

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

26 April 2022

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 (Gelfand and Smith 1990; Gilks and Roberts 1996; Congdon 2010) 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 data-generating process must ultimately be modelled in code using an error-prone process where model misspecification, variable indexing errors, prior distribution omissions, and other mismatches between desired model and coded model go easily unnoticed.

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

Modelling with Generative DAGs

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

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

```
graph = dag_create() %>%
  dag_node("Treatment Effect", "y",
    rhs = normal(theta, sigma),
    data = causact::schoolsDF$y) %>%
  dag_node("Std Error of Effect Estimates", "sigma",
    data = causact::schoolsDF$sigma,
    child = "y") %>%
  dag_node("Exp. Treatment Effect", "theta",
    child = "y",
    rhs = avgEffect + schoolEffect) %>%
  dag_node("Population Treatment Effect", "avgEffect",
    child = "theta",
    rhs = normal(0, 30)) %>%
  dag_node("School Level Effects", "schoolEffect",
    rhs = normal(0, 30),
    child = "theta") %>%
  dag_plate("Observation", "i", nodeLabels = c("sigma", "y", "theta")) %>%
  dag_plate("School Name", "school",
    nodeLabels = "schoolEffect",
    data = causact::schoolsDF$schoolName,
    addDataNode = TRUE)
graph %>% dag_render()
```

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

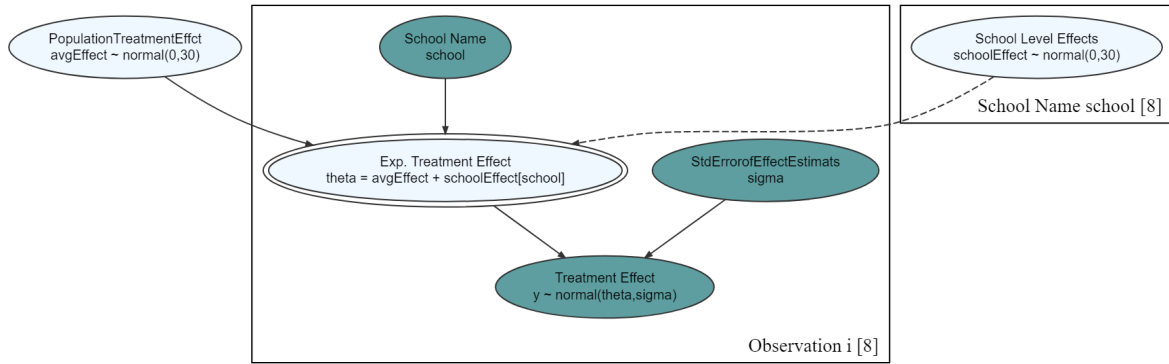


Figure 1: A generative DAG of the eight schools model.

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

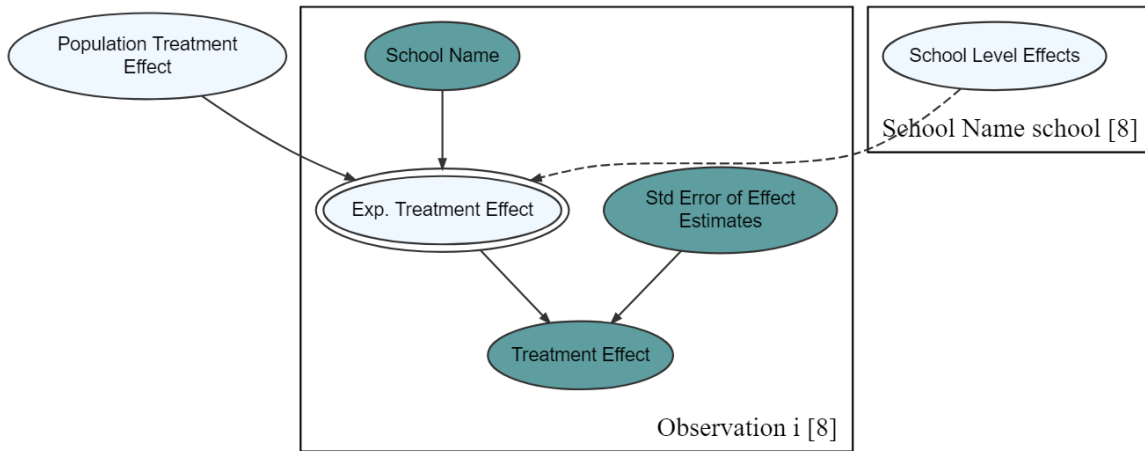


Figure 2: Hiding mathematical details to facilitate collaborations with domain experts.

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

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

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

```
## # A tibble: 4,000 x 9
##   avgEffect schoolEffect_Sc~ schoolEffect_Sc~ schoolEffect_Sc~
```

```
##           <dbl>           <dbl>           <dbl>           <dbl>
## 1      0.102           40.1           3.59           -4.51
## 2      4.59           23.6           4.43           -26.3
## 3     -0.451           18.5           24.3           16.5
## 4     18.9            8.07          -26.3           -28.6
## 5     17.3           -5.83          -4.25           -26.2
## 6      1.97           42.7           2.25           12.6
## 7     12.7          -11.2           5.31           -16.5
## 8      9.11          -17.4           9.09           -12.7
## 9     -3.74           71.5           1.82            23.6
## 10    -2.43           48.2          -13.3            2.89
## # ... with 3,990 more rows, and 5 more variables:
## #   schoolEffect_School4 <dbl>, schoolEffect_School5 <dbl>,
## #   schoolEffect_School6 <dbl>, schoolEffect_School7 <dbl>,
## #   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`:

```
graph %>% dag_greta(mcmc=FALSE)

## sigma <- as_data(causact::schoolsDF$sigma) #DATA
## y <- as_data(causact::schoolsDF$y) #DATA
## school <- as.factor(causact::schoolsDF$schoolName) #DIM
## school_dim <- length(unique(school)) #DIM
## schoolEffect <- normal(mean = 0, sd = 30, dim = school_dim) #PRIOR
## avgEffect <- normal(mean = 0, sd = 30) #PRIOR
## theta <- avgEffect + schoolEffect[school] #OPERATION
## distribution(y) <- normal(mean = theta, sd = sigma) #LIKELIHOOD
## gretaModel <- model(avgEffect,schoolEffect) #MODEL
## meaningfulLabels(graph)
## draws <- mcmc(gretaModel) #POSTERIOR
## drawsDF <- replaceLabels(draws) %>% as.matrix() %>%
##   dplyr::as_tibble() #POSTERIOR
## 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.

```
drawsDF %>% dagp_plot()
```

The code above makes the plot in Figure 3. For further posterior plotting, users would make their own plots using `ggplot2` (Wickham 2016), `ggdist` (Kay 2020), or similar. For further model validation, including

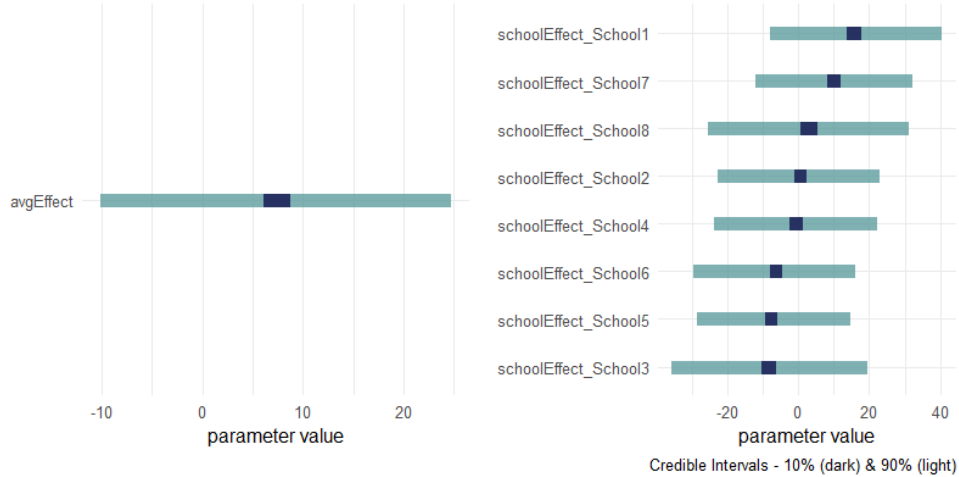


Figure 3: Credible intervals for the nine parameters of the eight schools model.

MCMC diagnostics, the user would use a package like **bayesplot** (Gabry et al. 2019) or **shinystan** (Gabry 2018). For users who prefer to work with an **mcmc** object, they can extract the **draws** object after running the generated **greta** code from `dag_greta(mcmc=FALSE)` or find the object in the **cacheEnv** environment after running `dag_greta(mcmc=FALSE)` using `get("draws",envir = causact:::cacheEnv)`.

Comparison to Other Packages

By focusing on generative DAG creation as opposed to PPL code, **causact** liberates users from the need to learn complicated probabilistic programming languages. As such, it is similar in spirit to any package whose goal is to make Bayesian inference accessible without learning a PPL. Perhaps the first such software was DoodleBUGS which provided a DAG-based graphical interface into WinBUGS (Lunn et al. 2000). In terms of leveraging more modern PPLs, **causact** is similar to **brms** (Bürkner 2017), **rstanarm** (Goodrich et al. 2020), and **rethinking** (McElreath 2020a) - three R packages which leverage **Stan** (Stan Development Team 2021) for Bayesian statistical inference with MCMC sampling. Like the **rethinking** package which is tightly integrated with a textbook (McElreath 2020b), a large motivation for developing **causact** was to make 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). As a point of contrast, the DAGitty package (Textor et al. 2017) also focuses on DAG creation/visualization, but DAGitty's intent is to help ensure consistency between the causal assumptions of a researcher and the dataset to which those assumptions should apply; DAGitty does not create PPL code for automating inference.

Conclusion

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

Acknowledgements

The **Stan** Development team has been inspirational for this work and has formed a wonderful Bayesian inference community around their powerful language. Additionally, the books of 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, Michael. 2020. “Towards a Principled Bayesian Workflow.” https://betanalpha.github.io/assets/case_studies/principled_bayesian_workflow.html.
- Bürkner, Paul-Christian. 2017. “Brms: An r Package for Bayesian Multilevel Models Using Stan.” *Journal of Statistical Software* 80 (1): 1–28. <https://doi.org/10.18637/jss.v080.i01>.
- Congdon, Peter D. 2010. *Applied Bayesian Hierarchical Methods*. CRC Press. <https://doi.org/10.1201/9781584887218>.
- Dillon, Joshua V, Ian Langmore, Dustin Tran, Eugene Brevdo, Srinivas Vasudevan, Dave Moore, Brian Patton, Alex Alemi, Matt Hoffman, and Rif A Saurous. 2017. “Tensorflow Distributions.” *arXiv Preprint arXiv:1711.10604*. <https://arxiv.org/abs/1711.10604>.
- Fleischhacker, Adam. 2020. *A Business Analyst’s Introduction to Business Analytics: Intro to Bayesian Business Analytics in the r Ecosystem*. Independently Published. <https://www.causact.com/>.
- Gabry, Jonah. 2018. *Shinystan: Interactive Visual and Numerical Diagnostics and Posterior Analysis for Bayesian Models*. <https://mc-stan.org/users/interfaces/shinystan>.
- Gabry, Jonah, Daniel Simpson, Aki Vehtari, Michael Betancourt, and Andrew Gelman. 2019. “Visualization in Bayesian Workflow.” *Journal of the Royal Statistical Society: Series A (Statistics in Society)* 182 (2): 389–402. <https://doi.org/10.1111/rssa.12378>.
- Gelfand, Alan E, and Adrian FM Smith. 1990. “Sampling-Based Approaches to Calculating Marginal Densities.” *Journal of the American Statistical Association* 85 (410): 398–409. <https://doi.org/10.2307/2289776>.
- Gelman, Andrew, John B Carlin, Hal S Stern, David B Dunson, Aki Vehtari, and Donald B Rubin. 2013. *Bayesian Data Analysis*. CRC press. <http://www.stat.columbia.edu/~gelman/book/>.
- Gilks, Walter R, and Gareth O Roberts. 1996. “Strategies for Improving MCMC.” *Markov Chain Monte Carlo in Practice* 6: 89–114.
- Golding, Nick. 2019. “Greta: Simple and Scalable Statistical Modelling in R.” *Journal of Open Source Software* 4 (40): 1601. <https://doi.org/10.21105/joss.01601>.
- Goodrich, Ben, Jonah Gabry, Imad Ali, and Sam Brilleman. 2020. “Rstanarm: Bayesian Applied Regression Modeling via Stan.” <https://mc-stan.org/rstanarm>.
- Iannone, Richard. 2020. *DiagrammeR: Graph/Network Visualization*. <https://github.com/rich-iannone/DiagrammeR>.
- Kay, Matthew. 2020. *ggdist: Visualizations of Distributions and Uncertainty*. <https://doi.org/10.5281/zenodo.3879620>.
- Kruschke, John. 2014. “Doing Bayesian Data Analysis: A Tutorial with r, JAGS, and Stan.” <https://sites.google.com/site/doingbayesiandataanalysis/>.
- Lunn, David J., Andrew Thomas, Nicky Best, and David Spiegelhalter. 2000. “WinBUGS - a Bayesian Modelling Framework: Concepts, Structure, and Extensibility.” *Statistics and Computing* 10: 325–37. <https://doi.org/10.1023/A:1008929526011>.
- McElreath, Richard. 2020a. *Rethinking: Statistical Rethinking Book Package*. <https://github.com/rmcelreath/rethinking>.
- . 2020b. *Statistical Rethinking: A Bayesian Course with Examples in R and Stan*. CRC press. <https://doi.org/10.1201/9780429029608>.
- Neal, Radford M. 1993. *Probabilistic Inference Using Markov Chain Monte Carlo Methods*. Department of Computer Science, University of Toronto Toronto, Ontario, Canada.

- Pfeffer, Avi. 2016. *Practical Probabilistic Programming*. Manning Publ.
- Stan Development Team. 2021. *Stan Modeling Language User's Guide and Reference Manual, Version 2.26*. <http://mc-stan.org/>.
- Sturtz, Sibylle, Uwe Ligges, and Andrew Gelman. 2005. "R2WinBUGS: A Package for Running WinBUGS from r." *Journal of Statistical Software, Articles* 12 (3): 1–16. <https://doi.org/10.18637/jss.v012.i03>.
- Textor, Johannes, Benito van der Zander, Mark S Gilthorpe, Maciej Liškiewicz, and George TH Ellison. 2017. "Robust causal inference using directed acyclic graphs: the R package 'dagitty'." *International Journal of Epidemiology* 45 (6): 1887–94. <https://doi.org/10.1093/ije/dyw341>.
- Wickham, Hadley. 2016. *Ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York. <https://ggplot2.tidyverse.org>.