# Soss

CHAD SCHERRER, RelationalAI

TAINE ZHAO, University of Tsukuba

We present Soss, a declarative probabilistic programming language embedded in the Julia language. Soss represents statistical models in terms of abstract syntax trees, and uses staged compilation for on-demand generation of "inference primitives" (random sampling, log-density, etc) without requiring casual users to worry about such details.

The approach taken by Soss makes it easy to extend to take advantage of other packages in the rapidly-growing Julia ecosystem. At the time of this writing, Soss users can choose from several inference back-ends and connect easily with larger systems SymPy and Gen.

## 1 INTRODUCTION

There are many approaches to building a Probabilistic Programming Language ("PPL"). Soss is distinct from most alternatives in a number of ways:

- Model specification in Soss is *declarative*; statements are represented not by the order entered by the user, but by the partial order given by their variable dependencies.
- Soss models are typically expressed in terms of *inference combinators* (like "For") that combine distributions in some way to arrive at a new distribution.
- Soss models are *first-class*; models can be passed as arguments to other models, or used in place of distribution functions like Normal within a model.
- Soss models are *function-like*; abstractly, a model may be considered a function from its arguments to a joint distribution. In particular, specification of "observed variables" is separated from model definition.
- Values and distributions in Soss are represented internally as *abstract syntax trees*, keeping the internal representation close to that given by the user for maximum flexibility.
- Soss uses *staged compilation* for inference primitives like rand and logpdf, via novel *generalized generated functions* from GeneralizedGenerated.jl.
- Soss supports *model transformations*, functions that take a model and return another model.
- Soss supports *symbolic simplifcation*, making it easy to inspect or manipulate a symbolic representation of the log-density, or to use it to generate optimized code.
- Soss is *extensible*; users can define new inference primitives or model transformations externally, and use them as if they had been included in Soss.

Finally, Soss uses Cockney rhyming slang to address the well-known challenge of naming things in computer science ("sauce pan" rhymes with "Stan"). **S**oss is **o**pen-**s**ource **s**oftware.

## 2 MODELING IN SOSS

Figure 1 shows a plate diagram and Soss implementation of a simplified two-way ANOVA model. The Soss model *m* has *arguments I* and *J*. Each line of a model is either a *sample statement* (v ~ rhs) or an *assign statement* (v = rhs). In either case, v must be a valid Julia variable name, and rhs must be a valid Julia expression. For a sample statement, rhs should also be "distribution-like", supporting logpdf and/or rand methods (depending on the inference algorithm to be used).

Authors' addresses: Chad Scherrer, RelationalAI, Seattle, WA, chad.scherrer@gmail.com; Taine Zhao, University of Tsukuba, Department of Computer Science.

$$i \in \{1, \cdots, I\}$$
$$j \in \{1, \cdots, J\}$$
$$\lambda \sim \text{Cauchy}_+(0, 1)$$
$$\sigma_I \sim \text{Normal}_+(0, \lambda)$$
$$\sigma_J \sim \text{Normal}_+(0, \lambda)$$
$$\alpha_i \sim \text{Normal}(0, \sigma_I)$$
$$\beta_j \sim \text{Normal}(0, \sigma_J)$$
$$\hat{y}_{ij} = \alpha_i + \beta_j$$
$$y_{ij} \sim \text{Normal}(\hat{y}_{ij}, \sigma)$$

(a)

(b)

```
m = @model I,J begin
    λ ~ HalfCauchy()
    σI ~ HalfNormal(λ)
    σJ ~ HalfNormal(λ)

    α ~ Normal(0, σI) |> iid(I)
    β ~ Normal(0, σJ) |> iid(J)
    yhat = α .+ β'

    σ ~ HalfNormal()
    y ~ For(I,J) do i,j
            Normal(yhat[i,j], σ)
        end
    end
```
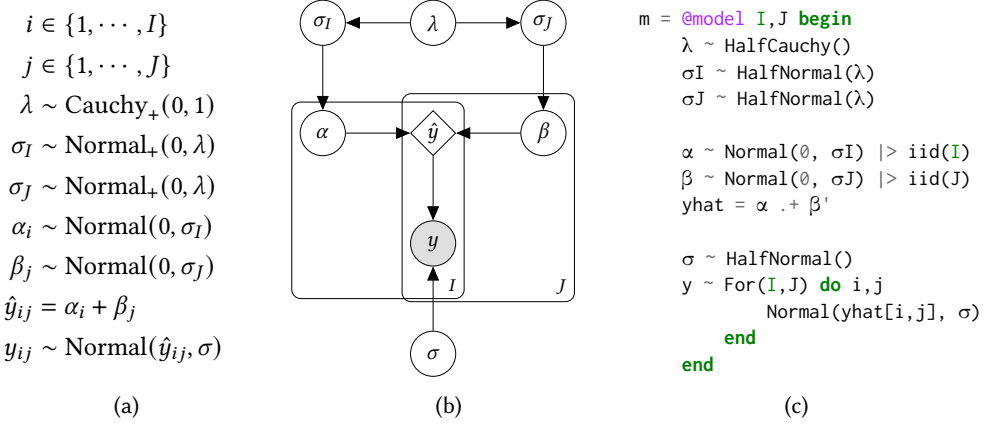
(c)

Fig. 1. A generative model (a), plate diagram (b), and Soss model (c) for a two-way analysis of variance model. Note the specification that $y$ will later be observed is not part of the model, but is instead given at inference time.

## 2.1 Joint Distributions

A Soss model is *function-like*, taking values for its arguments and returning a *joint distribution*. For example, given m from Figure 1, we can build the joint distribution d=m(I=2, J=3). This supports the usual rand and logpdf functions, as well as some others provided by Soss described in the following sections. For example,[1]

```
julia> d = m(I=2, J=3); rand(d) # or just rand(m(I=2, J=3))
(I = 2, J = 3, yhat = [-0.52 -1.5 0.0037; -0.81 -1.85 -0.28], σ = 2.5, λ = 1.3, σJ = 0.72, σI = 0.41,
 β = [-0.54, -1.5, -0.023], α = [0.026, -0.26], y = [-0.45 -2.0 2.2; -3.4 -3.9 2.3])
```

## 2.2 Inference Primitives

A given inference algorithm can be expressed in terms of a collection of functions on a model or joint distribution. For example, we often need to evaluate the log-density or its gradient or draw random samples. In Soss, these *inference primitives* generate model-specific code at execution time.

At this writing, the following inference primitive are available:

**rand** draws a random point or sample from a given distribution.

**particles** draws a *systematic sample*, as described in [Douc et al. 2005].

**logpdf** takes a distribution and a point, and returns the log-density at that point.

**xform** takes a distribution and some known data, and returns a bijection from $\mathbb{R}^n$ to the space of unknowns. This is useful for algorithms like Hamiltonian Monte Carlo [Neal 2011].

**symlogpdf** takes a model, and returns a SymPy [Meurer et al. 2017] representation of the log-density, with some additional Julia-based simplifications.

**codegen** calls symlogpdf, and uses the result to generate a more efficient version of logpdf.

In the above, a *distribution* refers to a joint distribution, while *point* refers to a named tuple in the space space of a given distribution.

cite MonteCarloMeasurements.jl and TransformVariables.jl

Users may add more

---

[1]This and future outputs are lightly edited for space, for example rounding floating point values.

## 2.3 Functions on Models

arguments sampled assigned parameters variables

   toposort poset digraph

## 2.4 Model Transformations

**Do** implements Pearl's `do` operator for causal intervention [Pearl 1995].[2]

**predictive** is like `Do`, but only returns the strongly-connected components containing the specified variables. This is useful for sampling from the posterior predictive distribution.

**markovBlanket** takes a model and a variable, and returns a model representing the Markov blanket at the given variable.

## 3 IMPLEMENTATION

```julia
struct Model{A,B,M}
    args  :: Vector{Symbol}
    vals  :: NamedTuple
    dists :: NamedTuple
    retn  :: Union{Nothing, Symbol, Expr}
end
```

## 4 PERFORMANCE

> Let's implement models from https://statisticalrethinkingjulia.github.io/MCMCBenchmarks.jl/latest/benchmarks/ for benchmark comparisons

> Can Soss really generate arbitrary code?

The ability to generate arbitrary Julia code makes it difficult to measure performance in Soss. When a performance bottleneck is found, developers or users can ask Soss to generate something different. This design places no a priori constraints on performance.

> Restructure my writing?

Besides the purely static code generation via regular macros happening in parsing time, Soss heavily uses a mechanism of "zero-cost" runtime code generation originated by Julia's generated functions [Bezanson et al. 2012], a.k.a staged functions in the original paper and earlier versions of Julia Programming Language.

The generated functions provide the capability of performing code generation during type inference, generating programs computed by the body of function(a.k.a, generator), once and only once for each combination of argument types.

Further, Julia enables type inference and compiler optimizations equivalent to non-runtime ones in runtime for the callsites of generated functions, hence we gain runtime code generation without losing performance.

To take advantage of this, Soss designs a system to encode sufficient information for generating actual codes, into the types, or objects that can be dispatched like types, which allows the use of generated functions here.

Unfortunately, in current stage, there're some implementation restrictions to Julia's generated functions, and generated functions cannot generate arbitrary Julia codes, where some advanced programming constructions are missing, such as generators and function-related stuffs like closures(including functions with no free variables), multiple dispatch, etc.

To address this, we introduced the works of easing the restrictions of Julia generated functions,

---

[2]Our Do operator is capitalized because do is a keyword in Julia.

> cite GG here?

which allow us to generate closures from generated functions, and make the programs sufficiently powerful.

Programs that can be generated by Soss are turing complete, as all constructs of Lambda Calculus can be generated.

## 5   EXTENDING SOSS

## 6   RELATED WORK

The idea of code generation and symbolic simplification in an embedded PPL goes back to *Passage* [Scherrer et al. 2012].

*Hakaru* [Narayanan et al. 2016] takes a more ambitious symbolic approach in a stand-alone PPL, allowing a wider variety of program transformations.

Soss began with a goal of representing models with continuous, fixed-dimensionality parameter spaces, inspired by *Stan* [Carpenter et al. 2017].

*Gen* [Cusumano-Towner et al. 2019] was developed independently from Soss, but takes a similar approach (and distinct from most PPLs) in its representation of a model as a function from its arguments to a "trace" (to use Gen's terminology). In both Soss and Gen's static DSL, this trace is a mapping from variable names to values. The similarity of these systems makes interoperability relatively straightforward, as demonstrated in *SossGen* (https://github.com/cscherrer/SossGen.jl).

*Turing* [Ge et al. 2018]

## ACKNOWLEDGMENTS

## REFERENCES

Jeff Bezanson, Stefan Karpinski, Viral B. Shah, and Alan Edelman. 2012. Julia: A Fast Dynamic Language for Technical Computing. arXiv:cs.PL/1209.5145

Bob Carpenter, Andrew Gelman, Matthew Hoffman, Daniel Lee, Ben Goodrich, Michael Betancourt, Marcus Brubaker, Jiqiang Guo, Peter Li, and Allen Riddell. 2017. Stan: A Probabilistic Programming Language. *Journal of Statistical Software* 76, 1 (2017).

Marco F. Cusumano-Towner, Feras A. Saad, Alexander K. Lew, and Vikash K. Mansinghka. 2019. Gen: A General-purpose Probabilistic Programming System with Programmable Inference. In *Proceedings of the 40th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI 2019)*. ACM, New York, NY, USA, 221–236. https://doi.org/10.1145/3314221.3314642

Randal Douc, Olivier Cappé, and Eric Moulines. 2005. Comparison of resampling schemes for particle filtering. *Image and Signal Processing and Analysis, 2005. ISPA 2005. Proceedings of the 4th International Symposium* 2005 (2005), 64–69. https://doi.org/10.1109/ispa.2005.195385 arXiv:cs/0507025

Hong Ge, Kai Xu, and Zoubin Ghahramani. 2018. Turing: a language for flexible probabilistic inference. In *International Conference on Artificial Intelligence and Statistics, AISTATS 2018, 9-11 April 2018, Playa Blanca, Lanzarote, Canary Islands, Spain*. 1682–1690. http://proceedings.mlr.press/v84/ge18b.html

Aaron Meurer, Christopher P. Smith, Mateusz Paprocki, Ondřej Čertík, Sergey B. Kirpichev, Matthew Rocklin, AMiT Kumar, Sergiu Ivanov, Jason K. Moore, Sartaj Singh, Thilina Rathnayake, Sean Vig, Brian E. Granger, Richard P. Muller, Francesco Bonazzi, Harsh Gupta, Shivam Vats, Fredrik Johansson, Fabian Pedregosa, Matthew J. Curry, Andy R. Terrel, Štěpán Roučka, Ashutosh Saboo, Isuru Fernando, Sumith Kulal, Robert Cimrman, and Anthony Scopatz. 2017. SymPy: symbolic computing in Python. *PeerJ Computer Science* 3 (Jan. 2017), e103. https://doi.org/10.7717/peerj-cs.103

Praveen Narayanan, Jacques Carette, Wren Romano, Chung-chieh Shan, and Robert Zinkov. 2016. Probabilistic inference by program transformation in Hakaru (system description). In *International Symposium on Functional and Logic Programming - 13th International Symposium, FLOPS 2016, Kochi, Japan, March 4-6, 2016, Proceedings*. Springer, 62–79. https://doi.org/10.1007/978-3-319-29604-3_5

Radford M. Neal. 2011. MCMC using hamiltonian dynamics. *Handbook of Markov Chain Monte Carlo* (2011), 113–162. https://doi.org/10.1201/b10905-6 arXiv:1206.1901

$$\langle\text{model}\rangle \models \texttt{@model }\langle\text{args}\rangle\texttt{ begin }\langle\text{statements}\rangle\texttt{ end}$$

$$\langle\text{args}\rangle \models \lambda \mid \langle\text{Symbol}\rangle \mid \langle\text{Symbol}\rangle, \langle\text{args}\rangle$$

$$\langle\text{statements}\rangle \models \langle\text{statement}\rangle \mid \langle\text{statement}\rangle \text{ } \textit{Julia Line Sep } \langle\text{statements}\rangle \mid \langle\text{statements}\rangle \langle\text{retn}\rangle$$

$$\langle\text{statement}\rangle \models \langle\text{assign}\rangle \mid \langle\text{sample}\rangle$$

$$\langle\text{assign}\rangle \models \langle\text{Symbol}\rangle \texttt{ = } \langle\text{Expr}\rangle$$

$$\langle\text{sample}\rangle \models \langle\text{Symbol}\rangle \texttt{ ~ } \langle\text{Measure}\rangle$$

$$\langle\text{retn}\rangle \models \texttt{return }\langle\text{Expr}\rangle$$

$$\langle\text{Symbol}\rangle \models \textit{Julia Symbol}$$

$$\langle\text{Expr}\rangle \models \textit{Julia Expr}$$

$$\langle\text{Measure}\rangle \models \textit{Probability measure (see text)}$$

Fig. 2. Backus-Naur Form representation for a user-specified model.

$$\ell = -3\log(2) - 5\log(\pi) - \sigma^2$$
$$- 4\log(\lambda) - 2\log(\lambda^2 + 1) - \frac{\sigma_I^2}{\lambda^2} - \frac{\sigma_J^2}{\lambda^2}$$
$$+ \sum_{i=1}^{I}\left(-\log(\sigma_I) - \frac{\log\pi}{2} - \frac{\log 2}{2} - \frac{\alpha_i^2}{2\sigma_I^2}\right)$$
$$+ \sum_{j=1}^{J}\left(-\log(\sigma_J) - \frac{\log\pi}{2} - \frac{\log 2}{2} - \frac{\beta_j^2}{2\sigma_J^2}\right)$$
$$+ \sum_{\substack{1 \le i \le I \\ 1 \le j \le J}}\left(-\log(\sigma) - \frac{\log\pi}{2} - \frac{\log 2}{2} - \frac{(y_{ij} - \hat{y}_{ij})^2}{2\sigma^2}\right)$$

(a) Before

$$\ell = -7.8 - 0.92I - 0.92J - 0.92IJ - \sigma^2$$
$$- 4\log\lambda - 2\log(\lambda^2 + 1) - IJ\log\sigma$$
$$- I\log\sigma_I - J\log\sigma_J - \frac{\sigma_I^2}{\lambda^2} - \frac{\sigma_J^2}{\lambda^2}$$
$$- \frac{1}{2\sigma_I^2}\sum_{i=1}^{I}\alpha_i^2 - \frac{1}{2\sigma_J^2}\sum_{j=1}^{J}\beta_j^2$$
$$- \frac{1}{2\sigma^2}\sum_{\substack{1 \le i \le I \\ 1 \le j \le J}}(y_{ij} - \hat{y}_{ij})^2$$

(b) After

Fig. 3. Symbolic log-density of $m$ before (a) and after (b) simplification steps. Reduction

Judea Pearl. 1995. Causal Diagrams of Empirical Research. *Biometrika* 82, 4 (1995), 669–688. https://escholarship.org/content/qt6gv9n38c/qt6gv9n38c.pdf

Chad Scherrer, Iavor Diatchki, Levent Erkök, and Matthew Sottile. 2012. Passage : A Parallel Sampler Generator for Hierarchical Bayesian Modeling. In *Neural Information Processing Systems workshop on Probabilistic Programming*. 1–4.

## A MODEL SYNTAX

## B SYMBOLIC SIMPLIFICATION

Before simplification:

After simplification:

## C EXTENDED EXAMPLE

Break this up! This is just a staging area

```
246    julia> using Revise, Soss, Random
247
248    julia> Random.seed!(1);
249
250    julia> m = @model I,J begin
251              λ ~ HalfCauchy()
252              σI ~ HalfNormal(λ)
253              σJ ~ HalfNormal(λ)
254              α ~ Normal(0, σI) |> iid(I)
255              β ~ Normal(0, σJ) |> iid(J)
256              yhat = α .+ β'
257              σ ~ HalfNormal()
258              y ~ For(I,J) do i,j
259                     Normal(yhat[i,j], σ)
260                  end
261          end;
262
263    julia> truth = rand(m(I=2,J=3)); pairs(truth);
264
265    julia> post = dynamicHMC(m(I=2,J=3), (y=truth.y,));
266
267    julia> postpred = predict(m(I=2,J=3), post) |> particles;
268
269    julia> postpred.yhat
270    2×3 Array{Particles{Float64,1000},2}:
271     0.241 ± 0.26  0.527 ± 0.3   -0.556 ± 0.38
272     0.22 ± 0.26   0.506 ± 0.32  -0.577 ± 0.37
273
274    julia> postpred.y
275    2×3 Array{Particles{Float64,1000},2}:
276     0.224 ± 0.58  0.529 ± 0.58  -0.58 ± 0.64
277     0.223 ± 0.56  0.483 ± 0.58  -0.58 ± 0.61
278
279    julia> m.args
280    2-element Array{Symbol,1}:
281     :I
282     :J
283
284    julia> pairs(m.vals)
285    pairs(::NamedTuple) with 1 entry:
286      :yhat => :(α .+ β')
287
288    julia> pairs(m.dists)
289    pairs(::NamedTuple) with 7 entries:
290      :λ  => :(HalfCauchy())
291      :σI => :(HalfNormal(λ))
292      :σJ => :(HalfNormal(λ))
293      :α  => :(Normal(0, σI) |> iid(I))
294      :β  => :(Normal(0, σJ) |> iid(J))
       :σ  => :(HalfNormal())
       :y  => :(For(I, J) do i, j...
```