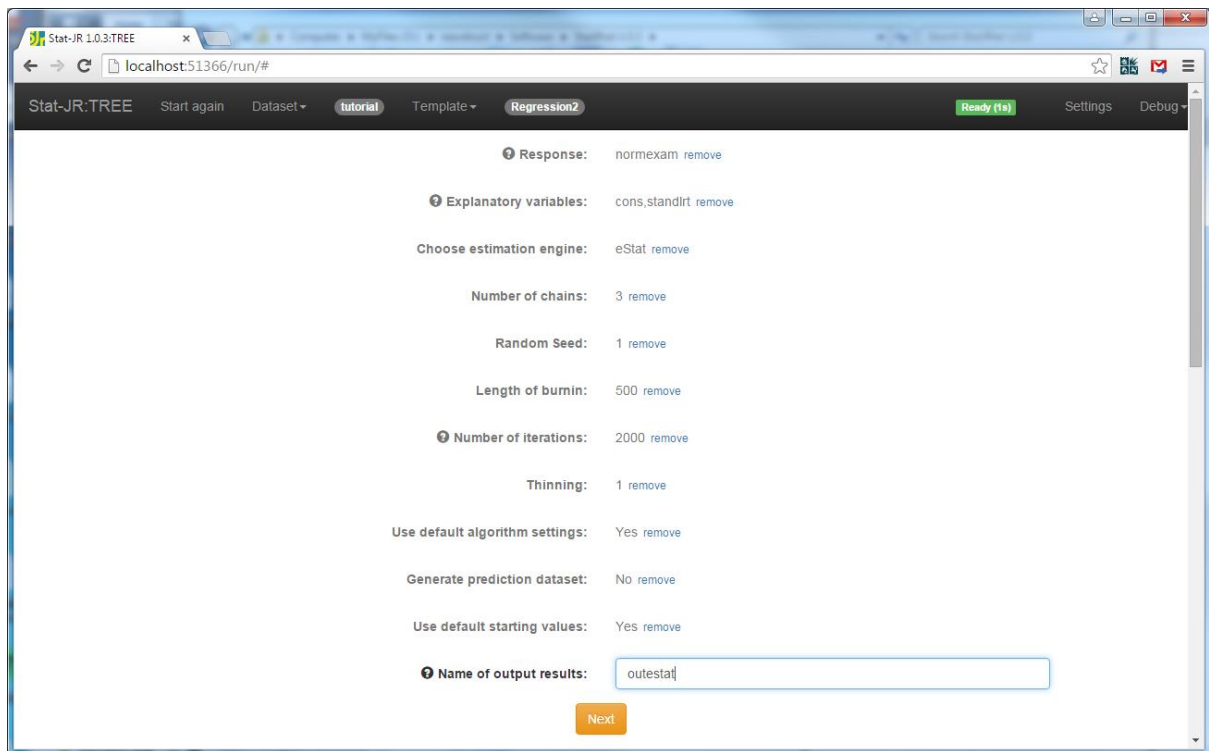
It also offers an easy way of comparing different software packages for a multitude of examples, and we will return to this in Section 4.4.4. Finally it can be thought of as a teaching tool, so that a user familiar with one package can use Stat-JR and directly compare the script files, etc., required for the package with which they are familiar to those required for an alternative package.

4.3.2 Regression in eStat revisited

In Section 4.2 we looked at fitting a few models to the **tutorial** dataset using the built-in eStat engine: a newly-developed estimation engine with the advantage of being transparent in that all the algebra, and even the program code, is available for inspection. It is an MCMC-based estimation method, but is also rather quick. In this chapter we will stick with one simple example, the initial linear regression model that we fitted to the ‘tutorial’ dataset that we considered in Section 4.2. We will need to use a new template, **Regression2**, as the **Regression1** template only supports the eStat engine.

We will begin by setting-up the model and running it in eStat:

From the top bar select **Regression2** as the template, and **tutorial** as the dataset using the **Choose** options on the pull down lists for templates and datasets and set-up the inputs as follows:



Click on **Next** and **Run** to fit the model.

Select **ModelResults** from the pull down list, and show this output in a new tab which should look as follows:

Stat-JR: TREE

Results

Parameters:

parameter	mean	sd	ESS	variable
tau	1.54160995074	0.0340065114631	5799	
beta_0	-0.00127835184871	0.0125770014327	5960	cons
beta_1	0.594959154334	0.012745358164	6129	standlrt
sigma2	0.648987956705	0.0143068971085	5784	
sigma	0.805548947358	0.00887975878981	5789	
deviance	9763.48848832	2.43302399601	6061	

Model:

Statistic	Value
Dbar	9763.48848832
D(thetabar)	9760.50978897
pD	2.97869934714
DIC	9766.46718766

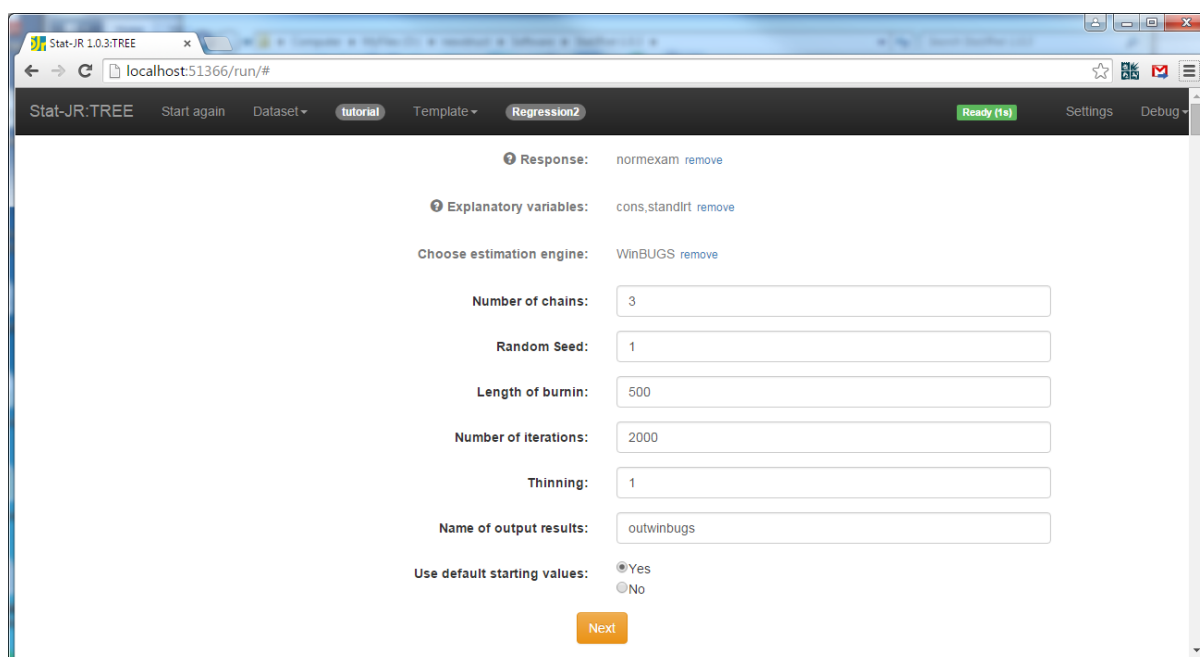
These results are identical to those we obtained using **Regression1** earlier, although we only looked at the plot for *beta_1* in Section 4.2.3. We will use this as a benchmark, contrasting these results with those we obtain from the other packages, although it is worth noting that all packages will have

good mixing and converge quickly for this simple linear regression model. You might like to explore differences between engines / packages for other models yourself after reading this chapter.

4.3.3 Interoperability with WinBUGS

WinBUGS (Lunn et al., 2000) is an MCMC-based package developed (as BUGS – Bayesian inference Using Gibbs Sampling) originally in the early 1990s by a team of researchers at the MRC Biostatistics Unit in Cambridge. It is a very flexible package and can fit, in a Bayesian framework, most statistical models, provided you can describe them in its model specification language. In Stat-JR we have borrowed much of this language for our own algebra system, and so many templates feature interoperability with WinBUGS.

To fit the current model using WinBUGS we can click on remove next to the **Choose estimation engine** input and set up the template inputs as follows:

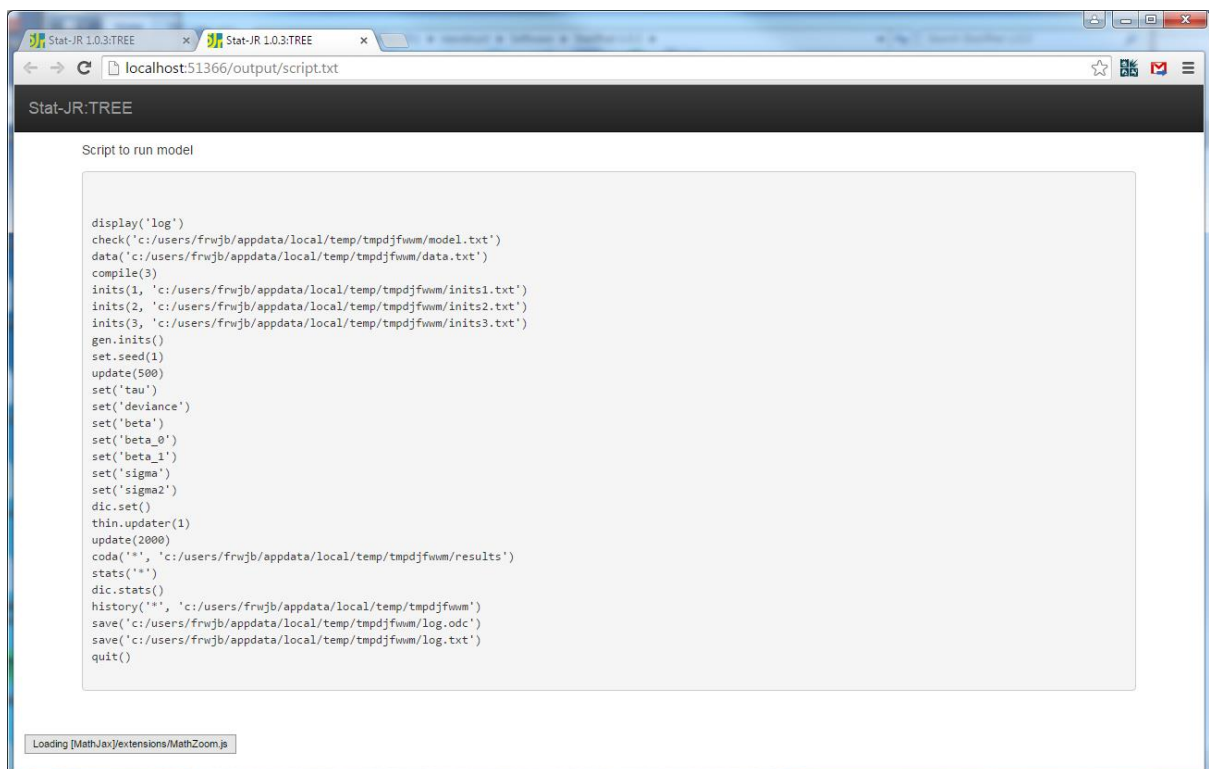


The screenshot shows the Stat-JR web interface in a browser window. The browser address bar shows 'localhost:51366/run/#'. The page title is 'Stat-JR: TREE'. The navigation bar includes 'Start again', 'Dataset', 'tutorial', 'Template', and 'Regression2'. A green 'Ready (1s)' indicator is visible. The main content area is titled 'Choose estimation engine: WinBUGS remove'. Below this, there are several input fields: 'Number of chains: 3', 'Random Seed: 1', 'Length of burnin: 500', 'Number of iterations: 2000', 'Thinning: 1', and 'Name of output results: outwinbugs'. At the bottom, there is a radio button for 'Use default starting values: Yes' (selected) and 'No'. An orange 'Next' button is located at the bottom center.

When we press **Next** the Stat-JR software will construct all the files required to run WinBUGS so for example we can choose *model.txt* from the list:



Here we see the model defined in the WinBUGS model specification language in the output pane. This file is almost identical to that used by eStat aside from the expression *length(normexam)* being replaced here by its value 4059. Selecting *script.txt* from the list and popping out to a new tab gives the following:



Here we see a list of the commands to be run in the WinBUGS language to fit the model. Note that this is done using a temporary directory and so this pathname appears in many commands.

Return to the tab containing the main page and click on the **Run** button.

The WinBUGS package then pops up in its own window, runs the above script, and returns control to Stat-JR when it has finished estimating the model. If we look at the *ModelResults* output from the list and pop it out to its own tab we will see the following:

The screenshot shows a web browser window with the URL `localhost:51366/output/ModelResults`. The page content is as follows:

Results

Parameters:

parameter	mean	sd	ESS
beta_0	-0.0010441346655	0.0126334002778	5728
beta_1	0.5947166	0.0127051159423	6665
deviance	9763.501	2.46481892506	6146
sigma	0.8055883833333	0.00890760237399	5758
sigma2	0.6490518333333	0.0143578582423	5761
tau	1.5414645	0.0340784693477	5748

Model:

Statistic	Value
Dbar_normexam	9763.5
Dhat_normexam	9760.51
pD_normexam	2.986
DIC_normexam	9766.48
Dbar_total	9763.5
Dhat_total	9760.51
pD_total	2.986
DIC_total	9766.48

These estimates, as one might expect, are very close to those from eStat, and again all ESS values are around 5,000-6,000. We can also look at the log file from WinBUGS:

Return to the template tab and choose *log.txt* in the outputs list.

Scroll the *log.txt* file down to the bottom, and the screen should look as follows:

```

inits(3,c:/users/frwjb/appdata/local/temp/tmpdjfwmm/inits3.txt)
model is initialized
gen.inits()
command #Bugs:gen.inits cannot be executed (is greyed out)
set.seed(1)
update(500)
set(tau)
set(deviance)
set(beta)
command #Bugs:set cannot be executed (is greyed out)
set(beta_0)
set(beta_1)
set(sigma)
set(sigma2)
dic.set()
thin.updater(1)
update(2000)
coda("c:/users/frwjb/appdata/local/temp/tmpdjfwmm/results")
stats(*)

Node statistics
  node   mean   sd   MC error   2.5%   median   97.5%   start   sample
beta_0 -0.001044 0.01263 1.641E-4 -0.02522 -0.001077 0.02372 501 6000
beta_1  0.5947  0.01271 1.472E-4  0.57  0.5946  0.6196  501 6000
deviance 9763.0 2.451 0.03236 9761.0 9763.0 9770.0 501 6000
sigma  0.8056  0.008908 1.00E-4 0.7885 0.8054 0.8232 501 6000
sigma2 0.6491  0.01436 1.74E-4 0.6218 0.6487 0.6777 501 6000
tau    1.541  0.03407 4.138E-4 1.476 1.542 1.608 501 6000

dic.stats()

DIC
Dbar = post.mean of -2logL; Dhat = -2LogL at post.mean of stochastic nodes
  Dbar  Dhat  pD  DIC
normexam 9763.500 9760.510 2.986 9766.480
total 9763.500 9760.510 2.986 9766.480
history("c:/users/frwjb/appdata/local/temp/tmpdjfwmm")

History

save(c:/users/frwjb/appdata/local/temp/tmpdjfwmm/log.odc)
save(c:/users/frwjb/appdata/local/temp/tmpdjfwmm/log.txt)

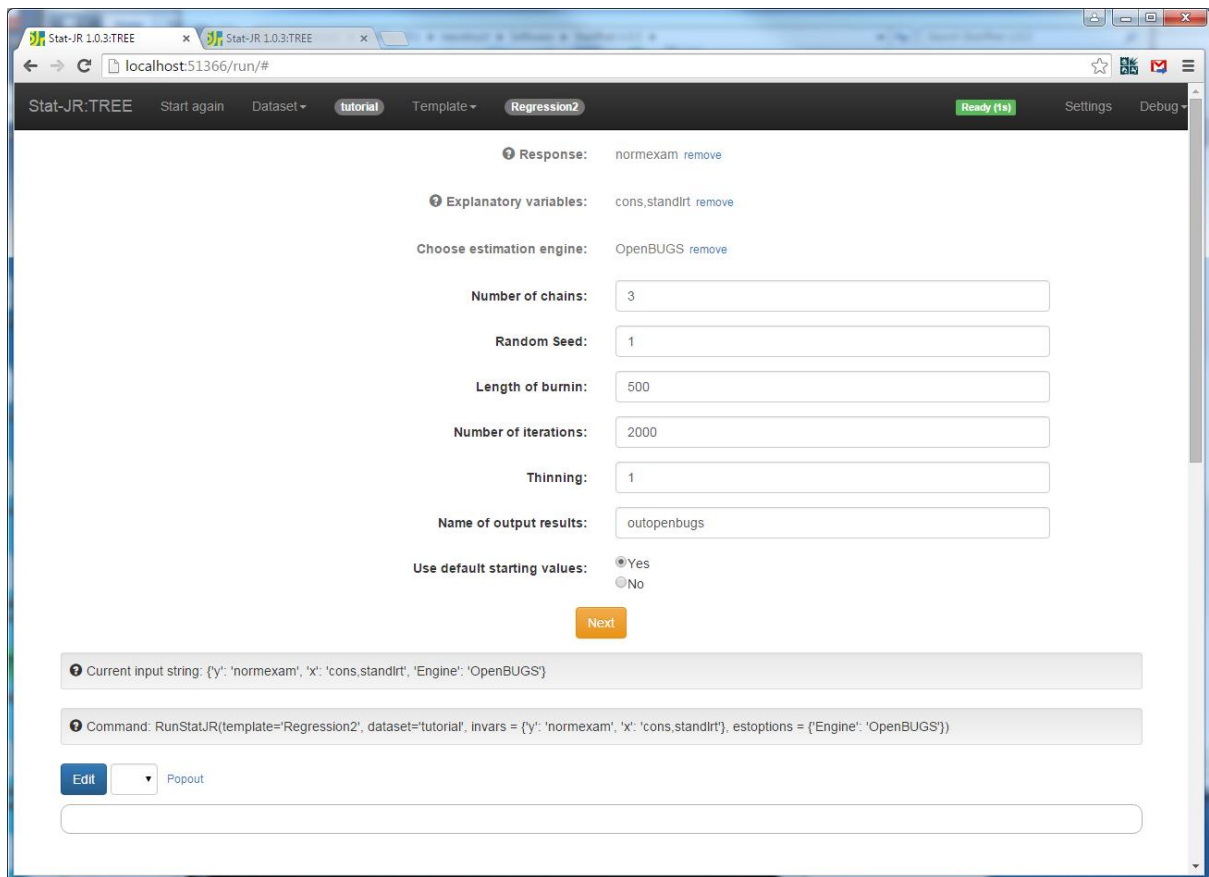
```

Here we see that the estimates and the DIC diagnostic are embedded in the log file, and take a similar value to eStat. WinBUGS required initial value files for each run (and these are stored in 3 text files beginning with *inits* and the chain number), together with a data file as well as the model and script files already shown. All of these are available to view and to use again, thus Stat-JR is useful for learning how these other packages, such as WinBUGS, work.

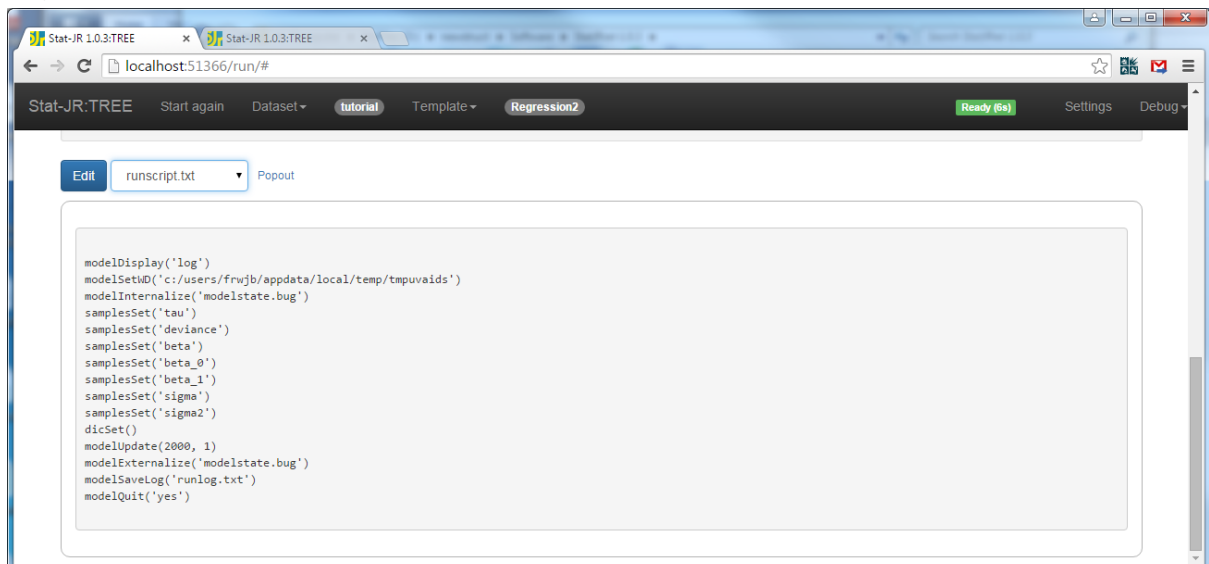
4.3.4 Interoperability with OpenBUGS

Our next package to consider is OpenBUGS (Lunn et al., 2009). OpenBUGS was developed by members of the same team who developed WinBUGS, but differs in that it is open source so other coders may get access to the source code, and in theory develop new features in the software.

To run OpenBUGS via Stat-JR click on the word **remove** next to the **Choose Estimation engine** input, set up the template as follows, and then click on **Next** :



This will have set-up the files required for OpenBUGS; these are similar, but not identical, to WinBUGS: the script file, in particular, is somewhat different and is split into 3 parts called *initscript.txt*, *runscript.txt* (shown below) and *resultsscript.txt*, (you can access this from the objects list):



OpenBUGS allows us to change the working directory, and so there is no need for other commands to include the temporary directory path. Unlike WinBUGS, OpenBUGS will run in the background, and so will not appear when we click run.

Clicking on **Run** and selecting **ModelResults** in its own tab gives the following:

Stat-JR: TREE

Results

Parameters:

parameter	mean	sd	ESS
beta_0	-0.001294676807	0.0126309741479	6018
beta_1	0.5950477	0.0128666030758	5858
deviance	9763.582	2.46413121945	5785
sigma	0.805422516667	0.00916667249695	5954
sigma2	0.648788483333	0.0147680896654	5961
tau	1.54212983333	0.0351074252826	5957

Model:

Statistic	Value
Dbar_normexam	9764.0
Dhat_normexam	9761.0
pD_normexam	3.071
DIC_normexam	9767.0
Dbar_total	9764.0
Dhat_total	9761.0
pD_total	3.071
DIC_total	9767.0

Again, these results are very similar in terms of parameter estimates and ESS values to the other software packages.

4.3.5 Interoperability with JAGS

The third standalone MCMC estimation engine available, via Stat-JR, is JAGS (Just Another Gibbs Sampler), developed by Martyn Plummer (Plummer, 2003). JAGS also uses WinBUGS model language, but has a few differences in terms of script files and data files.

To run JAGS via Stat-JR click on the **remove** text next to **Choose estimation engine** and set-up the template as follows, before clicking on **Next** :

Stat-JR: TREE

Start again Dataset **tutorial** Template **Regression2** Ready (1s) Settings Debug

Response: normexam **remove**

Explanatory variables: cons,standlrt **remove**

Choose estimation engine: JAGS **remove**

Number of chains:

Random Seed:

Length of burnin:

Number of iterations:

Thinning:

Name of output results:

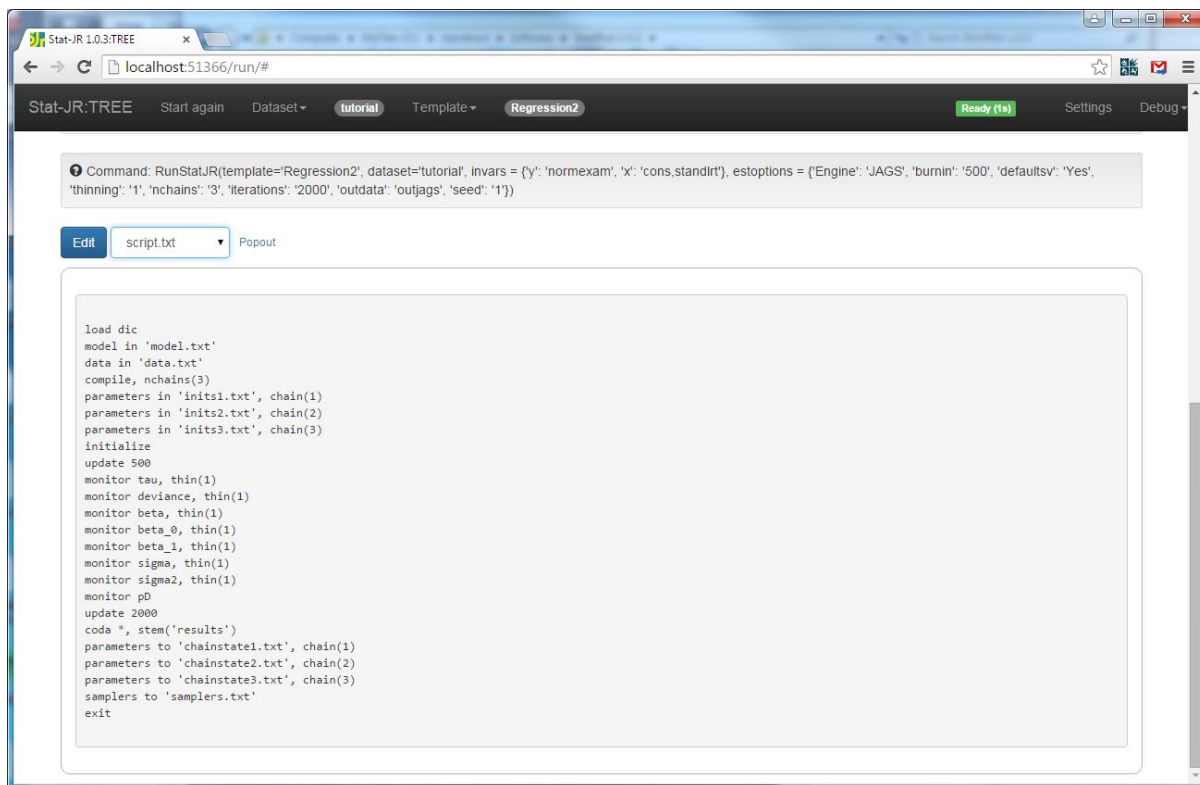
Use default starting values: Yes No

Next

Current input string: {y: 'normexam', x: 'cons,standlrt', Engine: 'JAGS'}

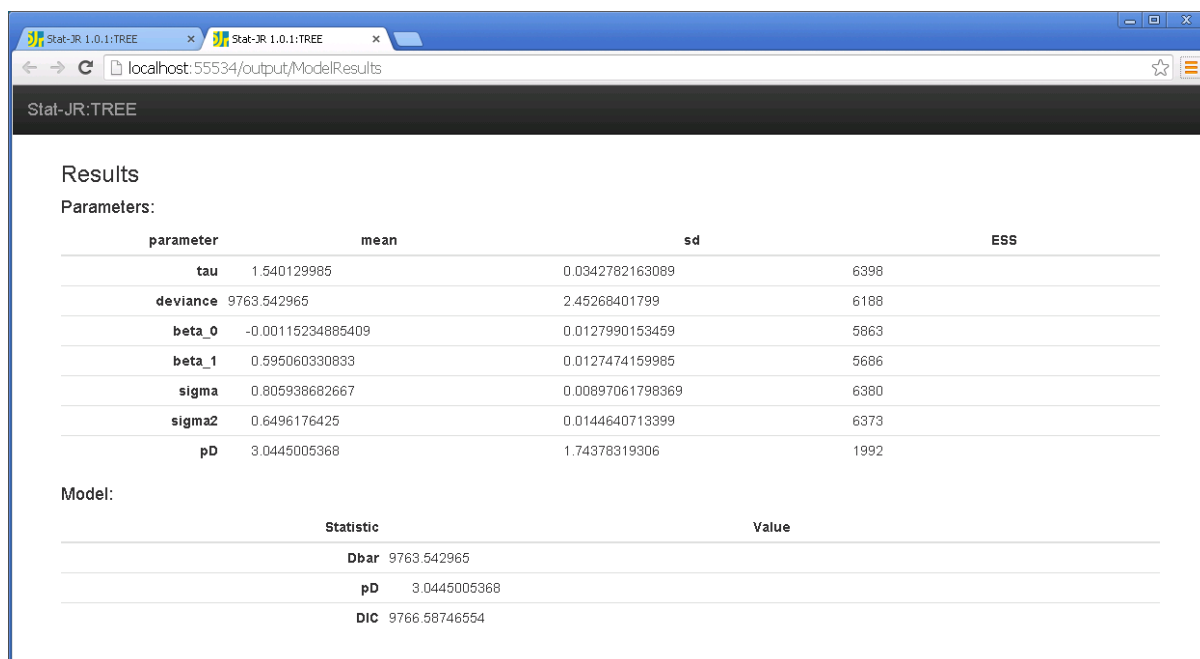
Command: RunStatJR(template='Regression2', dataset='tutorial', invars = {y: 'normexam', x: 'cons,standlrt'}, estoptions = {'Engine': 'JAGS'})

This will set-up the files required for JAGS; for example, here you can see the script file (*script.txt*) which show some differences to those for WinBUGS (as to the initial value file formats):



Like OpenBUGS, JAGS will run in the background (i.e. it will not open as a window on your screen).

Clicking on **Run**, and placing *ModelResults* in a new tab, gives the following:

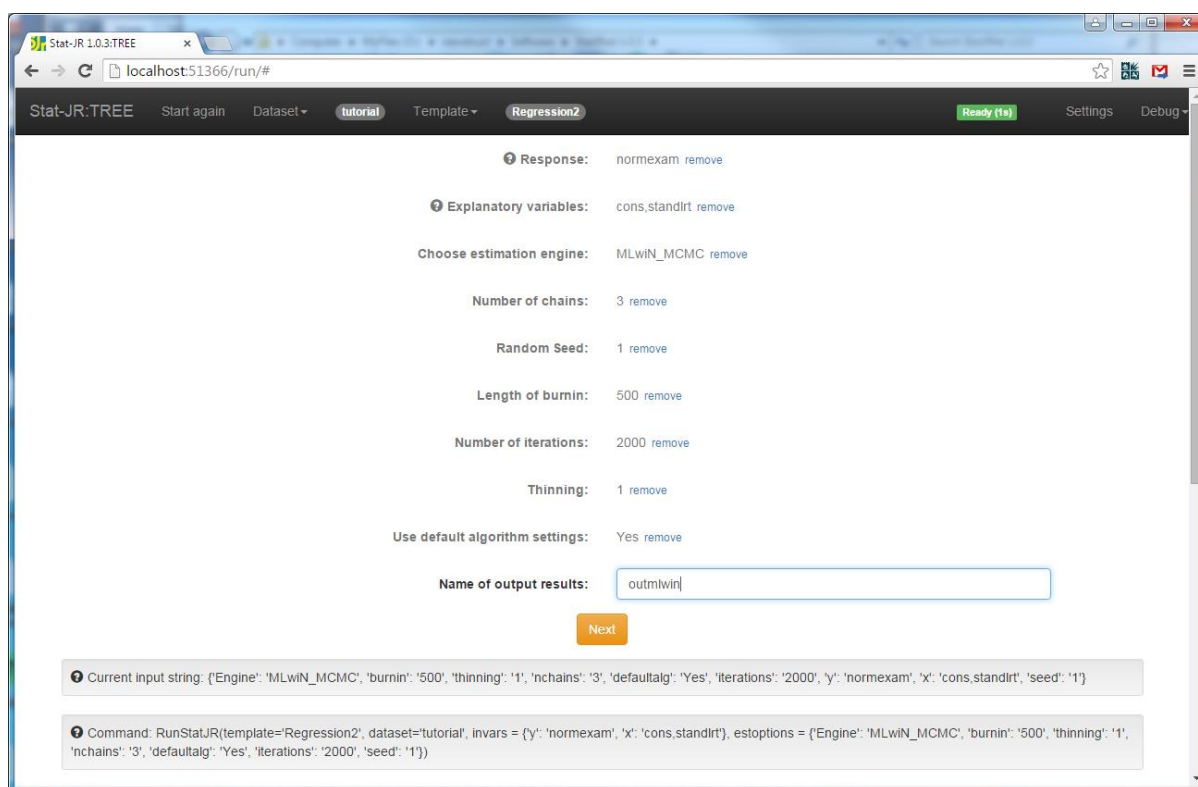


As you can see, we have similar estimates and effective sample sizes to the other estimation methods we've used. Whilst JAGS can be faster than WinBUGS and OpenBUGS, it fits a slightly smaller number of models.

4.3.6 Interoperability with MLwiN

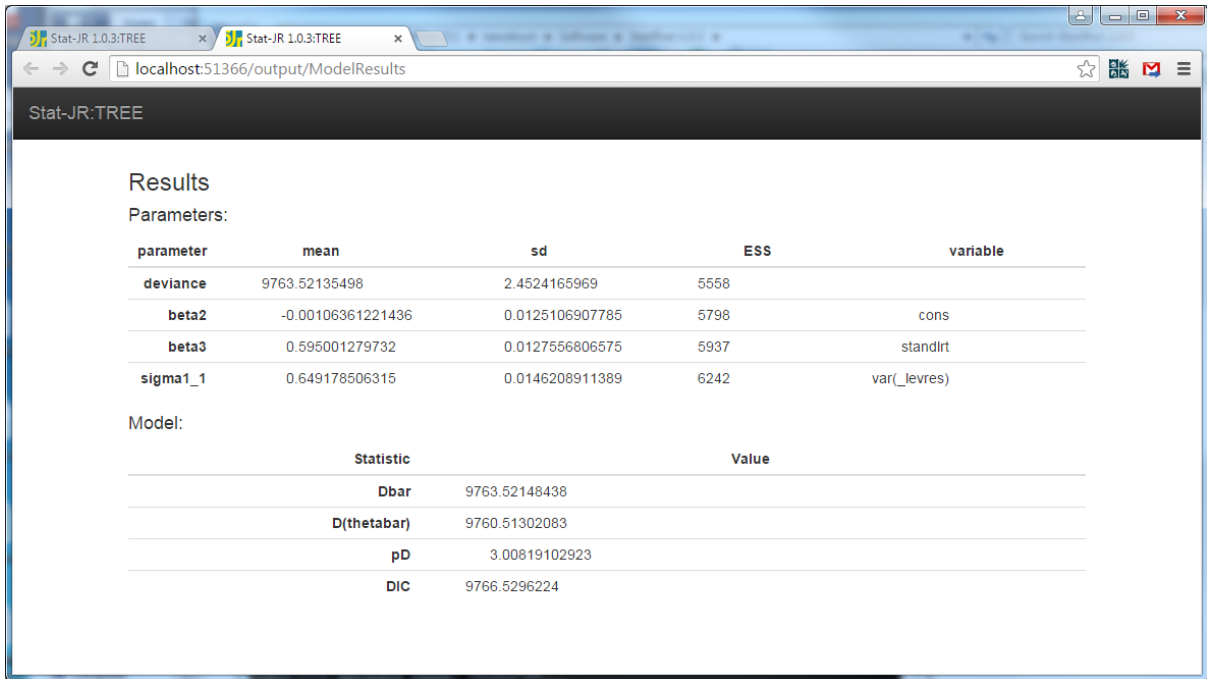
MLwiN (Rasbash et al. 2009) is a software package specifically written to fit multilevel statistical models. It features two estimation engines (for MCMC and likelihood-based (IGLS) methods, respectively) with a menu-driven, point-and-click user interface. It also has an underlying macro language, however, and this is what we use to interoperate with Stat-JR. We will first consider the MCMC engine. As it is limited in the scope of models it fits, this means it is generally quicker than the other MCMC packages. MLwiN is a single chain program, but can be made into a multiple chain engine with Stat-JR, since the latter can start-up three separate instances of MLwiN. At present these are given different random number seeds, but the same starting values, however we will try and change this in future.

To run MCMC in MLwiN, via Stat-JR, click on the **remove** text by **Choose estimation engine** input and set-up the template as follows before clicking on **Next** :

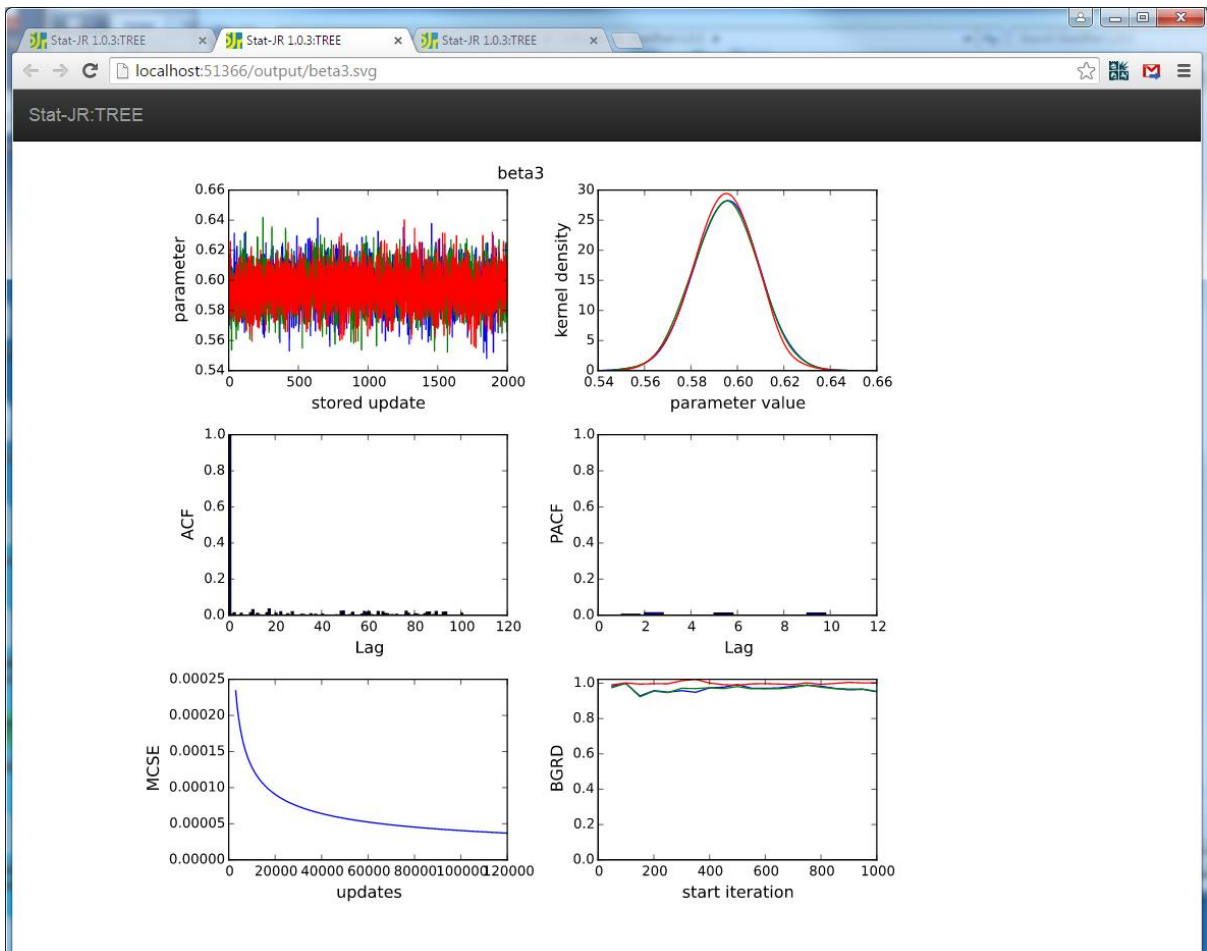


You can see, in the pulldown list the dataset (in .dta format) that is used by MLwiN. There are also several MLwiN script files for the multiple chains and the several stages of model fitting.

Clicking on the **Run** button will set off three instances of MLwiN (in the background) and Stat-JR will then collate the results together. Choosing *ModelResults*, and displaying them in a new tab, gives the following:

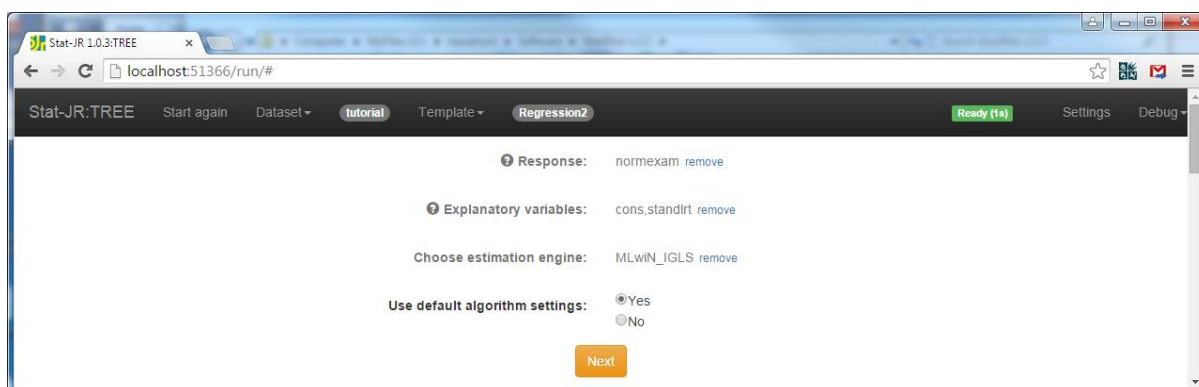


Once again here we have similar estimates, although the naming convention is slightly different for MLwiN. To show that we have multiple chains we can examine the chains for the slope (β_3), as shown below:

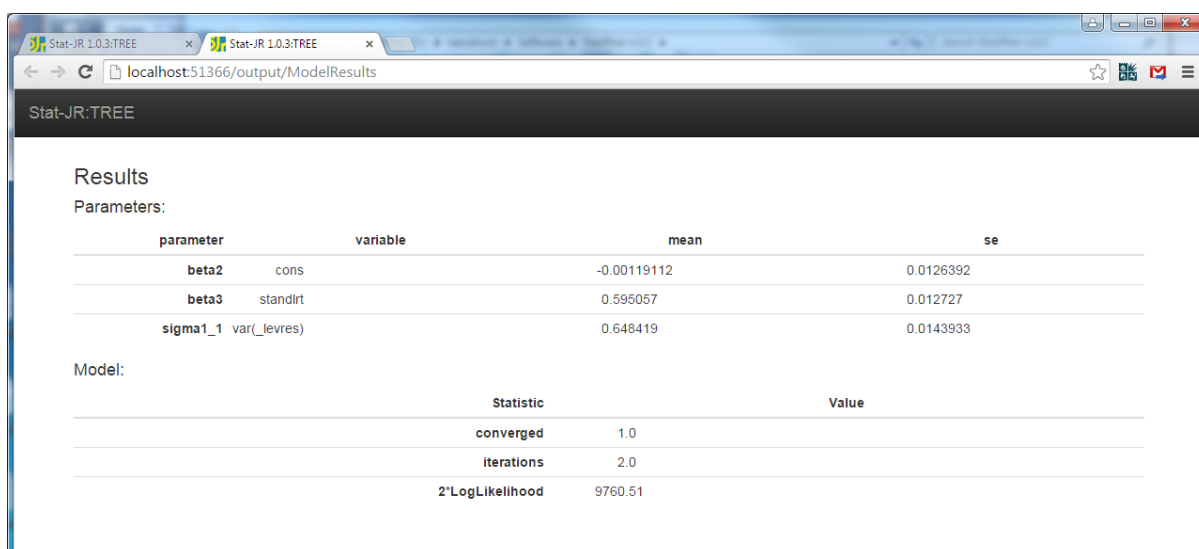


Stat-JR also offers the option of using the likelihood-based IGLS estimation engine in MLwiN.

To do this in MLwiN, via Stat-JR, click once again on the **remove** text next to the **Choose estimation engine** input and set-up the template as follows, before clicking on **Next**:



Again the dataset will appear in the output pane, and this time pressing **Run** will give the following in the **ModelResults** output:

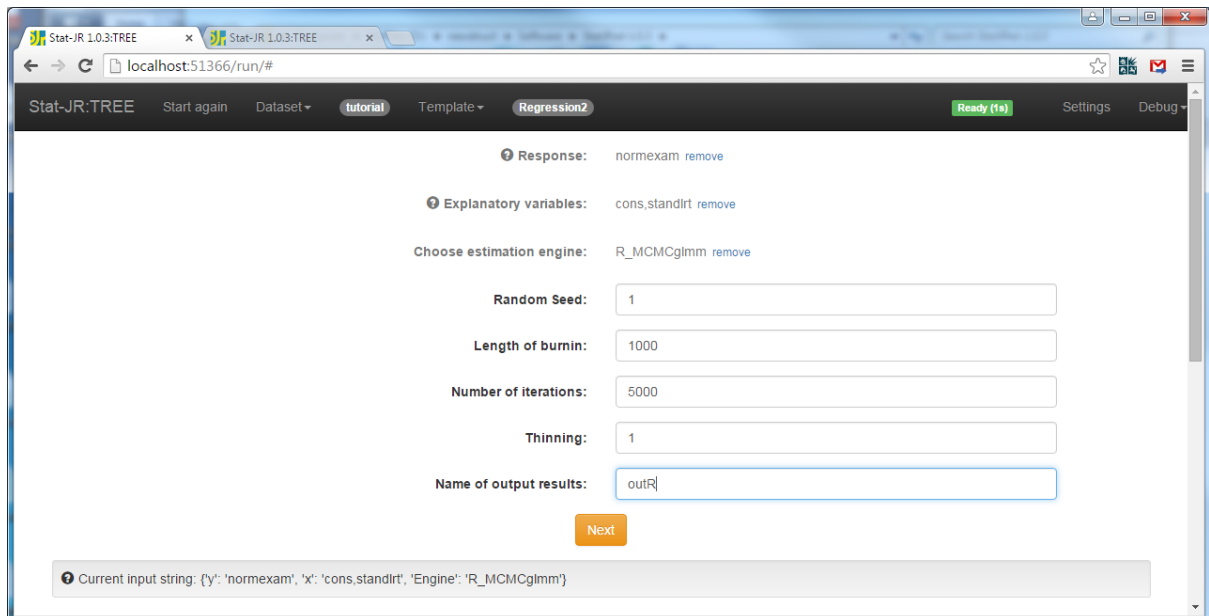


Here we get the *Deviance* ($-2 \times \text{Loglikelihood}$) value, together with parameter estimates with standard errors. The likelihood-based methods are far faster to run than the MCMC-based methods.

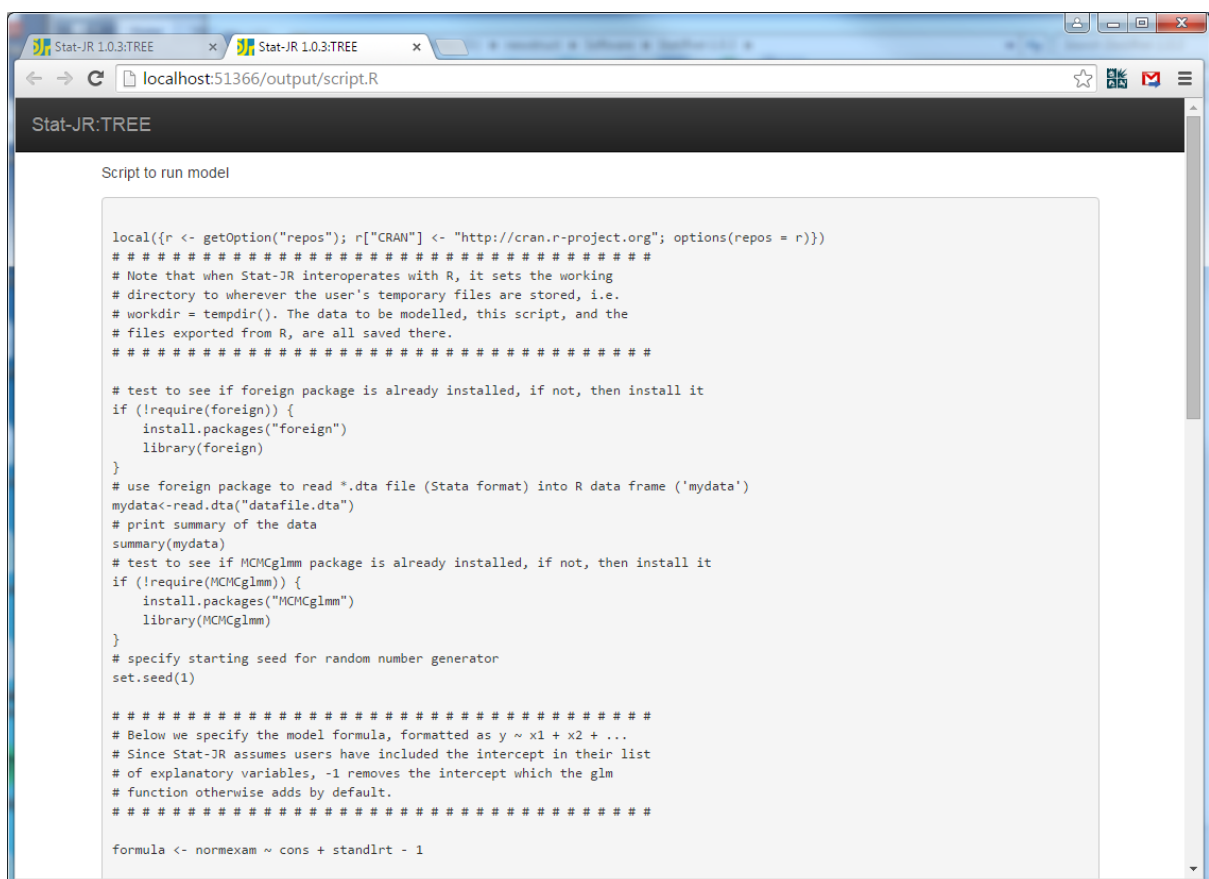
4.3.7 Interoperability with R

R (R Development Core Team, 2011) is another more general purpose package that can be used to fit many statistical models. R has many parallels with Stat-JR in that users can supply functions (like Stat-JR templates) which are then added to the library of R packages. We have thus far implemented interoperability features with R for several of these R functions; for example, for the template **Regression2**, we have implemented two R engines: *R_MCMCglmm*, which is MCMC-based, and *R_glm*, which is a standard regression modelling function. We will firstly demonstrate *MCMCglmm*.

To run MCMC in R, via Stat-JR, click on the remove text by the **Choose estimation engine** input and set-up the template as follows, and click on **Next**:



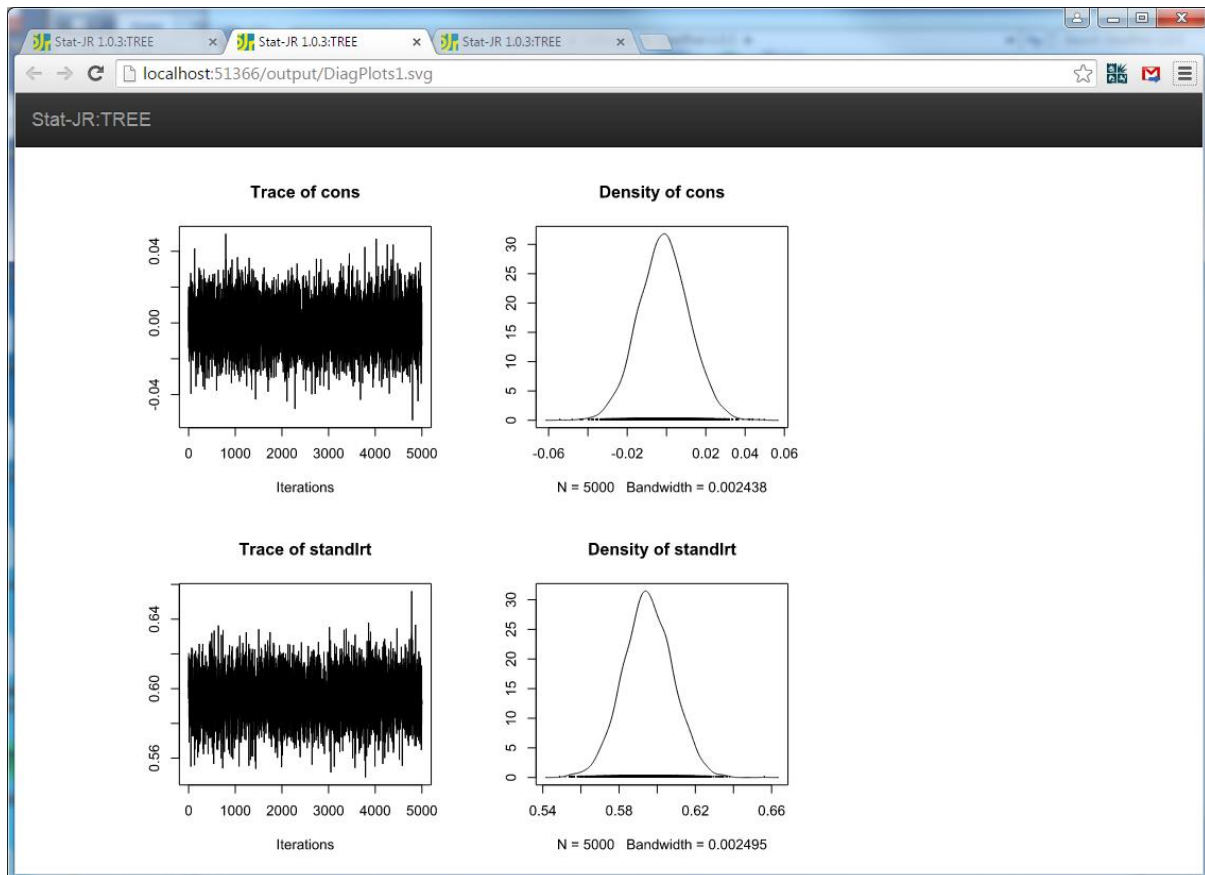
After pressing **Next**, if we look at the script file, *script.R*, which we can select from the outputs list, we see the following:



MCMCglmm can fit all forms of generalised linear mixed models, of which a linear regression is a rather trivial case. You will see that the script file contains some setup code which will actually download and install the *MCMCglmm* library the first time you execute the script (so ensure your machine is connected to the internet) before calling the *MCMCglmm* command and then producing summaries.

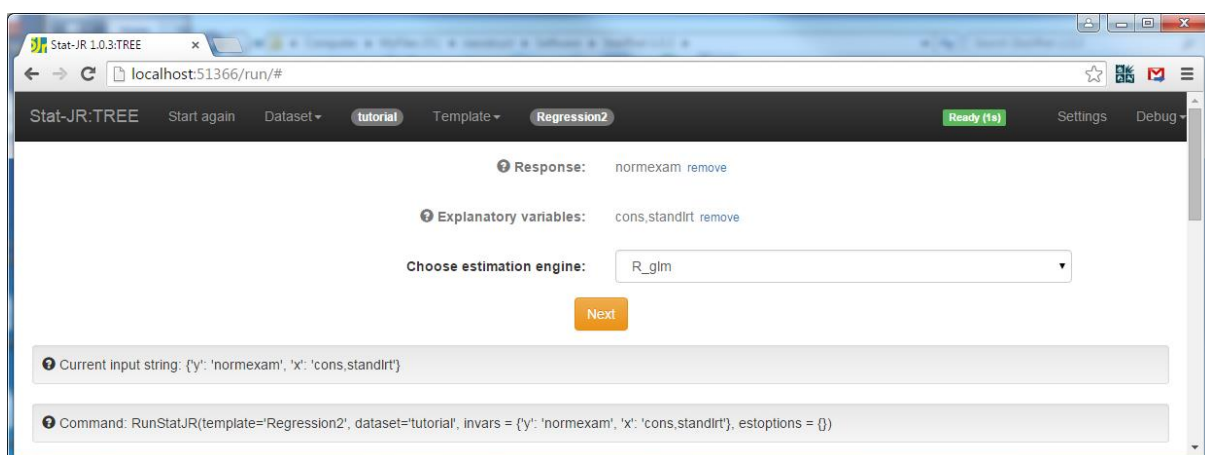
Clicking on **Run** in the main window will create several outputs.

The *ModelResults* are similar to other software but we can also look at diagnostics plots that are specific to R by selecting *DiagPlots1.png*:

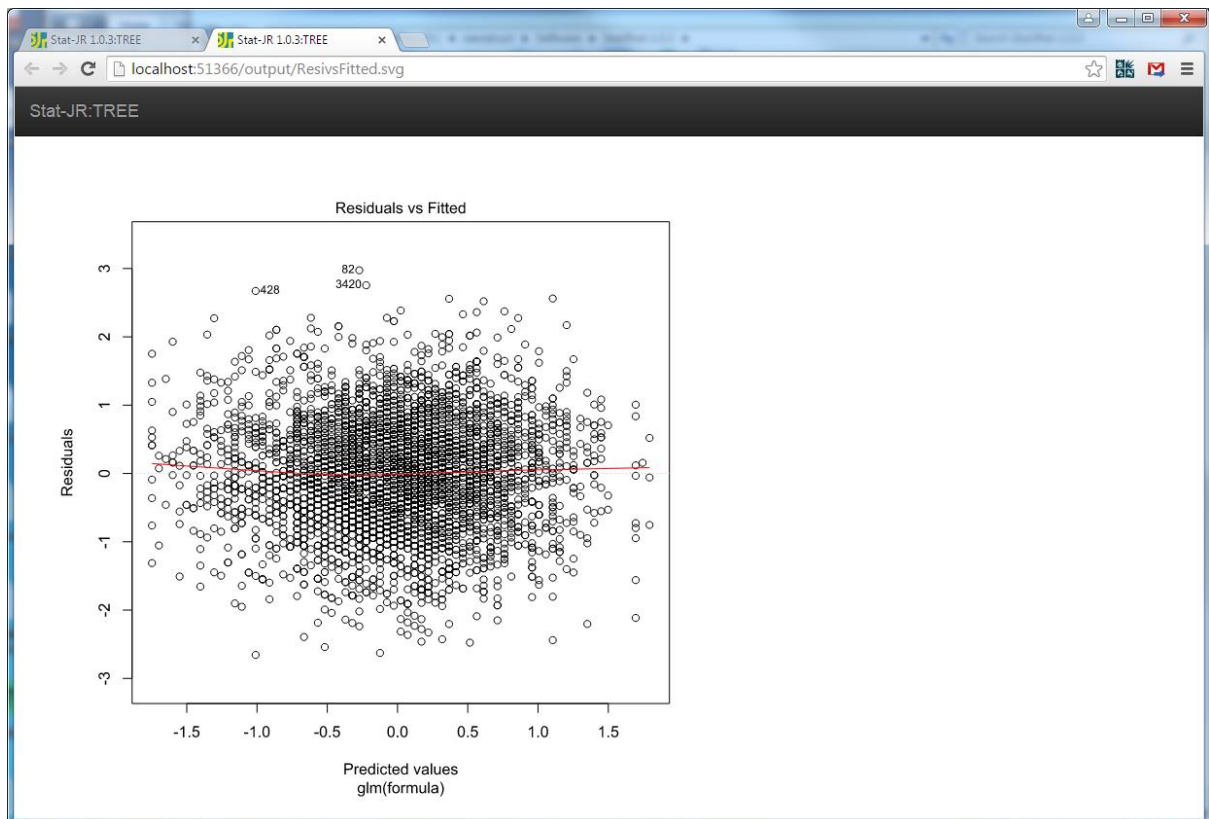


Here R gives trace plots and kernel density plots for both the intercept and the slope parameter.

Turning next to the *glm* package we can click on the **remove** text by **Choose estimation engine** and set-up the template as follows, before clicking on **Next** :



Clicking on **Run** will this time run the *MASS* package and give results in *ModelResults* as usual. There are additional graphical plots that come back from R; for example, below is a plot of residuals of the model fit against fitted values (*ResivsFitted.svg*).



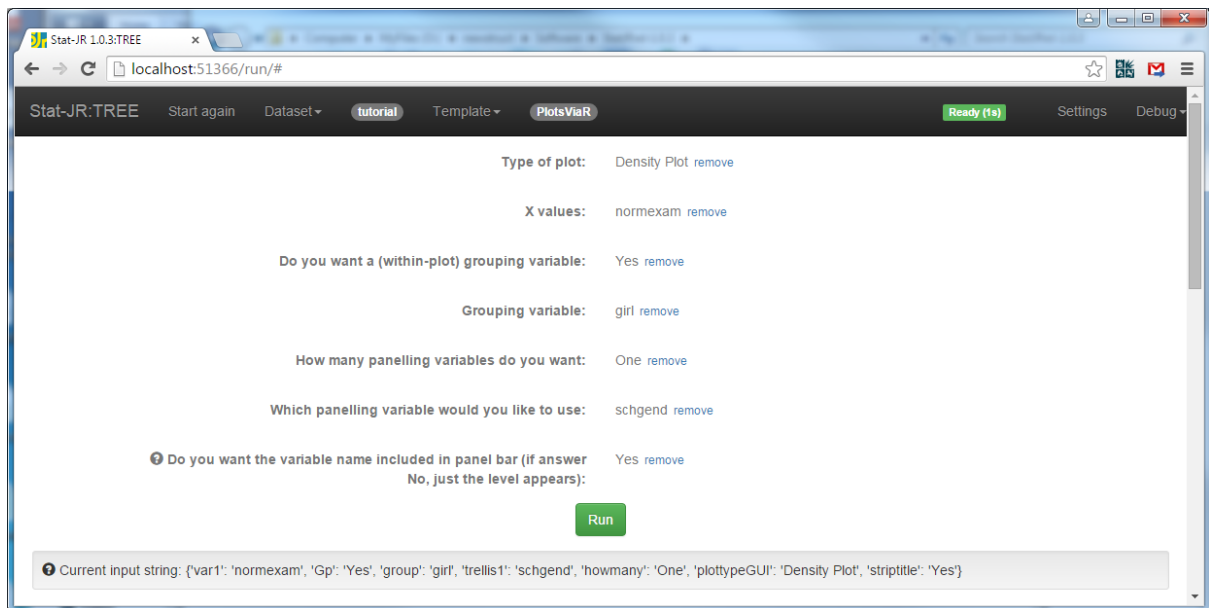
Before finishing with R, we will also demonstrate a non-model template developed with R called **PlotsViaR** that gives the Stat-JR user access to R's *lattice* graphics package through the Stat-JR interface.

Click on **Choose** from the **Template** pull down list at the top of the screen to get a list of all the templates. Note that the search cloud is useful with interoperability as it can be used to show which templates offer interoperability with a particular package (the engines are in red).

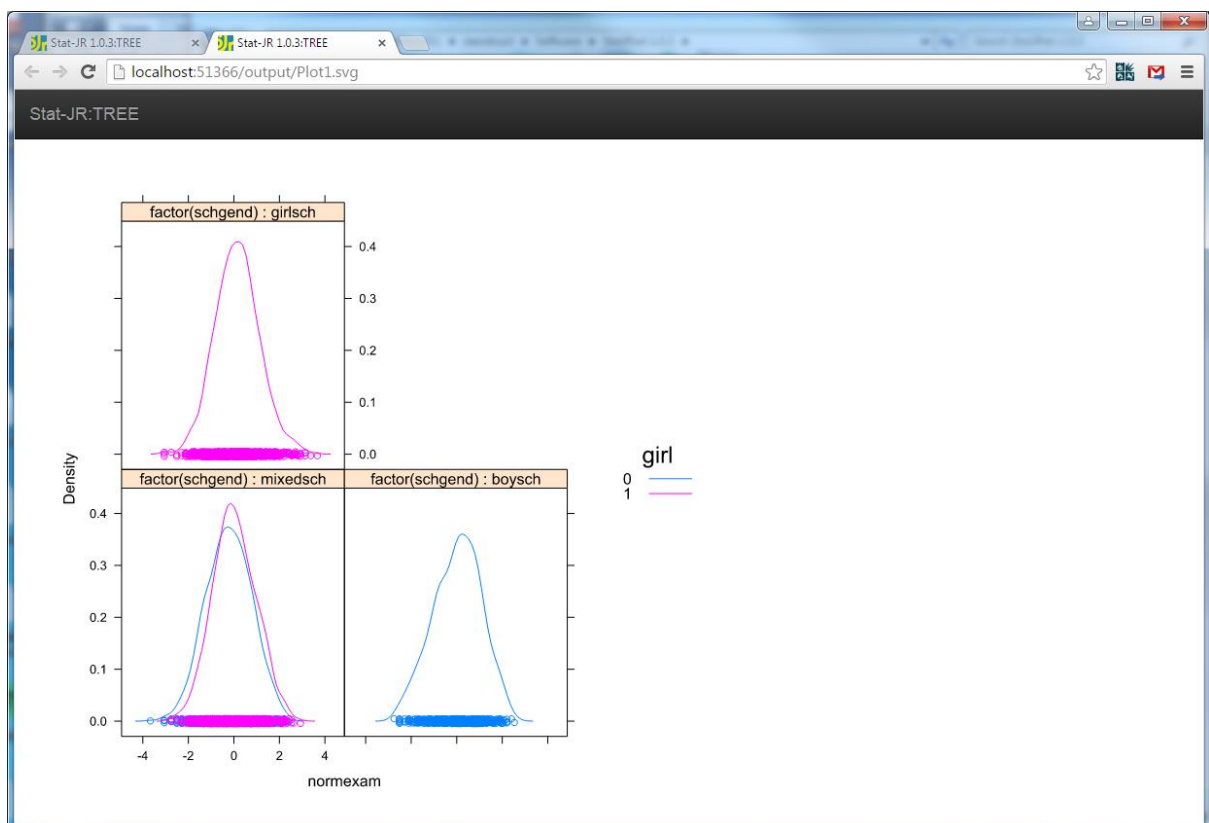
Click on **Plots** and also **R_script** in the blue tag cloud. You'll see that the list of templates, underneath, is accordingly reduced to just those that draw plots using R.

Select **PlotsViaR** from the list, and click **Use**.

Set up the template inputs as shown below:



These options will display kernel plots for the exam scores of pupils grouped by gender, with separate (panelled or trellise) plots for each school gender type. We can now **press Run** and show the plot (*Plot1.svg*) in a separate tab:



Here (by coincidence) we have blue for boys and pink for girls!

4.3.8 Interoperability with AML

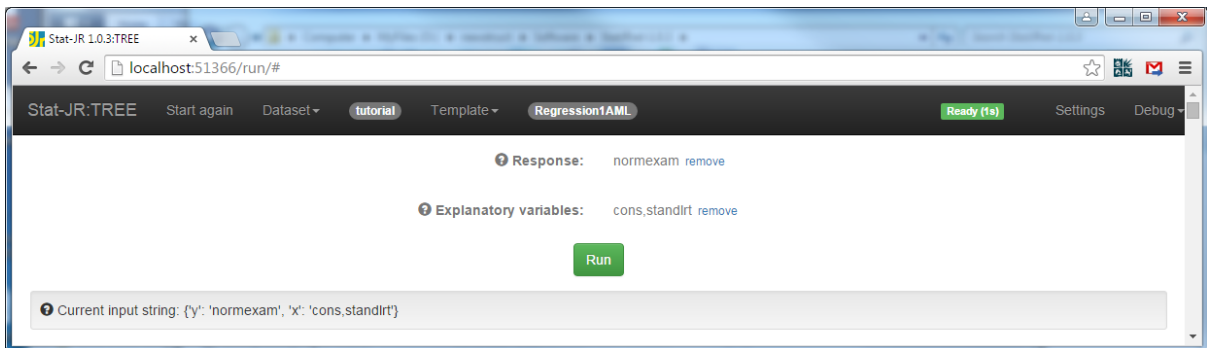
We will next look at another software package that can fit many statistical models via likelihood-based estimation. AML (Lillard & Panis, 2003) is very useful for fitting multi-process models, but as with other software packages can fit a simple regression as a special case. In our development work on Stat-JR we have written special templates for interoperability with AML as opposed to incorporating interoperability in the standard templates. We therefore need to do the following:

Click on the **Choose** option from the Template pull down list.

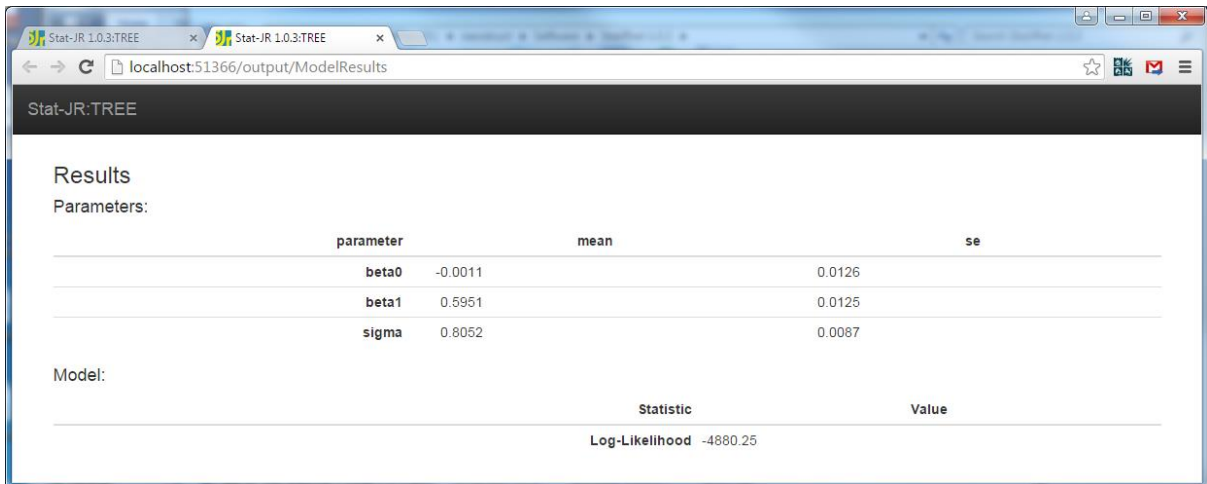
Select **Regression1AML** from the template list and click on **Use**, and stick with the **tutorial** dataset.

Note that if you have earlier clicked on **Plots** and **R_script** in the cloud of terms you will need to either unselect them or click on **[reset]** to see the required template.

Fill in the inputs as follows, and press **Next**:



Now click on **Run** to run the model in AML and select *ModelResults* from the list:



Here we see the model results are similar to other packages. AML has three input dataset (*amlfit.raw*, *amlfit.aml* and *amlfit.r2a*). There are also three additional output files from AML: *amlfit.out*, *amlfit.tab* and *amlfit.sum*. For more information on how AML works we recommend looking at the reference manual for the software.

We also have interoperability support for a variety of other packages, including GenStat, MATLAB, Minitab, Octave, Sabre, SAS, SPSS, Stan (via RStan) and Stata. These packages are either not installed

on the machine we are currently using, or are not supported by the **Regression2** template that is being demonstrated, and so Stat-JR realises this and does not offer them.

4.4 Application 2: Analysis of the Bangladeshi Fertility Survey dataset

4.4.1 The Bangladeshi Fertility Survey dataset

The Bangladeshi dataset (**bang1**) is an example dataset from the 1988 Bangladeshi Fertility Survey. It contains records from 1934 women based in 60 districts in Bangladesh, and we are planning to investigate variables that predict whether the women were using contraception or not at the time of the survey. Let us first look at the data and the variables we will consider.

Select **Choose** and pick **bang1** from the **Dataset** list and click on **Use**. Click on **View** from the **Dataset** list to view the data as follows:

The screenshot shows the Stat-JR interface with the dataset 'bang1' selected. The table below represents the data shown in the interface:

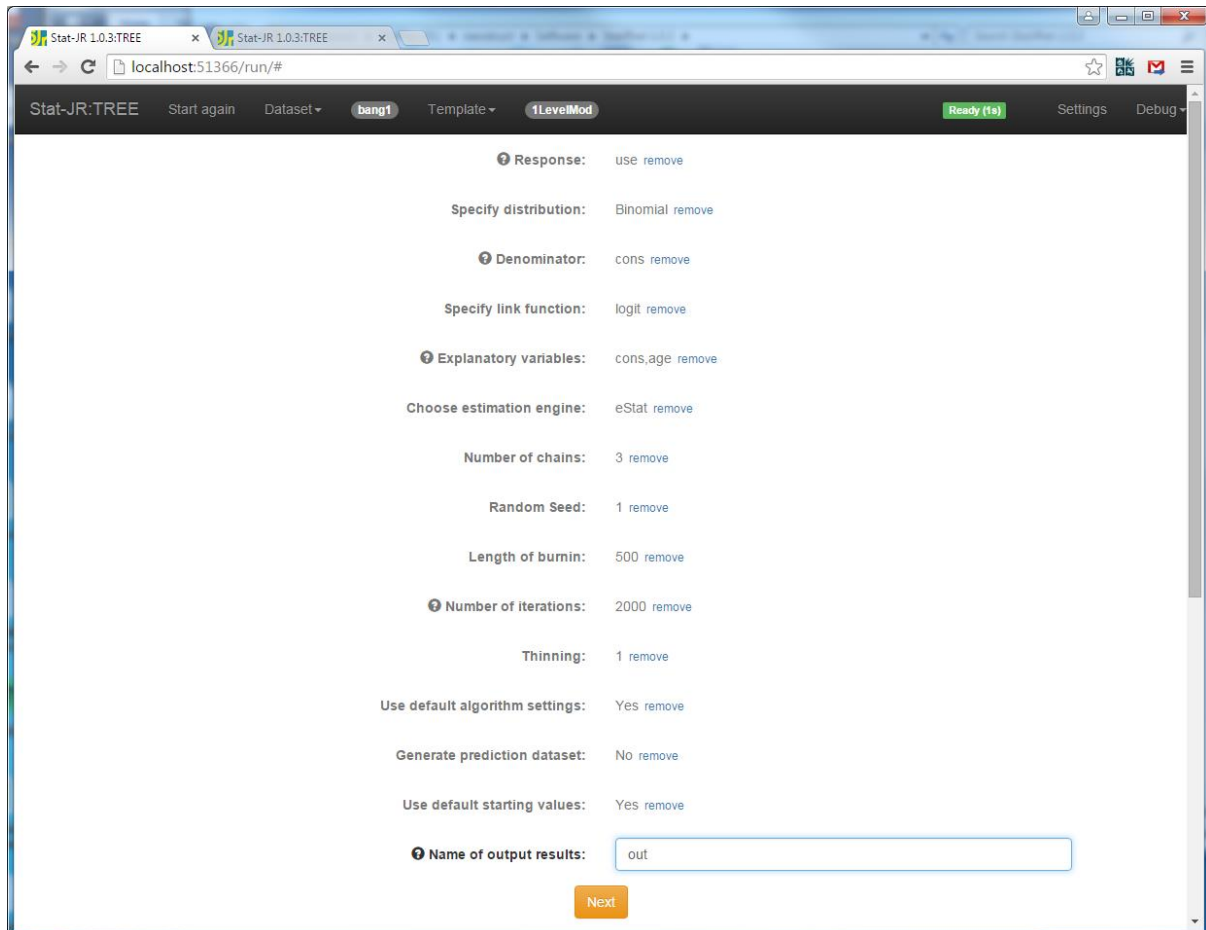
	woman	district	use	lc	age	urban	educ	hindu	d_illit	d_pray	cons
1	1	1	0	three+kids	18.44	1	1	0	0.58	0.64	1
2	2	1	0	nokids	-5.56	1	1	1	0.58	0.64	1
3	3	1	0	twokids	1.44	1	2	0	0.58	0.64	1
4	4	1	0	three+kids	8.44	1	1	0	0.58	0.64	1
5	5	1	0	nokids	-13.56	1	1	0	0.58	0.64	1
6	6	1	0	nokids	-11.56	1	1	0	0.58	0.64	1
7	7	1	0	three+kids	18.44	1	1	0	0.58	0.64	1
8	8	1	0	three+kids	-3.56	1	1	0	0.58	0.64	1
9	9	1	0	onekid	-5.56	1	1	0	0.58	0.64	1
10	10	1	0	three+kids	1.44	1	1	0	0.58	0.64	1
11	11	1	1	nokids	-11.56	1	1	0	0.58	0.64	1
12	12	1	0	nokids	-2.56	1	1	0	0.58	0.64	1
13	13	1	0	onekid	-4.56	1	1	0	0.58	0.64	1
14	14	1	0	three+kids	5.44	1	1	0	0.58	0.64	1
15	15	1	0	three+kids	-0.559999	1	1	0	0.58	0.64	1
16	16	1	1	three+kids	4.44	1	1	0	0.58	0.64	1
17	17	1	0	nokids	-5.56	1	1	0	0.58	0.64	1
18	18	1	1	three+kids	-0.559999	1	2	0	0.58	0.64	1
19	19	1	1	onekid	-6.56	1	4	0	0.58	0.64	1
20	20	1	0	twokids	-3.56	1	1	0	0.58	0.64	1
21	21	1	0	nokids	-4.56	1	3	0	0.58	0.64	1
22	22	1	0	nokids	-9.56	1	1	0	0.58	0.64	1
23	23	1	0	three+kids	2.44	1	2	0	0.58	0.64	1
24	24	1	1	twokids	2.44	1	4	0	0.58	0.64	1
25	25	1	1	onekid	-4.56	1	4	0	0.58	0.64	1
26	26	1	0	three+kids	14.44	1	4	0	0.58	0.64	1
27	27	1	1	nokids	-6.56	1	4	0	0.58	0.64	1

Here we see records for the first 24 women in district 1 displayed. The response variable *use* takes value 1 if the woman was using contraceptives during the time of the survey, and 0 if she was not. There are then several predictor variables, both woman-level and district-level. Here we will focus on just two: the number of living children (*lc*), which is a categorical variable with four categories (no kids, one kid, two kids, three+kids), and the respondents' *age*, which is measured to the nearest year and has been centred around its grand mean. We will now consider modelling the dataset.

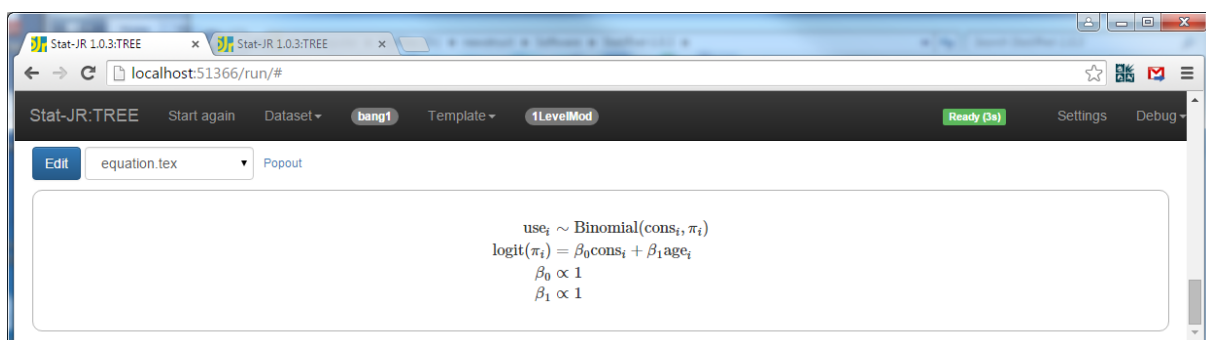
4.4.2 Modelling the data using logistic regression

We will firstly consider a simple linear regression model relating contraception use to the age of the woman.

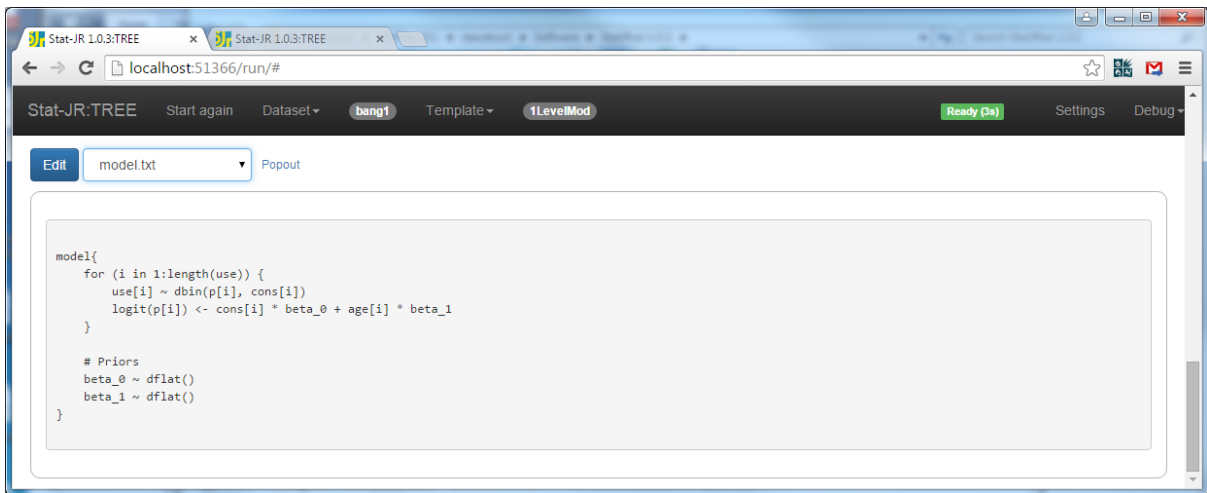
Choose the template **1LevelMod** from the **Template** list and click on **Use**.
Then setup the model with inputs as below.



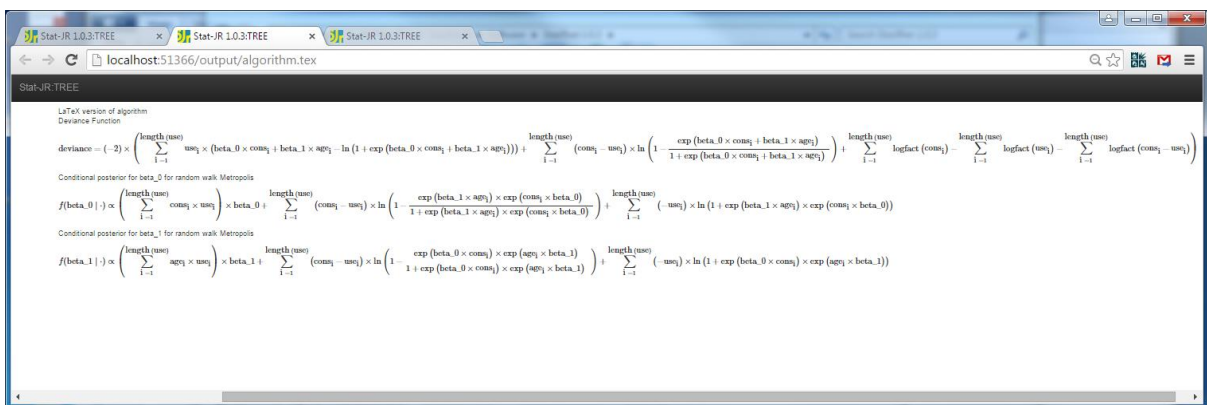
Clicking on **Next** and choosing *equation.tex* in the pull down list and we see the following:



Here we see the logistic regression model, in LaTeX, in the output pane. If we select *model.txt* we can then see the model code that the algebra system will interpret:

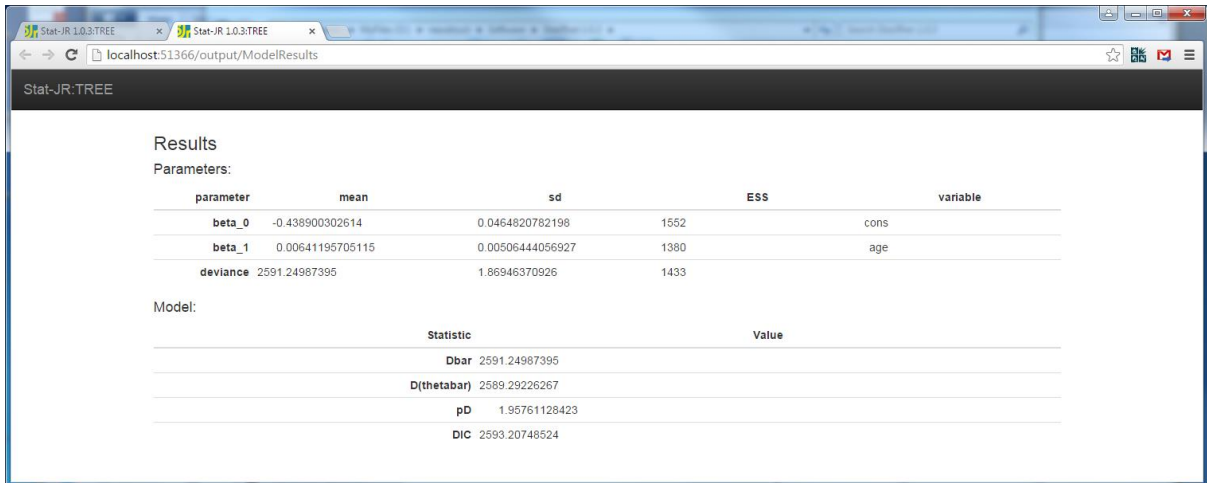


Now choosing **algorithm.tex** from the right-hand pane, and placing it in its own tab in the browser window, gives the following:

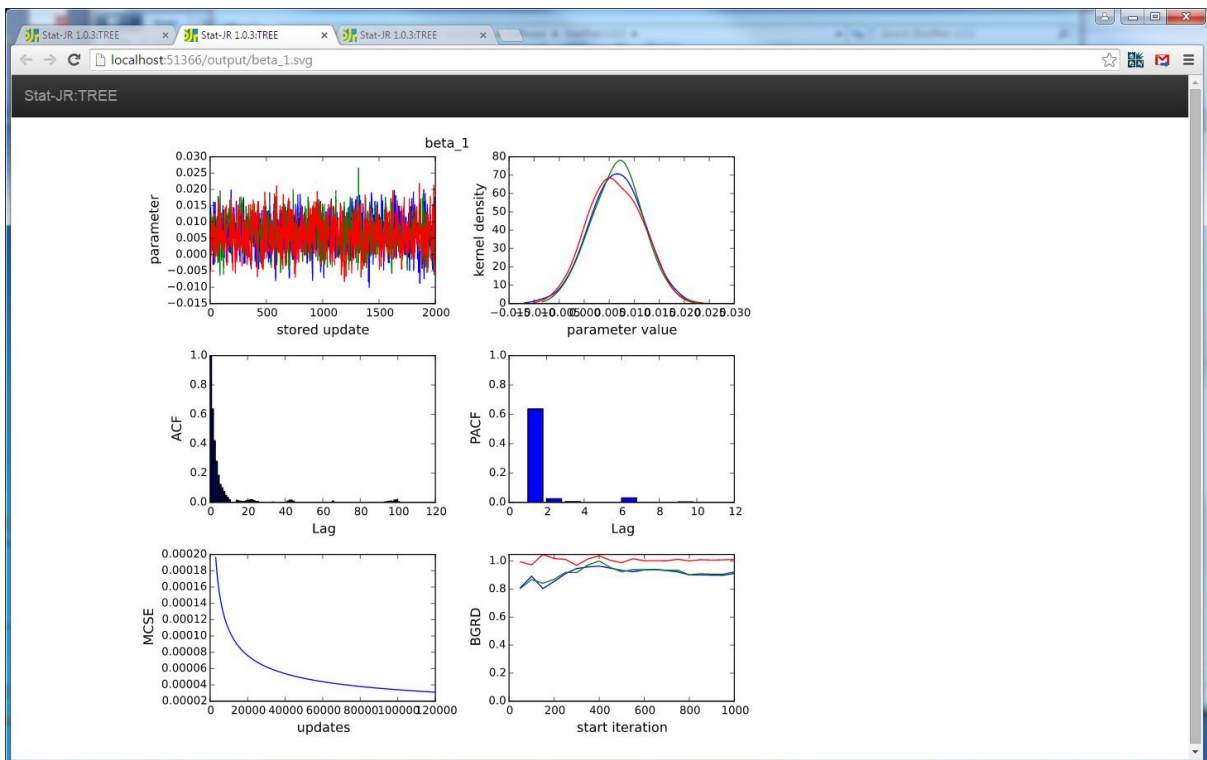


Here we see that the eStat engine uses a different MCMC method, random walk Metropolis, for the steps for the fixed effects (*beta0* and *beta1*) when fitting logistic regression models. We will come back to this modelling decision in Section 4.4.4 when we compare different software packages.

Returning to the main pane and clicking on **Run** will now run the model. Once it has finished, if we select *ModelResults* from the list, and look at it in a new tab, we get the following:



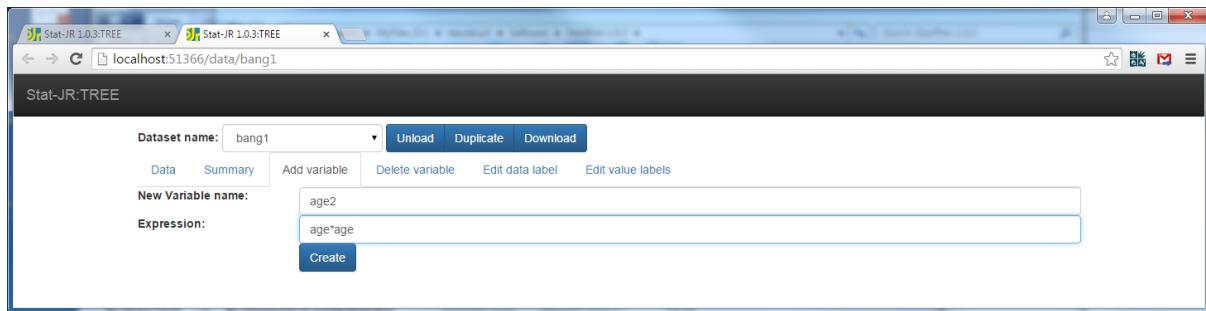
Perhaps disappointedly and surprisingly, age doesn't appear to have a significant effect (its estimate (0.0064) is similar in magnitude to its standard error (0.0051)). To see this more clearly we can look at the graph `beta_1.svg` in its own browser tab:



Here, whilst the values on the x-axis overlap and therefore aren't particularly clear, we can see that all three chains show strong support for the value 0.00 in the kernel density plot (i.e. it's comfortably within the distribution). It might be the case, however, that contraceptive use has a non-linear relationship with age (possibly quadratic) and this could also be confounded by how far through their own family-formation process the woman is, which we will model via the variable `lc`. We might also be interested in accounting for any clustering effects of having women nested within districts.

In order to fit a quadratic function to age we will need to construct the variable `age2` which we can easily do by viewing the Dataset and use the variable creation tool.

Return to the main screen and select **View** from the **Dataset** pull-down list at the top of the page
Click on the **Add Variable** tab and type the following:



Here we are going to overwrite the existing dataset (at least in temporary memory) with a version in which we have appended an additional column to it. Clicking on **Create** and looking at the data by (clicking on the **Data** tab) below gives the following:

	woman	district	use	lc	age	urban	educ	hindu	d_illit	d_pray	cons	age2
1	1	1	0	three+kids	18.44	1	1	0	0.58	0.64	1	340.033630371
2	2	1	0	nokids	-5.56	1	1	1	0.58	0.64	1	30.9135932922
3	3	1	0	twokids	1.44	1	2	0	0.58	0.64	1	2.0736014843
4	4	1	0	three+kids	8.44	1	1	0	0.58	0.64	1	71.2336120605
5	5	1	0	nokids	-13.56	1	1	0	0.58	0.64	1	183.873580933
6	6	1	0	nokids	-11.56	1	1	0	0.58	0.64	1	133.633590698
7	7	1	0	three+kids	18.44	1	1	0	0.58	0.64	1	340.033630371
8	8	1	0	three+kids	-3.56	1	1	0	0.58	0.64	1	12.6735943821
9	9	1	0	onekid	-5.56	1	1	0	0.58	0.64	1	30.9135932922
10	10	1	0	three+kids	1.44	1	1	0	0.58	0.64	1	2.0736014843
11	11	1	1	nokids	-11.56	1	1	0	0.58	0.64	1	133.633590698
12	12	1	0	nokids	-2.56	1	1	0	0.58	0.64	1	6.55359745026
13	13	1	0	onekid	-4.56	1	1	0	0.58	0.64	1	20.7935943604
14	14	1	0	three+kids	5.44	1	1	0	0.58	0.64	1	29.5936050415
15	15	1	0	three+kids	-0.559999	1	1	0	0.58	0.64	1	0.313599407673
16	16	1	1	three+kids	4.44	1	1	0	0.58	0.64	1	19.7136039734
17	17	1	0	nokids	-5.56	1	1	0	0.58	0.64	1	30.9135932922
18	18	1	1	three+kids	-0.559999	1	2	0	0.58	0.64	1	0.313599407673

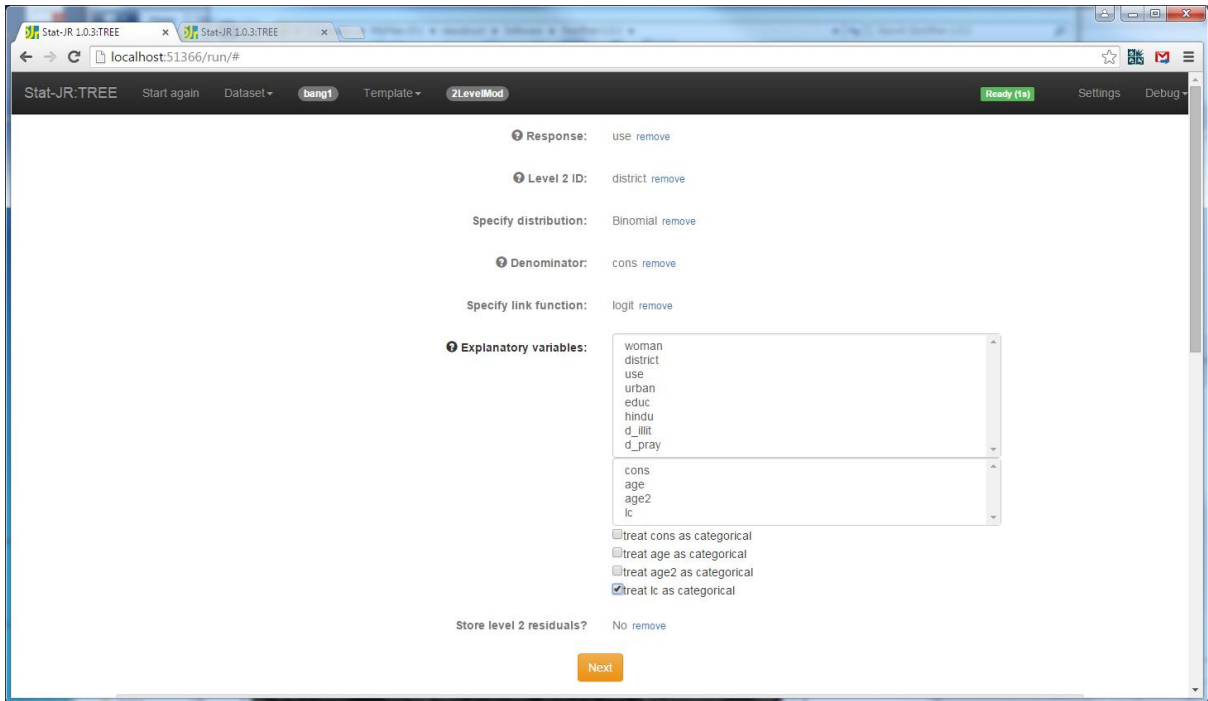
Here you see age2 (age²) appearing in the column on the far right. Whilst we could explore adding further explanatory variables to this 1-level model, we are going to move straight into fitting a 2-level model to account for districts.

4.4.3 Multilevel modelling of the data

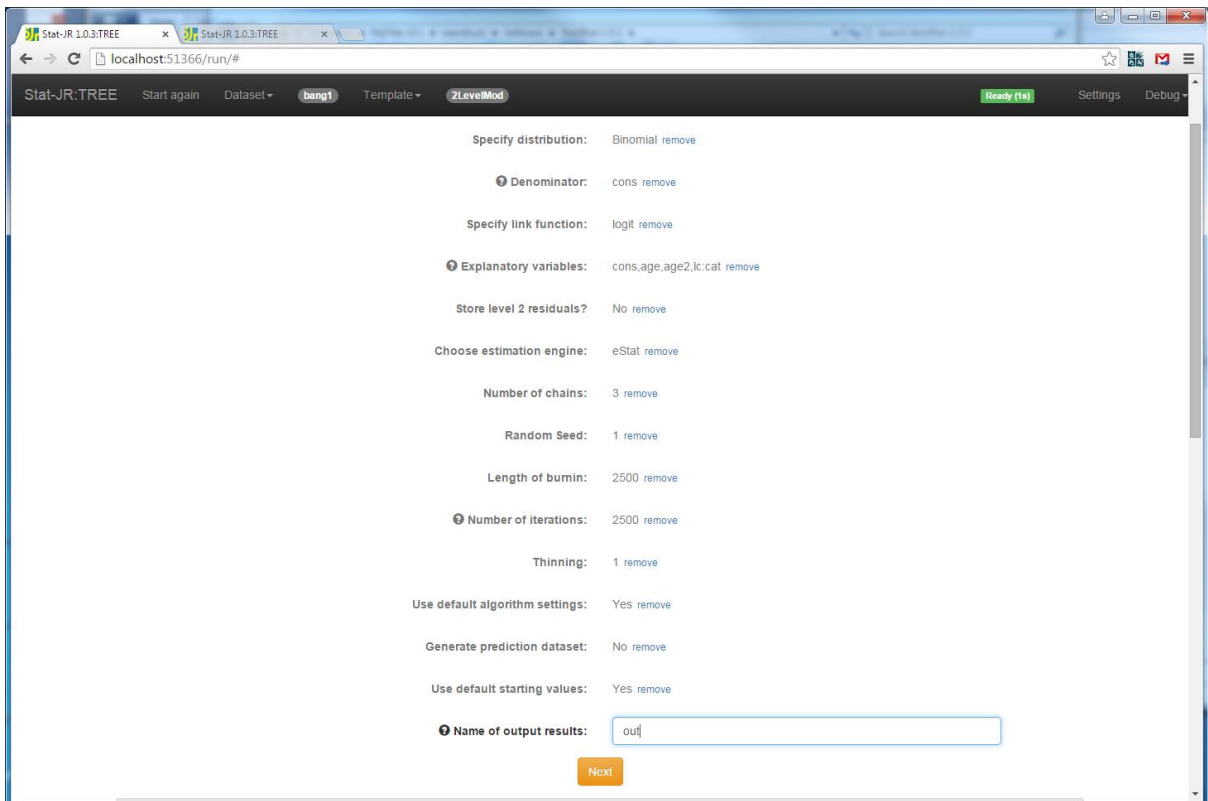
We will now require a template that will fit a 2-level logistic regression model to our dataset. In the earlier sections we looked at the template **2LevelMod** and we will once again use it here and also illustrate how to fit categorical predictor variables.

On the main tab, click on **Choose** in the **Template** pull down list and select **2LevelMod** and click on **Use** button to run this template.

Fill in the template inputs as follows:



Here we need to specify several extra inputs, including an input for the level 2 identifiers and also to let the software know which predictor variables are categorical. Continue with the inputs as follows:



Clicking on **Next** will run the algebra system and set up code to fit the model. If we select *model.txt* in the output list we will see the following:

```

model {
  for (i in 1:length(use)) {
    use[i] ~ dbin(p[i], cons[i])
    logit(p[i]) <- cons[i] * beta_0 + age[i] * beta_1 + age2[i] * beta_2 + lc_1[i] * beta_3 + lc_2[i] * beta_4 + lc_3[i] * beta_5 + u[district[i]]
  }

  for (j in 1:length(u)) {
    u[j] ~ dnorm(0, tau_u)
  }

  # Priors
  beta_0 ~ dflat()
  beta_1 ~ dflat()
  beta_2 ~ dflat()
  beta_3 ~ dflat()
  beta_4 ~ dflat()
  beta_5 ~ dflat()

  tau_u ~ dgamma(0.001000, 0.001000)
  sigma2_u <- 1 / tau_u
}

```

Here we see the more complicated model code for this 2-level model in the left-hand pane. Note that the *lc* predictor is treated as categorical and thus appears as 3 dummy variables (*lc_1* – *lc_3*)

If we select *tau_u.xml* in the output list we will see the following:

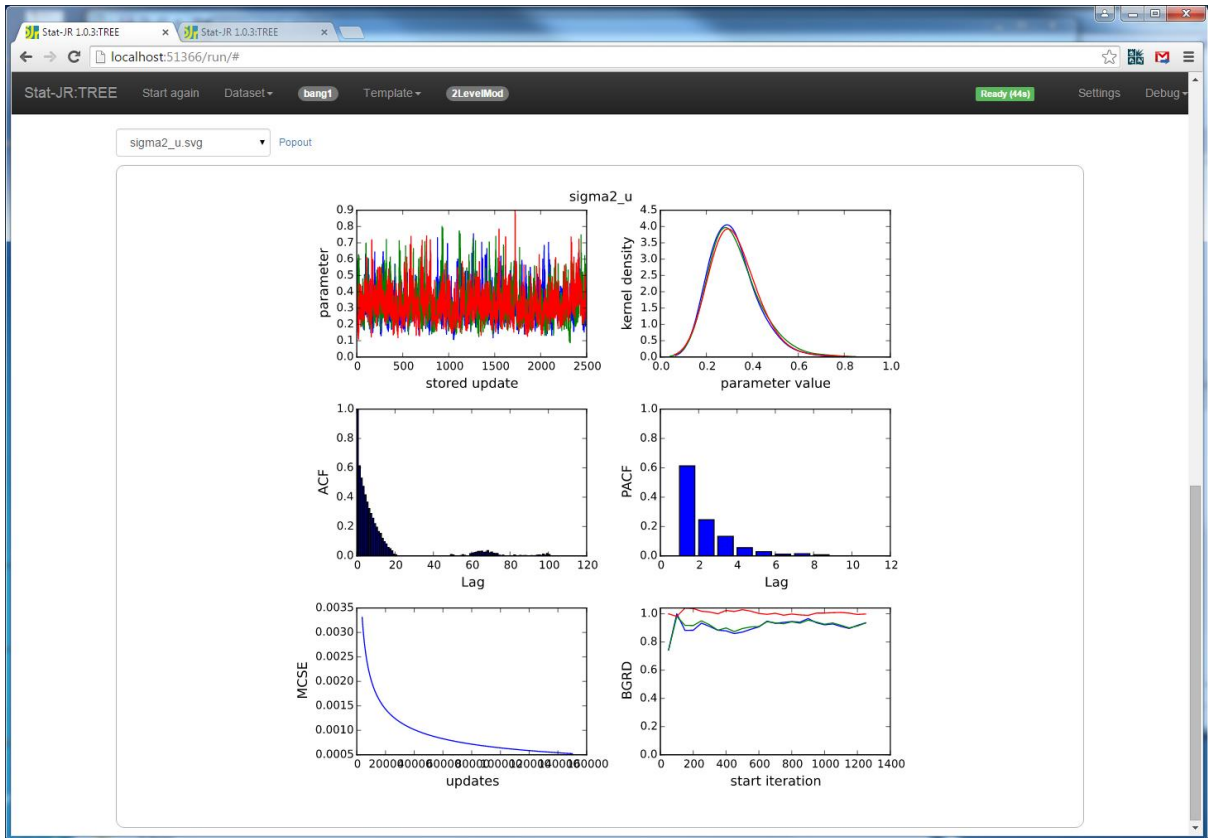
Use Gibbs sampling from conditional posterior for tau_u:

$$\tau_{u,i} \sim \Gamma \left(0.001 + 0.5 \times \text{length}(u), 0.001000 + \frac{\sum_{j=1}^{\text{length}(u)} u_j^2}{2} \right)$$

$$\tau_{u,i} \sim \Gamma \left(30.001, 0.001 + \left(\sum_{j=1.0}^{60.0} u_j^{2.0} \right) \times 0.5 \right)$$

Here we see the algorithm step for the parameter *tau_u*. Although most parameters in this model are updated by Random Walk Metropolis sampling, this parameter is updated by Gibbs Sampling as its conditional posterior distribution has a standard form.

If we now click on **Run** then after 52s (on a machine with Intel Core i7-3770S; this includes time for compiling and adapting) the model will have run and if we select *sigma2_u.svg* we will see the following:



Here we can see that convergence and mixing, for this parameter at least, are reasonable. In fact, if we look at the diagnostic plots for the other parameters, we see similar convergence there as well. Next we can look at *ModelResults* in its own tab to see the parameter estimates:

Stat-JR: TREE

Results

Parameters:

parameter	mean	sd	ESS	variable
σ^2_u	0.320003655857	0.099343397633	807	
β_0	-0.761482066363	0.183607015726	92	cons
β_1	0.00759569398534	0.00974610865938	183	age
β_2	-0.00488119664141	0.000743128200041	319	age ²
β_3	0.760961328671	0.164095856619	227	lc_onekid
β_4	0.808508671551	0.191526987545	169	lc_twokids
β_5	0.805088116084	0.191860340621	114	lc_three+kids
τ_u	3.42956763896	1.07988494419	760	
deviance	2351.0823043	11.2344889316	1242	

Model:

Statistic	Value
Dbar	2351.0823043
D(thetabar)	2308.12575411
pD	42.9565501939
DIC	2394.03885449

Here we see that β_2 is significant and negative (and larger than β_1) suggesting a quadratic fit to the age predictor. As the data is centred around its mean, this implies that contraceptive use is reduced the further from the mean age the woman is. We will look at this in more detail at the end of the chapter.

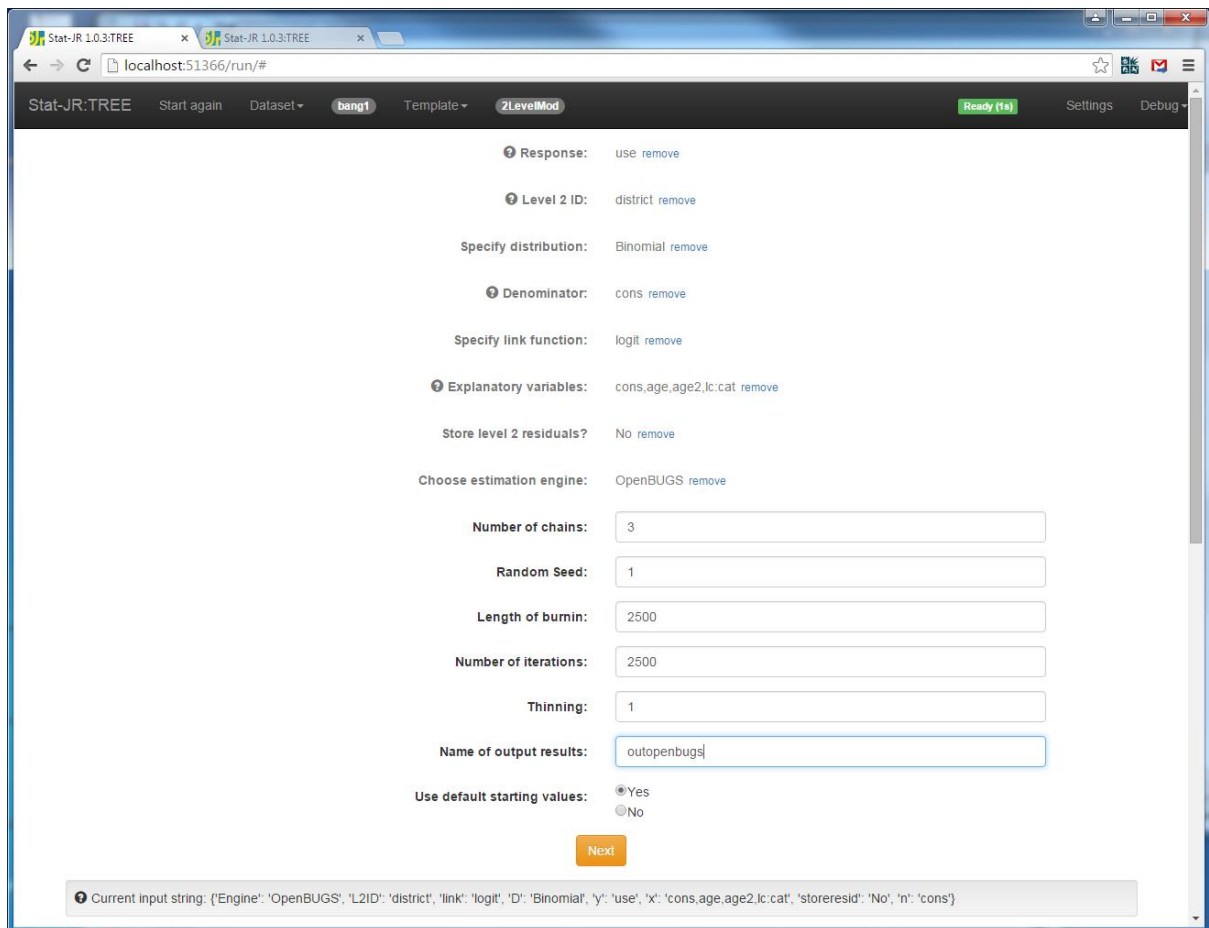
The parameters β_3 - β_5 are all significant, and positive (and of similar magnitude), which suggests that women with children are more likely to use contraceptives than those without. The parameter σ_u is fairly large, suggesting there are differences between districts in terms of contraceptive use.

What is slightly disappointing here are the ESS values for all the fixed parameters. We have run each chain, after burnin, for 2,500 iterations resulting in a total of 7,500 actual iterations (i.e. from 3 chains) but the effective sample sizes are of the order of 100-350. As this indicates, the default algorithm in eStat – random walk Metropolis – is not very efficient for this example. We will look at two possible solutions in the next two sections.

4.4.4 Comparison between software packages

Not all software packages fit the same MCMC algorithm for this model. So here we will show how to fit the same model in another package, OpenBUGS, which uses a different method: namely multivariate updating for the fixed effects in a GLMM, as developed by Gamerman (1997). This method results in slower estimation, but, as we will see, far better ESS. We will then look at a table comparing all the possible MCMC algorithms in the different packages for this model, which you can verify for yourselves.

To fit the model in OpenBUGS click on the **remove** text next to **Choose estimation engine** and set-up the model as follows:



Clicking on **Next** and **Run** will (after 2 min 18s on my machine) give the following, having selected *ModelResults* from the drop-down box above the output pane, and opening it in a new tab:

The screenshot shows a web browser window with two tabs for 'Stat-JR 1.0.1: TREE'. The address bar shows 'localhost:52622/output/ModelResults'. The page content is as follows:

Results

Parameters:

parameter	mean	sd	ESS
beta_0	-0.79044984	0.172520147789	2595
beta_1	0.0065719269272	0.00909771213816	5031
beta_2	-0.00481059466667	0.000726599414238	5057
beta_3	0.7824753212	0.162317868714	5291
beta_4	0.825564053333	0.18610036524	5181
beta_5	0.827532796	0.183875558079	4443
deviance	2351.19746667	11.5283450785	4441
sigma2_u	0.317104969333	0.100486311959	1753
tau_u	3.47796493333	1.13186027888	1645

Model:

Statistic	Value
Dbar_use	2351.0
Dhat_use	2309.0
pD_use	42.66
DIC_use	2394.0
Dbar_total	2351.0
Dhat_total	2309.0
pD_total	42.66
DIC_total	2394.0

Here we see far better effective sample size values, with runs of 7,500 iterations translating into ESS values of between 2,500 and 5,500 for the beta parameters.

We can repeat this analysis using WinBUGS, JAGS and MLwiN with the same run lengths. Note for JAGS you will need to edit the initial value files or it will not run. To do this view each in the output window and click on the **Edit** button. If you change the value for beta_2 (the fixed effect associated with age2) from 0.1 to 0.0 in all three initial values files and click **Save** each time then JAGS should run. It should also be noted here that results may vary a little if you have different versions of the third party software packages or have changed options in them.

We could also fit the model using the *MCMCglmm* package in R, although here we would need to run a single chain and logistic regression models for binary data are the one GLMM where the answers can be a little different as it assumes over-dispersion which is inappropriate in this case.

The table overleaf¹ details the results of fitting many of these options:

¹ This particular comparison used WinBUGS 1.4.3, OpenBUGS 3.2.3, JAGS 3.4.0, MLwiN 2.34, all run on Windows 64-bit machine with Intel Core i7-3770S; eStat times are of the form: *including compiling time (excluding compiling time)*.

Parameter	eStat	WinBUGS	OpenBUGS	JAGS	MLwiN	eStat orthogonal
Beta0	-0.761(0.183)	-0.789(0.170)	-0.790(0.173)	-0.776(0.177)	-0.835(0.170)	-0.784(0.180)
Beta0 ESS	92	396	2595	255	93	979
Beta1	0.0076(0.0097)	0.0068(0.0090)	0.0066(0.0091)	0.0069(0.0090)	0.0050(0.0089)	0.0068(0.0096)
Beta1 ESS	183	951	5031	550	247	1792
Beta2	-0.0049(0.00074)	-0.0048(0.00072)	-0.0048(0.00073)	-0.0048(0.00071)	-0.0047(0.00071)	-0.0048(0.00073)
Beta2 ESS	319	1286	5057	926	315	1799
Beta3	0.761(0.164)	0.779(0.160)	0.782(0.162)	0.778(0.163)	0.799(0.162)	0.779(0.165)
Beta3 ESS	227	1117	5291	630	268	1686
Beta4	0.809(0.192)	0.822(0.181)	0.826(0.186)	0.818(0.188)	0.856(0.183)	0.823(0.190)
Beta4 ESS	169	780	5181	477	196	1726
Beta5	0.805(0.192)	0.824(0.180)	0.828(0.184)	0.817(0.187)	0.863(0.177)	0.823(0.191)
Beta5 ESS	114	547	4443	329	131	1666
Sigma2u	0.320(0.099)	0.318(0.100)	0.317(0.100)	0.317(0.100)	0.328(0.103)	0.322(0.101)
Sigma2u ESS	807	1764	1753	1445	733	756
Pd	42.96	42.44	42.66	42.16	43.08	43.21
DIC	2394.03	2393.36	2394.0	2393.39	2393.65	2394.51
Time (s)	40 (24)	195	138	176	7	36 (22)

In summary we see that MLwiN is by far the fastest of the packages, with eStat quicker than the other three as well. Both MLwiN and eStat use the simple random walk Metropolis algorithm, which is not the best method for this model and gives fairly poor ESS. Interestingly, both WinBUGS and OpenBUGS use the Gamerman method, but in this case OpenBUGS performs better in terms of time taken and ESS. This is somewhat puzzling as when each is run with a single chain, their performance is almost identical. Finally, for this example, JAGS has a similar speed intermediate to the two BUGS packages but its performance is relatively disappointing with regard ESS; however, there have been many comparisons between JAGS and BUGS for different models, and which method is better varies from model to model, so we shouldn't dismiss it based on just this one example. The final column shows another eStat method which we will discuss next.

4.4.5 Orthogonal parameterisation.

The reason eStat (and MLwiN) perform badly in terms of ESS in this instance is that they are performing single-site updating, and the parameters are correlated. So here we will consider a reparameterisation method that aims to fit parameters that are less correlated, and then translates them back to the original parameters. For this we construct a set of orthogonal vectors from the original predictor variables (see Browne et al. 2009 for details).

We will therefore now look at the **NLevelOrthogParamRS** template in order to use orthogonalisation on our model. This template actually fits a larger family of models: those with any number of higher levels/classifications (hence 'NLevel'), allowing for the possibility of random slopes at each of these levels (hence 'RS'), and so our 2-level random intercept model is perhaps the simplest case that the template fits.

Click on the **Template** pull down list and click **Choose** then select **NLevelOrthogParamRS** from the template list.

Click on **Use** and fill in the template inputs as follows:

The screenshot shows the Stat-JR 1.0.1: TREE web interface. The browser address bar shows localhost:52622/run/#. The interface has a dark header with 'Stat-JR: TREE', 'Start again', 'Dataset' (set to 'bang1'), 'Template' (set to 'NLevelOrthogParamRS'), and a green 'Ready (1s)' button. Below the header, the configuration for the 'NLevelOrthogParamRS' template is displayed with various settings and 'remove' links:

- Number of Classifications: 1 remove
- Classification 1: district remove
- Response: use remove
- Specify distribution: Binomial remove
- Denominator: cons remove
- Specify link function: logit remove
- Explanatory variables: cons, age, age2, Ic: cat remove
- Explanatory variables random at district classification: cons remove
- Do you want to use orthogonal parameterisation?: Yes remove
- Type: Orthogonal remove
- Store residuals?: No remove
- Choose estimation engine: eStat remove
- Number of chains: 3 remove
- Random Seed: 1 remove
- Length of burnin: 2500 remove
- Number of iterations: 2500 remove
- Thinning: 1 remove
- Use default algorithm settings: Yes remove
- Generate prediction dataset: Yes remove
- Use default starting values: Yes remove
- Name of output results: [empty text box]

Giving a name for the results and clicking on **Next** and selecting *equation.tex* in the pull down list (we've opened it in a new tab) will show the following:

```

use_i ~ Binomial(cons_i, pi_i)
logit(pi_i) = beta_0^*orthcons_i + beta_1^*orthage_i + beta_2^*orthage2_i + beta_3^*orthlc_1_i + beta_4^*orthlc_2_i + beta_5^*orthlc_3_i + v_{0,district}^{(2)} cons_i
v_{0,district}^{(2)} ~ N(0, sigma_u2^2)
tau_u2 ~ Gamma(0.001, 0.001)
sigma_u2^2 = 1/tau_u2
beta_0^* <math>\propto 1</math>
beta_1^* <math>\propto 1</math>
beta_2^* <math>\propto 1</math>
beta_3^* <math>\propto 1</math>
beta_4^* <math>\propto 1</math>
beta_5^* <math>\propto 1</math>
beta_0 = 1.0*beta_0^* - 0.00204810386113*beta_1^* - 81.1914651166*beta_2^* - 0.214084151662*beta_3^* - 0.276389332062*beta_4^* - 0.678464606418*beta_5^*
beta_1 = 0.0*beta_0^* + 1.0*beta_1^* - 3.98134713968*beta_2^* + 0.00731689215009*beta_3^* - 0.00244525658859*beta_4^* - 0.0356343746801*beta_5^*
beta_2 = 0.0*beta_0^* + 0.0*beta_1^* + 1.0*beta_2^* + 0.000382130301133*beta_3^* + 0.000959517791234*beta_4^* + 0.00134115598842*beta_5^*
beta_3 = 0.0*beta_0^* + 0.0*beta_1^* + 0.0*beta_2^* + 1.0*beta_3^* + 0.217129714164*beta_4^* + 0.469386306038*beta_5^*
beta_4 = 0.0*beta_0^* + 0.0*beta_1^* + 0.0*beta_2^* + 0.0*beta_3^* + 1.0*beta_4^* + 0.627080532383*beta_5^*
beta_5 = 0.0*beta_0^* + 0.0*beta_1^* + 0.0*beta_2^* + 0.0*beta_3^* + 0.0*beta_4^* + 1.0*beta_5^*

```

Here we see that the model code is actually fitting a different set of predictors, each with the prefix 'orth' and a corresponding set of coefficients. There is then a set of deterministic statements that translate these coefficient values to the coefficient values for the original predictors (again, see Browne et al. (2009) for details)

Clicking on the **Run** button will run the model (which took 36s on this particular machine, including compiling), after which selecting *ModelResults* from the pull down list, and popping out into a new tab, gives the following:

Results

Parameters:

parameter	mean	sd	ESS	variable
sigma2_u0_1	0.322611997838	0.100988344452	756	
deviance	2351.30603584	11.7512785968	1293	
betaort_0	-0.584422193533	0.0940982100777	315	
betaort_1	0.00925227105254	0.00611737929645	1863	
betaort_2	-0.00633183225716	0.000669841236497	1679	
betaort_3	0.325993846063	0.129887324437	1799	
betaort_4	0.307001409516	0.142090095981	1850	
betaort_5	0.823501468221	0.191451499533	1667	
beta_0	-0.783622684065	0.180131036563	979	cons
beta_1	0.00675107940155	0.00956597217846	1792	age
beta_2	-0.00480826241604	0.000730012649254	1799	age2
beta_3	0.779101816497	0.165132052505	1686	lc_onekid
beta_4	0.823279406876	0.189574245211	1726	lc_twokids
beta_5	0.823437482593	0.191418149556	1666	lc_three+kids
tau_u0_1	3.41856462849	1.133235441	689	

Model:

Statistic	Value
Dbar	2351.30603584
D(thetabar)	2308.09827984
pD	43.2077560009
DIC	2394.51379184

The estimates, their ESS, and the time taken to run the model are all added to the end of the software comparison table we looked at above. It indicates that, compared to the other method we employed to fit the model in eStat, there is no obvious overhead incurred when performing the orthogonalising algorithm, and it is much faster than OpenBUGS, and the ESS are now much better (if still not as good as OpenBUGS). We therefore have two ways of fitting the model that are reasonably comparable in terms of ESS/s, with little to choose between them. This orthogonalising approach is also available in MLwiN: this will be faster again, and should have similar ESS to the method in eStat, and therefore may be the best overall in terms of ESS/s, but we leave this for the reader to investigate.

4.4.6 Predictions from the model

When we ran this model we discussed some interpretation of the fit, but it would be nice to plot some predictions from the model as well. In this latest version of Stat-JR we have added the option to store predictions when fitting the model. So hopefully in the last model fit you will have ticked yes to the generate prediction dataset question. This will generate a new dataset named *prediction_datafile* which contains the original data and several prediction columns formed from the model fit.

To use this dataset we need to select **Choose** on the dataset list and select *prediction_datafile* from the list and click **Use**.

In fact the dataset has a full prediction column called *pred_full* but this also contains the district random effects. We would here like to simply predict from the fixed part of the model so we can construct the variable *pred_fixed* as follows:

Click on **View** from the **Dataset** menu, then choose **Add variable**, and input the new variable *pred_fixed* as indicated below.

Click on **Create** to create the variable

Dataset name: prediction_datafile Unload Duplicate Download

Data Summary Add variable Delete variable Edit data label Edit value labels

New Variable name:

pred_fixed

Expression:

pred_full - pred_u0_0

Create

This has created a variable on the fixed predictor scale but as we are fitting a logistic regression we need to take an anti-logistic transform to convert these predictions to probabilities. This can be done by creating another column in the dataset as shown below:

Dataset name: prediction_datafile Unload Duplicate Download

Data Summary Add variable Delete variable Edit data label Edit value labels

New Variable name:

fitprob

Expression:

$\exp(\text{pred_fixed}) / (1 + \exp(\text{pred_fixed}))$

Create

In order to plot separate fitted curves for the various numbers of living children we can use the template **XYGroupPlot** as shown below:

Stat-JR: TREE Start again Dataset prediction_datafile Template XYGroupPlot Ready (1s) Settings Debug

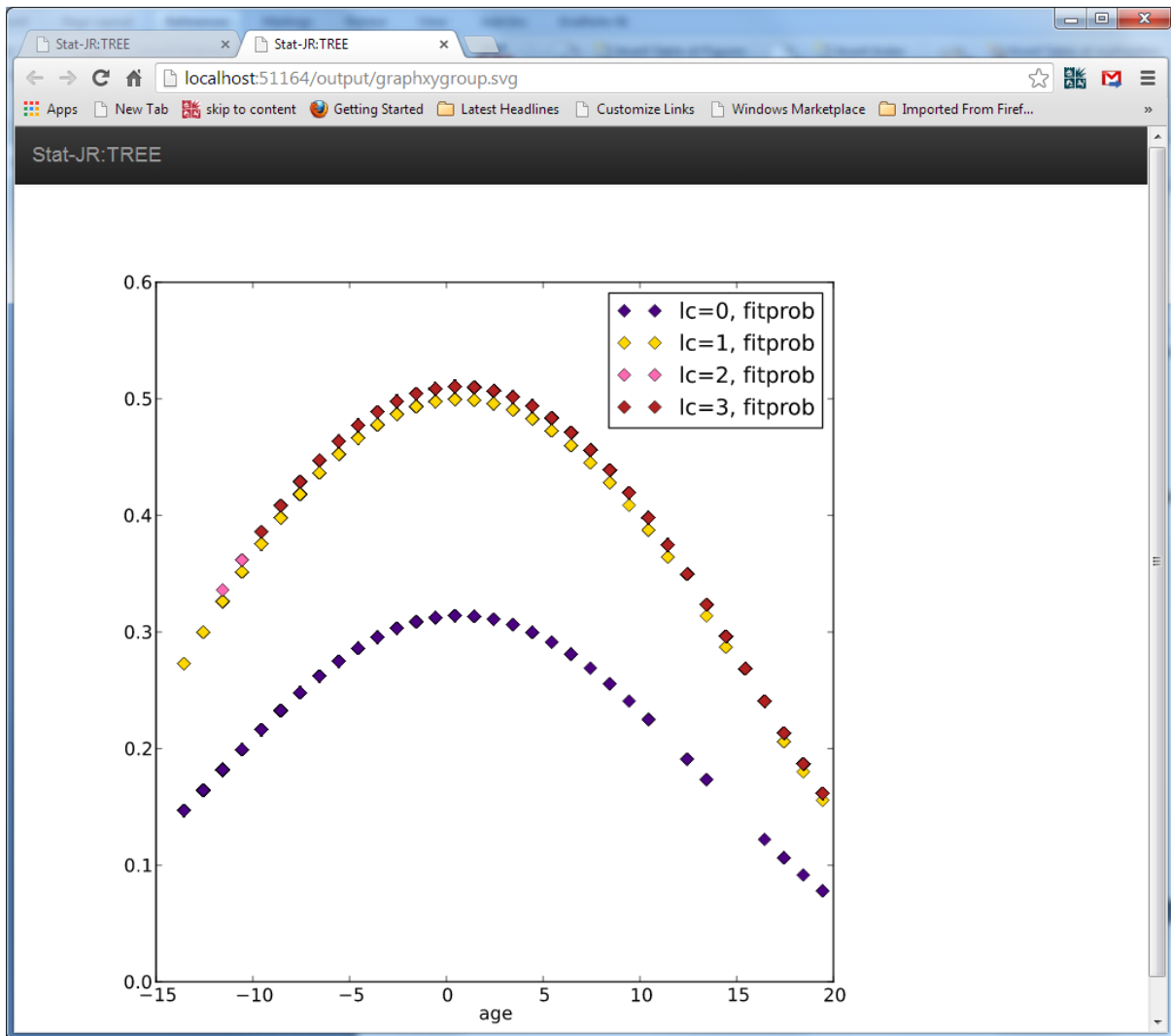
X values: age remove

Y values: fitprob remove

Grouped by: lc remove

Run

Clicking on Run and popping out *graphxygroup.svg* gives the following:



Here we see the four curves (although three of them are very close together) which clearly showing that the women with children have higher probabilities of using contraceptives, and that the peak for each group is around the average age of the sample, as discussed earlier.

Hopefully this section has shown firstly that Stat-JR can fit models other than Normal response models; in fact there are a vast number of model templates which fit lots of other model classes. Secondly, we hope we've shown its utility in terms of comparing model-fitting across different software packages for different models, accessing each from a common hub.

4.5 Miscellaneous other topics e.g. Data Input/Export

Stat-JR works with datasets saved in Stata format, i.e. with a *.dta* extension. It looks for these in the...*\datasets* folder of the Stat-JR install, and also in a folder saved, by default, under your user name, e.g. *C:\Users\YourName\.statjr\datasets* (you can change the path via **Settings** in the black bar at the top of the browser window in the TREE interface).

If your dataset is already in *.dta* format (see below), then you can upload it, in TREE, via (i) **Dataset > Upload** (menu options in the black bar at the top of the browser window), which will upload it into the temporary memory cache, or by (ii) saving your dataset in the *StatJR\datasets* folder, and then selecting **Debug > Reload datasets** (again, accessible via the black bar at the top of the browser window). If, instead, you have it (iii) saved as a *.txt* file, you can use Stat-JR's *LoadTextFile* template

to save it into the temporary memory cache (the template *LoadTextFileMoreOptions* allows the user to specify more particulars, and can also handle string variables).

In the case of options (i) and (iii) the dataset will be available for use in the current session, but you then need to download it (as a *.dta* file) via **Dataset > Download** (e.g. saving it into the *StatJR\datasets* folder) for use in the future sessions too.

So, via option (iii) (and downloading), Stat-JR will save your dataset as a *.dta* file, but you can also create *.dta* files via Stata, MLwiN and R (e.g. the foreign package in R).

5 References

Brooks, S.P. and Gelman, A. (1998) General Methods for Monitoring Convergence of Iterative Simulations. *Journal of Computational and Graphical Statistics*, **7**: 434-455.

Browne, W.J. (2012) *MCMC Estimation in MLwiN, v2.26*. Centre for Multilevel Modelling, University of Bristol.

Browne, W.J., Steele F., Golalizadeh, M., and Green M.J. (2009) The use of simple reparameterizations to improve the efficiency of Markov chain Monte Carlo estimation for multilevel models with applications to discrete time survival models *Journal of Royal Statistical Society, Series A*. **172**: 579-598

Gamerman, D. (1997). Sampling from the posterior distribution in generalized linear mixed models. *Statistics and Computing*, **7**, 57-68.

Gelfand, A.E., Hills, S.E., Racine-Poon, A. and Smith, A.F.M. (1990). Illustration of Bayesian inference in Normal data models using Gibbs Sampling. *Journal of the American Statistical Association*, **85**:972-985.

Lillard, L.A. and Panis C.W.A. (2003) *aML Multilevel Multiprocess Statistical Software, Version 2.0*. EconWare, Los Angeles, California.

Lunn, D.J., Thomas, A., Best, N., and Spiegelhalter, D. (2000). WinBUGS - a Bayesian modelling framework: concepts, structure, and extensibility. *Statistics and Computing*, **10**: 325--337.

Lunn, D., Spiegelhalter, D., Thomas, A. and Best, N. (2009). The BUGS project: Evolution, critique, and future directions, *Statistics in Medicine*, **28**, 3049-3067.

Plummer, Martyn (2003). JAGS: A Program for Analysis of Bayesian Graphical Models Using Gibbs Sampling, *Proceedings of the 3rd International Workshop on Distributed Statistical Computing (DSC 2003)*, March 20–22, Vienna, Austria. ISSN 1609-395X.

R Development Core Team (2011). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. ISBN 3-900051-07-0

Rasbash, J., Charlton, C., Browne, W.J., Healy, M. and Cameron, B. (2009). *MLwiN Version 2.1*. Centre for Multilevel Modelling, University of Bristol.

Spiegelhalter, D.J., Best, N.G., Carlin, B.P. and van der Linde, A. (2002). Bayesian measures of model complexity and fit (with discussion). *Journal of the Royal Statistical Society, Series B*, **64**: 191-232.

6 Appendix: List of Third Party Software that are used by Stat-JR

Stat-JR makes use of several third party software products that are included within the distributed code or (in the case of MinGW) need to be downloaded separately. These software products each have a license file that can be viewed from the links in the table below and/or in the licences subdirectory of the installed code.

Package	Link	Licence terms
beautifulsoup	http://bazaar.launchpad.net/~leonardr/beautifulsoup/bs4/view/head:/COPYING.txt	MIT
BLAS	http://www.netlib.org/blas/faq.html#2	Own licence (Netlib)
Bootstrap	https://github.com/twitter/bootstrap/blob/master/LICENSE	MIT
cssselect	http://www.opensource.org/licenses/bsd-license.php	BSD
cx_freeze	http://cx-freeze.readthedocs.org/en/latest/license.html	PSF
cycler	https://opensource.org/licenses/BSD-3-Clause	BSD
dateutil	http://opensource.org/licenses/BSD-2-Clause	Simplified BSD
decorator	https://micheles.googlecode.com/hg/decorator/documentation.html#license	BSD
keepalive	https://github.com/wikier/keepalive/blob/master/LICENSE	LGPL
html5lib	https://github.com/html5lib/lib/html5lib-python/blob/master/LICENSE	MIT
isodate	http://www.opensource.org/licenses/bsd-license.php	BSD
jqgrid	http://www.trirand.com/blog/?page_id=87	Dual MIT/GPL(v2)
jquery	http://jquery.org/license	MIT

jquery-cookie	https://github.com/carhartl/jquery-cookie/blob/master/jquery.cookie.js	MIT
jQuery File Upload	http://opensource.org/licenses/MIT	MIT
jQuery text align	http://www.opensource.org/licenses/bsd-license.php	BSD
jquery-treeview	https://github.com/jzaefferer/jquery-treeview	Dual MIT/GPL
jquery-ui	http://jquery.org/license	MIT
jQuery-xpath	http://opensource.org/licenses/MIT	MIT
LAPACK	http://www.netlib.org/lapack/LICENSE.txt	Modified BSD
lxml	http://lxml.de/index.html#license	BSD
mako	http://www.opensource.org/licenses/mit-license.php	MIT
markupsafe	http://www.opensource.org/licenses/bsd-license.php	BSD
MathJax	http://cdn.mathjax.org/mathjax/2.0-latest/LICENSE	Apache
matplotlib	http://matplotlib.sourceforge.net/users/license.html	Modified BSD
MinGW	http://www.mingw.org/license	Not distributed with software directly
networkx	http://networkx.github.io/documentation/development/reference/legal.html	BSD
numexpr	http://www.opensource.org/licenses/mit-license.php	MIT
numpy	http://numpy.scipy.org/license.html#license	BSD
pandas	http://pandas.pydata.org/pandas-docs/stable/overview.html#license	Modified BSD
patsy	https://github.com/pydata/patsy/blob/master/LICENSE.txt	BSD
ply	http://www.dabeaz.com/ply/README.txt	BSD
prov	https://github.com/trungdong/prov/blob/master/LICENSE	MIT
provpy	http://opensource.org/licenses/BSD-2-Clause	BSD

pyparsing	http://www.opensource.org/licenses/mit-license.php	MIT
pyquery	http://www.opensource.org/licenses/bsd-license.php	BSD
Python	http://docs.python.org/license.html	PSF
pytz	https://pypi.python.org/pypi/pytz/	MIT
rdflib	http://www.opensource.org/licenses/bsd-license.php	BSD
reset-fonts-grids	http://yuilibrary.com/license/	BSD
scipy	http://www.scipy.org/License_Compatibility	BSD
setuptools	http://docs.python.org/license.html /	PSF
six	https://bitbucket.org/gutworth/six/src/e3da7fd96039a6ed89493f89d121c4f3797e6713/LICENSE?at=default	MIT
sparqlwrapper	http://www.w3.org/Consortium/Legal/2002/copyright-software-20021231	W3C
statsmodels	https://github.com/statsmodels/statsmodels/blob/master/LICENSE.txt	Modified BSD
tinymce	https://github.com/tinymce/tinymce/blob/master/LICENSE.TXT	LGPL
weave	http://projects.scipy.org/scipy/browser/trunk/Lib/weave/LICENSE.txt?rev=1511	BSD
web.py	https://github.com/webpy/webpy/blob/master/LICENSE.txt	PSF