

Generative DAGs as an Interface Into Probabilistic Programming with the R Package causact

Adam J. Fleischhacker¹ and Thi Hong Nhung Nguyen²

¹ Adam Fleischhacker, JP Morgan Chase Faculty Fellow, University of Delaware, Newark, DE 19716
² Institute for Financial Services Analytics, University of Delaware, Newark, DE 19716

DOI: [10.xxxxxx/draft](https://doi.org/10.xxxxxx/draft)

Software

- [Review](#) ↗
- [Repository](#) ↗
- [Archive](#) ↗

Editor: [Open Journals](#) ↗

Reviewers:

- [@openjournals](#)

Submitted: 01 January 1970

Published: unpublished

License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](#)).

Summary

The causact package provides R functions for visualizing and running inference on generative directed acyclic graphs (DAGs). Once a generative DAG is created, the package automates Bayesian inference via the greta package ([Golding, 2019](#)) and TensorFlow ([Dillon et al., 2017](#)). The package eliminates the need for three separate versions of a model: 1) the narrative describing the problem, 2) the statistical model representing the problem, and 3) the code enabling inference written in a probabilistic programming language. Instead, causact users create one unified model, a generative DAG, using a visual representation.

Statement of Need

Bayesian data analysis mixes data with domain knowledge to quantify uncertainty in unknown outcomes. Its beautifully-simple theoretical underpinnings are deployed in three main steps ([Gelman et al., 2013](#)):

- **Modelling:** Joint probability distributions are specified to encode domain knowledge about potential data generating processes.
- **Conditioning:** Bayes rule is used to reallocate plausibility among the potential data generating processes to be consistent with both the encoded domain knowledge and the observed data. The conditioned model is known as the posterior distribution.
- **Validation:** Evidence is collected to see whether the specified model as well as the computational implementation of the model and conditioning process are to be trusted or not.

Algorithmic advances in the *conditioning* step of Bayesian data analysis have given rise to a new class of programming languages called probabilistic programming languages (PPLs). Practical and complex statistical models which are analytically intractable can now be solved computationally using inference algorithms. In particular, Markov Chain Monte Carlo (MCMC) algorithms ([Congdon, 2010](#); [Gelfand & Smith, 1990](#); [Gilks & Roberts, 1996](#)) handle arbitrarily large and complex models via highly effective sampling processes that quickly detect high-probability areas of the underlying distribution [Kruschke \(2014\)](#).

The causact package, presented in this paper, focuses on solving a three-language problem that occurs during Bayesian data analysis. First, there is the language of the domain expert which we refer to as the *narrative* of how data is generated. Second, there is the language of *math* where a statistical model, amenable to inference, is written. Lastly, there is the language of *code*, where a PPL language supports computational inference from a well-defined statistical model. The existence of these three languages creates friction as diverse stakeholders collaborate to yield insight from data; often mistakes get made in both communicating and translating between the three languages. Prior to causact, any agreed upon narrative of a

41 data-generating process must ultimately be modelled in code using an error-prone process
 42 where model misspecification, variable indexing errors, prior distribution omissions, and other
 43 mismatches between desired model and coded model go easily unnoticed.

44 To unify inference-problem narratives, the statistical models representing those narratives, and
 45 the code implementing the statistical models, *causact* introduces a modified visualization
 46 of *directed acyclic graphs* (DAGs), called the *generative DAG*, to serve as a more intuitive
 47 and collaborative interface into probabilistic programming languages and to ensure faithful
 48 abstractions of real-world data generating processes.

49 Modelling with Generative DAGs

50 Generative DAGs pursue two simultaneous goals. One goal is to capture the narrative by
 51 building a conceptual understanding of the data generating process that lends itself to statistical
 52 modelling. And two, gather all the mathematical elements needed for specifying a complete
 53 Bayesian model of the data generating process. Both of these goals will be satisfied by
 54 iteratively assessing the narrative and the narrative's translation into rigorous mathematics
 55 using *causact* functions.

56 Capturing the narrative in code uses some core *causact* functions like `dag_create()`, `dag_node`
 57 `e()`, `dag_edge()`, and `dag_plate()` with the chaining operator `%>%` used to build a DAG from
 58 the individual elements. `dag_render()` or `dag_greta()` are then used to visualize the DAG
 59 or run inference on the DAG, respectively. The simplicity with which generative DAGs are
 60 constructed belies the complexity of models which can be supported. For example, multi-level
 61 or hierarchical models are easily constructed as shown here in code for constructing and
 62 visualizing an oft-cited Bayesian example known as eight schools (Sturtz et al., 2005) whose
 63 data is included in *causact* (`causact::schoolsDF`). The example is a study of coaching effects
 64 on test scores where students from eight schools were put into coached and uncoached groups.

```
65 R> graph = dag_create() %>%
66 +   dag_node("Treatment Effect", "y",
67 +     rhs = normal(theta, sigma),
68 +     data = causact::schoolsDF$y) %>%
69 +   dag_node("Std Error of Effect Estimates", "sigma",
70 +     data = causact::schoolsDF$sigma,
71 +     child = "y") %>%
72 +   dag_node("Exp. Treatment Effect", "theta",
73 +     child = "y",
74 +     rhs = avgEffect + schoolEffect) %>%
75 +   dag_node("Population Treatment Effect", "avgEffect",
76 +     child = "theta",
77 +     rhs = normal(0, 30)) %>%
78 +   dag_node("School Level Effects", "schoolEffect",
79 +     rhs = normal(0, 30),
80 +     child = "theta") %>%
81 +   dag_plate("Observation", "i", nodeLabels = c("sigma", "y", "theta")) %>%
82 +   dag_plate("School Name", "school",
83 +     nodeLabels = "schoolEffect",
84 +     data = causact::schoolsDF$schoolName,
85 +     addDataNode = TRUE)
86 R> graph %>% dag_render()
```

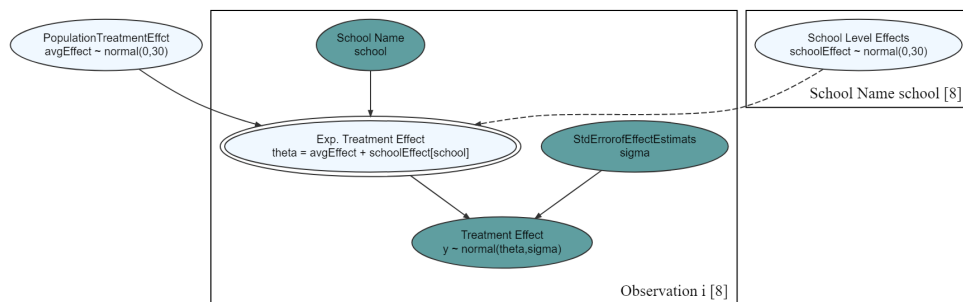


Figure 1: A generative DAG of the eight schools model.

87 **Figure 2** replicates **Figure 1** without math for less intimidating discussions with domain experts
88 about the model using the `shortLabel = TRUE` argument (shown below). `causact` does not
89 require a complete model specification prior to rendering the DAG, hence, `causact` facilitates
90 qualitative collaboration on the model design between less technical domain experts and the
91 model builder.

92 `R> graph %>% dag_render(shortLabel = TRUE)`

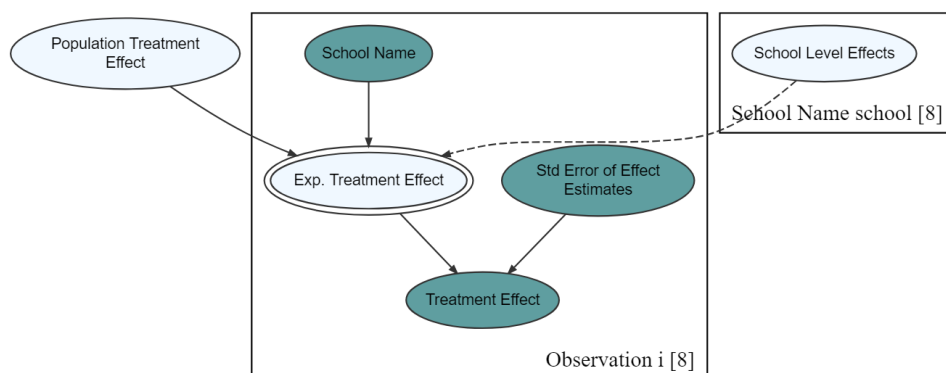


Figure 2: Hiding mathematical details to facilitate collaborations with domain experts.

93 All visualizations, including **Figure 1** and **Figure 2**, are created via `causact`'s calls to the
94 `DiagrammeR` package (Iannone, 2020). The `dag_diagrammer()` function can convert a `causact`
95 `t_graph` to a `dgr_graph` (the main object when using `DiagrammeR`) for further customizing of
96 a visualization using the `DiagrammeR` package.

97 Sampling from the posterior of the eight schools model (**Figure 1**) does not require a user to
98 write PPL code, but rather a user will simply pass the generative DAG object to `dag_greta()`
99 and then inspect the data frame of posterior draws:

```
100 R> library(greta) ## greta uses TensorFlow to get sample
101 R> drawsDF = graph %>% dag_greta()
102 R> drawsDF

103 # A tibble: 4,000 x 9
104   avgEffect schoolEffect_Sc~ schoolEffect_Sc~ schoolEffect_Sc~
105   <dbl>          <dbl>          <dbl>          <dbl>
106 1      0.102        40.1          3.59         -4.51
```

```

107 2      4.59      23.6      4.43      -26.3
108 3     -0.451     18.5     24.3      16.5
109 4     18.9      8.07    -26.3     -28.6
110 5     17.3     -5.83    -4.25    -26.2
111 6      1.97     42.7      2.25     12.6
112 7     12.7    -11.2      5.31    -16.5
113 8      9.11    -17.4      9.09    -12.7
114 9     -3.74     71.5      1.82     23.6
115 10    -2.43     48.2    -13.3      2.89

```

```

116 # ... with 3,990 more rows, and 5 more variables:
117 #   schoolEffect_School4 <dbl>, schoolEffect_School5 <dbl>,
118 #   schoolEffect_School6 <dbl>, schoolEffect_School7 <dbl>,
119 #   schoolEffect_School8 <dbl>

```

Behind the scenes, `causact` creates the model's code equivalent using the `greta` PPL, but this is typically hidden from the user. However, for debugging or further customizing a model, the `greta` code can be printed to the screen without executing it by setting the `mcmc` argument to `FALSE`:

```

124 R> graph %>% dag_greta(mcmc=FALSE)

125 sigma <- as_data(causact::schoolsDF$sigma) #DATA
126 y <- as_data(causact::schoolsDF$y) #DATA
127 school <- as.factor(causact::schoolsDF$schoolName) #DIM
128 school_dim <- length(unique(school)) #DIM
129 schoolEffect <- normal(mean = 0, sd = 30, dim = school_dim) #PRIOR
130 avgEffect <- normal(mean = 0, sd = 30) #PRIOR
131 theta <- avgEffect + schoolEffect[school] #OPERATION
132 distribution(y) <- normal(mean = theta, sd = sigma) #LIKELIHOOD
133 gretaModel <- model(avgEffect, schoolEffect) #MODEL
134 meaningfulLabels(graph)
135 draws <- mcmc(gretaModel) #POSTERIOR
136 drawsDF <- replaceLabels(draws) %>% as.matrix() %>%
137   dplyr::as_tibble() #POSTERIOR
138 tidyDrawsDF <- drawsDF %>% addPriorGroups() #POSTERIOR

```

The produced `greta` code is shown in the above code snippet. The code can be difficult to digest for some and exemplifies the advantages of working visually using `casuact`. The above code is also challenging to write without error or misinterpretation. Indexing is particularly tricky in PPL's with indexing based on meaningless numbers (e.g. 1,2,3,...). To facilitate quicker interpretation `causact` abbreviates posterior parameters using human-interpretable names.

The output of `dag_greta()` is in the form of a data frame of draws from the joint posterior. To facilitate a quick look into posterior estimates, the `dagp_plot()` function creates a simple visual of 90% credible intervals. It is the only core function that does not take a graph as its first argument. By grouping all parameters that share the same prior distribution and leveraging the meaningful parameter names constructed using `dag_greta()`, it allows for quick comparisons of parameter values.

```

151 R> drawsDF %>% dagp_plot()

```

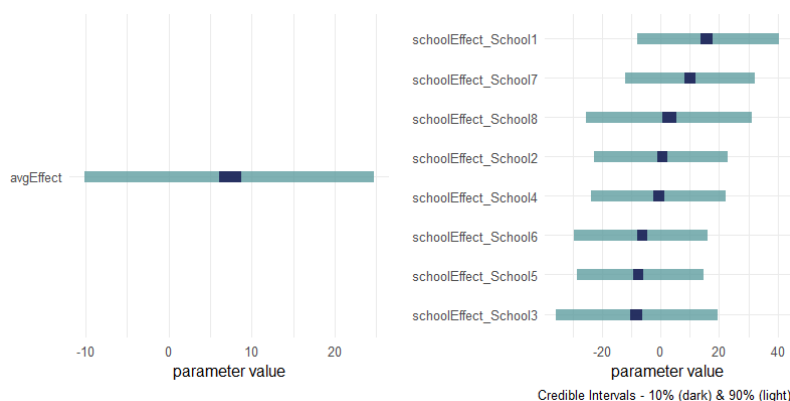


Figure 3: Credible intervals for the nine parameters of the eight schools model.

152 The code above makes the plot in [Figure 3](#). For further posterior plotting, users would
 153 make their own plots using `ggplot2` ([Wickham, 2016](#)), `ggdist` ([Kay, 2020](#)), or similar. For
 154 further model validation, including MCMC diagnostics, the user would use a package like
 155 `bayesplot` ([Gabry et al., 2019](#)) or `shinystan` ([Gabry, 2018](#)). For users who prefer to work
 156 with an `mcmc` object, they can extract the draws object after running the generated greta code
 157 from `dag_greta(mcmc=FALSE)` or find the object in the `cacheEnv` environment after running
 158 `dag_greta(mcmc=FALSE)` using `get("draws",envir = causact:::cacheEnv)`.

159 Comparison to Other Packages

160 By focusing on generative DAG creation as opposed to PPL code, `causact` liberates users
 161 from the need to learn complicated probabilistic programming languages. As such, it is similar
 162 in spirit to any package whose goal is to make Bayesian inference accessible without learning
 163 a PPL. In this vein, `causact` is similar to `brms` ([Bürkner, 2017](#)), `rstanarm` ([Goodrich et al., 2020](#)),
 164 and `rethinking` ([McElreath, 2020a](#)) - three R packages which leverage Stan ([Stan Development Team, 2021](#))
 165 for Bayesian statistical inference with MCMC sampling. Like the `rethinking` package which is tightly
 166 integrated with a textbook ([McElreath, 2020b](#)), a large motivation for developing `causact` was to make
 167 learning Bayesian inference easier. The package serves a central role in a textbook titled *A Business Analyst's Introduction to Business Analytics: Intro to Bayesian Business Analytics in the R Ecosystem*. ([Fleischhacker, 2020](#)).

170 Conclusion

171 The `causact` modelling syntax is flexible and encourages modellers to make bespoke models. The
 172 long-term plan for the `causact` package is to promote a Bayesian workflow that philosophically
 173 mimics the Principled Bayesian Workflow outlined by Betancourt ([2020](#)). The structure of a
 174 generative DAG is sure to be much more transparent and interpretable than most other modern
 175 machine learning workflows; this is especially true when models are made accessible to those
 176 without statistical or coding expertise. For this reason, generative DAGs can help facilitate
 177 effective communication between modelers and domain users both during the designing process
 178 of the models and when explaining the results returned by the models.

179 Acknowledgements

180 The Stan Development team has been inspirational for this work and has formed a wonderful
 181 Bayesian inference community around their powerful language. Additionally, the books of
 182 Kruschke ([2014](#)) and McElreath ([2020b](#)) are tremendous resources for learning Bayesian data

analysis and their pedagogy is aspirational. This work would not be possible without the greta dev team and special thanks to Nick Golding and Nick Tierney. Lastly, thanks to the University of Delaware students, MBAs and PhDs, who have contributed time, code, testing, and enthusiasm for this project from its beginning.

References

- Betancourt, M. (2020). *Towards a principled Bayesian workflow*. https://betanalpha.github.io/assets/case_studies/principled_bayesian_workflow.html.
- Bürkner, P.-C. (2017). Brms: An r package for Bayesian multilevel models using stan. *Journal of Statistical Software*, 80(1), 1–28.
- Congdon, P. D. (2010). *Applied Bayesian hierarchical methods*. CRC Press.
- Dillon, J. V., Langmore, I., Tran, D., Brevdo, E., Vasudevan, S., Moore, D., Patton, B., Alemi, A., Hoffman, M., & Saurous, R. A. (2017). Tensorflow distributions. *arXiv Preprint arXiv:1711.10604*.
- Fleischhacker, A. (2020). *A business analyst's introduction to business analytics: Intro to Bayesian business analytics in the r ecosystem*. Independently Published. ISBN: 9798667128175
- Gabry, J. (2018). *Shinystan: Interactive visual and numerical diagnostics and posterior analysis for Bayesian models*. <https://CRAN.R-project.org/package=shinystan>
- Gabry, J., Simpson, D., Vehtari, A., Betancourt, M., & Gelman, A. (2019). Visualization in Bayesian workflow. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, 182(2), 389–402.
- Gelfand, A. E., & Smith, A. F. (1990). Sampling-based approaches to calculating marginal densities. *Journal of the American Statistical Association*, 85(410), 398–409.
- Gelman, A., Carlin, J. B., Stern, H. S., Dunson, D. B., Vehtari, A., & Rubin, D. B. (2013). *Bayesian data analysis*. CRC press.
- Gilks, W. R., & Roberts, G. O. (1996). Strategies for improving MCMC. *Markov Chain Monte Carlo in Practice*, 6, 89–114.
- Golding, N. (2019). Greta: Simple and scalable statistical modelling in r. *Journal of Open Source Software*, 4(40), 1601.
- Goodrich, B., Gabry, J., Ali, I., & Brilleman, S. (2020). *Rstanarm: Bayesian applied regression modeling via Stan*. <https://mc-stan.org/rstanarm>
- Iannone, R. (2020). *DiagrammeR: Graph/network visualization*. <https://github.com/rich-iannone/DiagrammeR>
- Kay, M. (2020). *ggdist: Visualizations of distributions and uncertainty*. <https://doi.org/10.5281/zenodo.3879620>
- Kruschke, J. (2014). *Doing Bayesian data analysis: A tutorial with r, JAGS, and stan*.
- McElreath, R. (2020a). *Rethinking: Statistical rethinking book package*.
- McElreath, R. (2020b). *Statistical rethinking: A Bayesian course with examples in r and stan*. CRC press.
- Neal, R. M. (1993). *Probabilistic inference using markov chain monte carlo methods*. Department of Computer Science, University of Toronto Toronto, Ontario, Canada.
- Pfeffer, A. (2016). *Practical probabilistic programming*. Manning Publ.

- 225 Stan Development Team. (2021). *Stan Modeling Language User's Guide and Reference*
226 *Manual, Version 2.26*. <http://mc-stan.org/>
- 227 Sturtz, S., Ligges, U., & Gelman, A. (2005). R2WinBUGS: A package for running WinBUGS
228 from r. *Journal of Statistical Software, Articles*, 12(3), 1–16. [https://doi.org/10.18637/](https://doi.org/10.18637/jss.v012.i03)
229 [jss.v012.i03](https://doi.org/10.18637/jss.v012.i03)
- 230 Wickham, H. (2016). *ggplot2: Elegant graphics for data analysis*. Springer-Verlag New York.
231 ISBN: 978-3-319-24277-4

DRAFT