

Philipp Gabler, BSc

Automatic Graph Tracking in Dynamic Probabilistic Programs via Source Transformations

Master's Thesis

to achieve the university degree of
Master of Science

submitted to
Graz University of Technology

Supervisor
Univ.-Prof. Dipl.-Ing. Dr. mont. Franz Pernkopf

Co-supervisor
Dipl.-Ing. Dr. Martin Trapp, BSc

Institute of Signal Processing and Speech Communication

Faculty of Electrical and Information Engineering

Graz, XXXX 2020

AFFIDAVIT

I declare that I have authored this thesis independently, that I have not used other than the declared sources/resources, and that I have explicitly indicated all material which has been quoted either literally or by content from the sources used. The text document uploaded to TUGRAZonline is identical to the present master's thesis.

Date

Signature

This work is licensed under a
[Creative Commons Attribution-ShareAlike 4.0 International License](https://creativecommons.org/licenses/by-sa/4.0/).



All code samples, unless otherwise noted or cited from other sources,
are also available under an [MIT license](#):

The MIT License (MIT)

Copyright (c) 2020 Philipp Gabler

Permission is hereby granted, free of charge, to any person obtaining a copy of this software and associated documentation files (the "Software"), to deal in the Software without restriction, including without limitation the rights to use, copy, modify, merge, publish, distribute, sublicense, and/or sell copies of the Software, and to permit persons to whom the Software is furnished to do so, subject to the following conditions:

The above copyright notice and this permission notice shall be included in all copies or substantial portions of the Software.

THE SOFTWARE IS PROVIDED "AS IS", WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT. IN NO EVENT SHALL THE AUTHORS OR COPYRIGHT HOLDERS BE LIABLE FOR ANY CLAIM, DAMAGES OR OTHER LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE, ARISING FROM, OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE OR OTHER DEALINGS IN THE SOFTWARE.

The ~~La~~^{TeX} source of this document is available at
<https://github.com/philpsgabler/master-thesis>
or upon request from the author*.

*pgabler@student.tugraz.at

ABSTRACT

This thesis presents a novel approach for the implementation of a tracking system to facilitate program analysis, based on program transformations. The approach is then applied to a specific problem in the field of probabilistic programming.

The main contribution is a general system for the extraction of rich computation graphs in the Julia programming language, based on a transformation of the intermediate representation (IR) used by the compiler. These graphs contain a slice of the whole recursive structure of any Julia program in terms of executed IR instructions. The system is flexible enough to be used for multiple purposes that require dynamic program analysis or abstract interpretation, such as automatic differentiation or dependency analysis.

The second part of the thesis describes the application of this graph tracking system to probabilistic programs written for Turing, a probabilistic programming system implemented as an embedded language within Julia. Through this, an executed Turing model can be analyzed, and the dependency structure of involved random variables be extracted from it. Given this structure, analytical Gibbs conditionals can be calculated for a large set of models and passed to Turing's inference mechanism, where they are used in Markov-Chain Monte Carlo samplers approximating the modelled distribution.

Contents

Notation	xi
1 Introduction	1
1.1 Related Work	2
2 Background	3
2.1 Bayesian Inference and MCMC methods	3
2.2 Probabilistic Programming	6
2.3 Compilation and Metaprogramming in Julia	6
2.4 Computation Graphs and Automatic Differentiation	6
3 Implementation of Dynamic Graph Tracking in Julia	7
3.1 Automatic Graph Tracking and Extended Wengert Lists	7
4 Graph Tracking in Probabilistic Models	9
4.1 Dependency Analysis in Dynamic Models	9
4.2 JAGS-Style Automatic Calculation of Gibbs Conditionals	9
4.3 Evaluation	9
5 Discussion	11
5.1 Future Work	11
Bibliography	13
List of Algorithms	15

Notation

$$\mathbb{P}[\Theta \in A \mid X = x]$$

Random variables and their realizations will usually be denoted by upper and lower case letters, respectively (with occasional exceptions for Greek variable names). Sets are also named by uppercase letters.

$$\mathbb{E}[X], \mathbb{V}_X[f(X, Y)]$$

Expectation and variance; if necessary, the variable with respect to which the moment is taken is indicated as a subscript.

$$\phi(x), f_Z(x)$$

Density functions are named using letters commonly used for functions, with an optional subscript indicating the random variable they belong to. Densities always come with implied base measures depending on the type of the random variable.

$$p(x, y \mid z)$$

The usual abuse of notation with the letter “p” standing for any density indicated by the names of the variables given to it is used when no confusion arises (in this case, $f_{X,Y|Z}$ is implied). A q may be used as well, mostly for proposal distributions or unnormalized densities.

$$X_i \sim \text{Normal}(\mu, \sigma)$$

The tilde notation for describing random variables is used throughout, without explicitly specifying dependence or independence, where understood from context. Named distributions that are not themselves random variables are spelled out in upright script.

$$Y \sim q(\cdot, X_{i-1})$$

The same notation is used when a random variable is specified to be sampled from a given, possibly unnormalized, density. In this context and elsewhere, the midpoint is employed to denote anonymous functions of one variable given by partial application.

$$y \mapsto p(x \mid y, z)$$

Anonymous functions are distinguished from function evaluation; this is crucial to differentiate between probability densities and likelihoods, for example.

$$\int p(x) \, dx = 1$$

Integrals over the whole domain of a density or measure are written as indefinite integrals, where the usage is clear.

$$[x, y, z] = \begin{bmatrix} x \\ y \\ z \end{bmatrix}$$

For consistency with Julia code, vectors (arrays of rank 1) are written in brackets, with elements separated by commas. Thereby, the form written in a row denotes a column vector; actual row vectors are written as transposed column vectors.

$$\Theta^{(k)} = [\Theta_1^{(k)}, \dots, \Theta_N^{(k)}]$$

Superscript indices in parentheses are used for series or sequences of variables, and subscript indices for components of multivariate variables.

$$z_{-i} = [z_1, \dots, z_{i-1}, z_{i+1}, \dots, z_N]$$

Negative indices denote all components of a variable without the negated one.

$$f.(x, 1) = [f(x_1, 1), \dots, f(x_N, 1)]$$

Function application with a period indicates vectorized application, as in Julia code^{*}: the function is applied over all elements of the input arrays individually, whereby arrays of lower rank or scalars are “broadcasted” along dimensions as necessary.

$$\text{f}(x) = \text{rand}(x)$$

Julia code (including identifiers mention in the text) is always typeset in typewriter font.

^{*}See <https://docs.julialang.org/en/v1/manual/functions/#man-vectorized-1>

1 Introduction

This chapter gives an overview over the scope of the thesis and existing approaches in the literature. A preliminary version of this work has already been presented in Gabler et al. (2019), which forms the basis of the introduction.

MANY METHODS in the field of machine learning work on computation graphs of programs that represent mathematical expressions. One example are forms of automatic differentiation (AD) which derive an “adjoint” expression from a expression that usually represents a loss function, to calculate its gradient (Gebremedhin; Walther 2020; Griewank; Walther 2008). Another one are message passing algorithms (Minka 2005), which use the graph as the basic data structure for the operation they perform: passing values between nodes, representing random variables that depend on each other (in fact, message passing generalizes various other methods, including AD). But also in more or less unrelated fields, such as program analysis or program transformation (cf. XXX), the same requirements might occur through the need to derive abstract graphs of program flow from a given program.

examples

abstract interpretation

There are several options how to provide the computation graph in question to an application, many of which are already established in the AD community (see Baydin et al. (2018) for a survey on AD methods). For one, graphs can be required to be written out explicitly by the user, by defining a custom input format or a library to build graphs “by hand” through an API (e.g., PyTorch (Paszke et al. 2017) or TensorFlow (Abadi et al. 2015)). Such APIs are called *operator overloading* in AD language, because they extend existing operations to additionally track the computation graph at runtime on so-called tapes or Wengert lists (Bartholomew-Biggs et al. 2000). This kind of tracking is dynamic, in the sense that a new tape is recorded for every execution. However, being implemented on a library level, it usually requires the programmer to use non-native constructs instead of language primitives, leading to cognitive overhead. Furthermore, there are additional runtime costs due to the separate interpretation of derivatives stored on the tape.

example

Alternatively, an implementation can allow the user to write out computations as a “normal” program in an existing programming language (or possibly a restricted subset of it), and use metaprogramming techniques to extract graphs from the input program. Such metaprograms, known under the name *source transformations*, can in turn operate on plain source code (cf. Tapenade (Tapenade developers 2019)), or on another, more abstracted notion used by the programming language infrastructure,

like the abstract syntax tree (AST), or an intermediate representation (IR). They operate on the syntactic structure of the whole program, during or before compilation. Unlike in operator overloading, it is hence possible to inspect and exploit control structures directly. This can lead to more efficient results, compared to operator overloading, since the transformation is done only once per program and eligible for compiler optimisations. Additionally, the user is not restricted to the domain specific language provided by a library, and can use regular language constructs, data structures, and custom functions rather freely. But in this approach, no records of the actual execution paths are constructed explicitly – purely static information is used only at compile time, and cannot be accessed for further analysis or transformation during execution.

IN A VARIETY of domains, though, the execution path of programs can drastically change at each run. Examples of this from machine learning are models with non-uniform data, such as parse trees (Socher et al. 2011) or molecular graphs (Bianucci et al. 2000), Bayesian nonparametric models (Hjort et al. 2010), or simply the occurrence of stochastic control flow in any probabilistic model. Such programs we call dynamic models. The lack of an explicit, unique graph structure makes it impossible, or at least difficult, to apply source transformation approaches on them. Operator overloading is the more direct way for supporting dynamic models, since it automatically records a new tape for each input. In fact, many state-of-the-art machine learning libraries are based on dynamic graphs using operator overloading in some form.

However, relying on operator overloading makes it impossible to take advantage of the benefits of source transformations, such as utilizing information about the control flow, integrating with optimizations at compile time, or exploiting the source model structure. The source transformation approach based on intermediate representations has recently gained popularity in machine learning.

1.1 RELATED WORK

2 Background

This chapter provides the background for the concepts used later in chapters 3 and 4. Initially, it gives a quick overview of Bayesian inference and probabilistic programming in general, necessary to understand the requirements and usual approaches of probabilistic programming systems.

Consequently, the machinery and language used to develop the graph tracking system forming the main part of the work are described. This consists firstly of a short introduction to graph tracking and source-to-source automatic differentiation, which contain many ideas and terminology that will be used later, and often provided inspiration. Secondly, the basic notions and techniques of the Julia compilation process as well as the language’s metaprogramming capabilities are described, which form the basis of the implementation.

2.1 BAYESIAN INFERENCE AND MCMC METHODS

Generative modelling is a technique for modelling phenomena based on the assumption that observables can be fully described through some stochastic process. When we assume this process to belong to a specified family of processes, the estimation of the “best” process is a form of learning: if we have a good description of how observations are generated, we can make summary statements about the whole population (descriptive statistics) or predictions about new observations. When observations come in pairs of independent and dependent variables, learning the conditional model of one given the other solves a regression or classification problem.

Within a Bayesian statistical framework, we assume that the family of processes used is specified by random variables related through conditional distributions with densities, which describe how the observables would be generated: some *unobserved variables* are generated from *prior distributions*, and the *observed data* are generated conditionally on the unobserved variables. The goal is to learn the *posterior distribution* of the parameters given the observations, which is a sort of “inverse” of how the problem is specified.

As an example, consider image classification: if we assume that certain percentages of an image data set picture cats and dogs, respectively, the distribution of these labels forms the prior. Given the information which kind of animal is depicted on it, an image can then be generated as a matrix of pixels based on a distribution of images conditioned on labels. The posterior distribution is then conditional distribution of

the label given an image. When we have this information, we can, for example, build a Bayesian classifier, by returning for a newly observed image that label which has the higher probability under the posterior.

This kind of learning is called Bayesian inference since, in the form of densities, the form of the model can be expressed using Bayes' theorem as the conditional distribution¹

$$\overbrace{p(\theta | x)}^{\text{posterior}} = \frac{\overbrace{p(x | \theta)}^{\text{likelihood}} \overbrace{p(\theta)}^{\text{prior}}}{p(x)}, \quad (2.1)$$

where x are the observed data, and θ are the unobserved parameters. The posterior represents the distribution of the unobserved variables as a combination of the prior belief updated by what has been observed (Congdon 2006). (In practice, not all of the unobserved variables have to be model parameters we are actually interested in; these can be integrated out).

Going beyond simple applications like the classifier mentioned above, handling the posterior gets difficult, though. Simply evaluating the posterior density $\theta \mapsto p(\theta | x)$ at single points is not enough in a Bayesian setting for usages such as prediction, parameter estimation, or evaluation of probabilities of continuous variables. The problem is that almost all of the relevant quantities depend on some sort of expectation over the posterior density, an integral of the form

$$\mathbb{E}[f(\theta) | X = x] = \int f(\theta) p(\theta | x) d\mu(\theta), \quad (2.2)$$

for a suitable f (with the base measure μ depending on the type of θ). This in turn involves calculating the normalizing marginal

$$p(x) = \int p(x, \tilde{\theta}) d\mu(\tilde{\theta}). \quad (2.3)$$

in equation 2.1 (also called the “evidence”).

When the distributions involved form a sufficiently “nice” combination, e.g., a conjugate pair, the integration can be performed analytically, since the posterior density has a closed form for a certain known distribution, or at least is a known integral. In general, however, this is not tractable, not even numerically, and approximations have to be made. Even for discrete variables, combinatorial explosion limits the applicability of simple summation.

DIFFERENT TECHNIQUES for posterior approximation are available: among them are distribution-based methods like message passing and variational inference for general graphical models. The methods described in this thesis, however, fall into the category of Monte Carlo methods, based on sampling. Their fundamental idea is

¹Note the abuse of notation regarding $p(\cdot)$; see page xi on notation.

Algorithm 1 General scheme for the Metropolis-Hastings algorithm.

1. Start from an arbitrary $\Theta^{(1)}$.
 2. For each $k \geq 1$:
 1. Propose $\hat{\Theta}^{(k)} \sim q(\Theta^{(k-1)}, \cdot)$.
 2. With probability $\alpha(\hat{\Theta}^{(k)}, \Theta^{(k-1)})$, set $\Theta^{(k)} = \hat{\Theta}^{(k)}$; else, keep $\Theta^{(k)} = \Theta^{(k-1)}$.
-

to derive, from a specified density, a sampling procedure with a consistent estimator for expectations:

$$I^{(k)}(f) \rightarrow \mathbb{E}[f(\Theta) | X = x], \quad \text{as } k \rightarrow \infty \quad (2.4)$$

in some appropriate stochastic convergence (usually convergence in probability is enough) (Vihola 2020, chapter 1). Most of these methods are defined in a form that samples a sequence of individual random variables $\Theta^{(k)}$, called a *chain*, for which a law of large numbers (LLN) holds:

$$I^{(k)}(f) = \frac{1}{n} \sum_{i=1}^n f(\Theta^{(k)}) \rightarrow \mathbb{E}[f(\Theta) | X = x] \quad (2.5)$$

When we can sample $\Theta^{(k)} \sim p(\cdot | x)$ exactly, they are i.i.d. and the LLN holds trivially; such samplers exist, but might also be difficult to derive or not possess good enough convergence properties (especially in high dimensions). Another large class of samplers is formed by *Markov Chain Monte Carlo* (MCMC) methods, which, instead of sampling exactly from the density, define $\Theta^{(k)}$ via a Markov chain: by choosing the transition the right way, the resulting Markov chain is ergodic with the target density as the unique stationary distribution. The advantage of MCMC methods is that they apply equally well to structurally complex models, and mostly treat the density as a black box, without requiring formal manipulations of a specific factorization.

overthink advantages

Frequently, MCMC methods use variations of the Metropolis-Hastings algorithm (MH), which requires the definition of a Markov transition kernel by means of two helper functions: a proposal distribution with density q , and an acceptance rate α , both depending on the old value in the chain; see algorithm 1. There exist many MH-based schemes with different properties and requirements: the classical random-walk Metropolis algorithm with Gaussian proposals, Reversible Jump MCMC, or gradient-informed methods like Metropolis Adjusted Langevin and Hamiltonian Monte Carlo (HMC). I refer to Vihola (2020, chapter 6) and Murphy (2012, chapters 24 and following) for an introduction to MCMC theory and algorithms.

nicer algorithm formatting

When we have a multi-component structure $\Theta = [\Theta_1, \dots, \Theta_N]$, a full transition kernel can be hard to find, and we can instead use a family of componentwise updates, given by conditional kernels q_i operating on only one component of Θ , with the

others fixed:

$$\begin{aligned}\hat{\Theta}_{-i}^{(k)} &= \Theta_{-i}^{(k-1)} \\ \hat{\Theta}_i^{(k)} &\sim q_i(\Theta_i^{(k-1)}, \cdot \mid \Theta_{-i}^{(k-1)})\end{aligned}\tag{2.6}$$

Components can be scalars or multivariate blocks, and the kernel may itself be any valid transition kernel. This allows one to freely mix different MCMC methods.

This “within-Gibbs” sampler bears its name because it is a generalization of the classical Gibbs sampling algorithm: the conditional densities $\Theta_i \mapsto p(\Theta_i \mid \Theta_{-i}, x)$ can directly be used as component proposals for a within-Gibbs sampler, leading to a cancelling acceptance rate of $\alpha \equiv 1$. This sampler has the advantage that it is in many cases rather easy to derive, even manually, from a given joint density.

2.2 PROBABILISTIC PROGRAMMING

Probabilistic programming is a means of describing probabilistic models through the syntax of a programming language. Probabilistic programs distinguish themselves from normal programs by the possibility of being sampled from conditionally, with some of the internal variables fixed to observed values. While probabilistic programming systems are often implemented as separate, domain-specific languages, they can also be embedded into “host” programming languages with sufficient syntactic flexibility. The latter is advantageous if one wants to use regular general-purpose programming constructs or interact with other functionalities of the host language.

See van de Meent et al. (2018) for a general introduction into the implementation of PPLs. Goodman; Stuhlmüller (2014) gives an in-depth overview of the implementation and usage of one specific, continuation-based implementation called WebPPL.

2.3 COMPILATION AND METAPROGRAMMING IN JULIA

Singer (2018)

2.4 COMPUTATION GRAPHS AND AUTOMATIC DIFFERENTIATION

3 Implementation of Dynamic Graph Tracking in Julia

3.1 AUTOMATIC GRAPH TRACKING AND EXTENDED WENGERT LISTS

4 Graph Tracking in Probabilistic Models

4.1 DEPENDENCY ANALYSIS IN DYNAMIC MODELS

4.2 JAGS-STYLE AUTOMATIC CALCULATION OF GIBBS CONDITIONALS

4.3 EVALUATION

5 Discussion

5.1 FUTURE WORK

Bla.¹

¹These ideas have already been informally described by me online at <https://github.com/phipsgabler/probability-ir>.

Bibliography

- Abadi, M. et al. (2015). “TensorFlow: Large-scale machine learning on heterogeneous systems”. Preliminary White Paper. Preliminary White Paper. URL: <https://www.tensorflow.org/> (visited on 2020-07-29).
- Bartholomew-Biggs, M. et al. (2000). “Automatic differentiation of algorithms”. In: *Journal of Computational and Applied Mathematics*. Numerical Analysis 2000. Vol. IV: Optimization and Nonlinear Equations 124.1, pp. 171–190. DOI: 10.1016/S0377-0427(00)00422-2. URL: <http://www.sciencedirect.com/science/article/pii/S0377042700004222> (visited on 2019-04-22).
- Baydin, A. G. et al. (2018). “Automatic differentiation in machine learning: a survey”. In: *Journal of Machine Learning Research* 18.153, pp. 1–43. URL: <http://jmlr.org/papers/v18/17-468.html>.
- Bianucci, A. M. et al. (2000). “Application of Cascade Correlation Networks for Structures to Chemistry”. In: *Applied Intelligence* 12.1, pp. 117–147. DOI: 10.1023/A:1008368105614.
- Congdon, P. (2006). *Bayesian statistical modelling*. 2nd ed. Wiley Series in Probability and Statistics. Chichester, England ; Hoboken, NJ: John Wiley & Sons. 573 pp.
- Gabler, P. et al. (2019). “Graph Tracking in Dynamic Probabilistic Programs via Source Transformations”. In: 2nd Symposium on Advances in Approximate Bayesian Inference. Vancouver. URL: <https://openreview.net/forum?id=r1eAFknEKr> (visited on 2020-07-07).
- Gebremedhin, A. H.; A. Walther (2020). “An introduction to algorithmic differentiation”. In: *WIREs Data Mining and Knowledge Discovery* 10.1, e1334. DOI: 10.1002/widm.1334. URL: <https://onlinelibrary.wiley.com/doi/abs/10.1002/widm.1334> (visited on 2020-07-29).
- Goodman, N. D.; A. Stuhlmüller (2014). *The Design and Implementation of Probabilistic Programming Languages*. URL: <http://dippl.org> (visited on 2019-10-15).
- Griewank, A.; A. Walther (2008). *Evaluating derivatives: principles and techniques of algorithmic differentiation*. 2nd ed. Philadelphia: Society for Industrial and Applied Mathematics. 438 pp.

- Hjort, N. L. et al. (2010). *Bayesian nonparametrics*. Cambridge Series in Statistical and Probabilistic Mathematics 28. Cambridge: Cambridge University Press.
- Minka, T. (2005). *Divergence Measures and Message Passing*. Technical Report MSR-TR-2005-173. Microsoft Research. URL: <https://www.microsoft.com/en-us/research/publication/divergence-measures-and-message-passing/> (visited on 2019-10-09).
- Murphy, K. P. (2012). *Machine learning: a probabilistic perspective*. Adaptive Computation and Machine Learning Series. Cambridge, MA: MIT Press. 1067 pp.
- Paszke, A. et al. (2017). “Automatic differentiation in PyTorch”. In: NIPS 2017 Workshop Autodiff.
- Singer, J. (2018). “Introduction”. In: *Single Static Assignment Book*. URL: <http://ssabook.gforge.inria.fr/latest/book-full.pdf> (visited on 2020-07-30).
- Socher, R. et al. (2011). “Parsing Natural Scenes and Natural Language with Recursive Neural Networks”. In: *Proceedings of the 28th International Conference on International Conference on Machine Learning*. ICML’11. USA: Omnipress, pp. 129–136. URL: <http://dl.acm.org/citation.cfm?id=3104482.3104499>.
- Tapenade developers (2019). *The Tapenade A.D. engine*. URL: <https://www-sop.inria.fr/tropics/tapenade.html> (visited on 2019-10-09).
- Van de Meent, J.-W. et al. (2018). “An Introduction to Probabilistic Programming”. In: arXiv: 1809.10756 [cs, stat]. URL: <http://arxiv.org/abs/1809.10756> (visited on 2019-03-08).
- Vihola, M. (2020). *Lectures on stochastic simulation*. University of Jyväskylä. URL: <http://users.jyu.fi/~mviholastochsim/>.

List of Algorithms

1	General scheme for the Metropolis-Hastings algorithm.	5
---	---	---

COLOPHON

This document was typeset using the pdf^{La}TeX typesetting system, with the memoir document class. The body text is set in 11 pt Linux Libertine, enhanced by the microtype package. Other fonts include Biolinum and Inconsolata.

The document source has been written in Emacs with AU^CTeX mode, using TeXworks as PDF viewer.