
Automatic Differentiation Equipped Variable Elimination for Sensitivity Analysis on Probabilistic Inference Queries

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Probabilistic Models are a natural framework for describing the stochastic relation-
2 ships between variables in a system to perform inference tasks, such as estimating
3 the probability of a specific set of conditions or events. In application it is often
4 appropriate to perform sensitivity analysis on a model, for example, to assess the
5 stability of analytical results with respect to the governing parameters. However,
6 typical programming language are cumbersome for encoding and reasoning with
7 complex models and current approaches to sensitivity analysis on probabilistic
8 models are not scalable, as they require repeated computation or estimation of the
9 derivatives of complex functions. To overcome these limitations, and to enable effi-
10 cient sensitivity analysis with respect to arbitrary model queries, e.g., $P(X|Y = y)$,
11 we propose to use Automatic Differentiation to extend the Probabilistic Program-
12 ming Language Figaro.

13 1 Introduction

14 In reasoning about an uncertain system, Probabilistic Models (PMs) can help understand how the
15 system will behave even though aspects of it are stochastic or unknown. For example, a Bayesian
16 network is a directed acyclic graph, which encodes local probability relationships, through the graph's
17 structure [4]. Variable's relationships are defined through their probability density functions (pdfs),
18 and the parameters that define them. Once the pdfs are in place, it is natural to ask questions on the
19 model such as: the probability of a specific set of conditions of the system, the most probable state of
20 variables, generating the likelihood of events, or asking these queries with evidence asserted. For
21 instance, an analyst could wish to perform a query on the system such as: what is the probability X is
22 true given we observe y (this can be written as $P(X|Y = y)$)? For a concrete case which will serve
23 as a running example for this paper, consider the system graphically shown in Fig 1(a). It has random
24 variables representation the occurrence of an earthquake and a burglary. These variables influence
25 whether a burglar alarm sounds, which in turn influences whether a neighbor calls. In this example
26 all the variables are Boolean. A query on this model could be: what is the probability that the alarm
27 is tripped, given a call was received? Solving these inference tasks can be done by exact methods
28 such as variable elimination, or approximate methods such as belief propagation, Monte Carlo, Gibbs
29 Sampling, etc.). We envision supporting a wide range of diverse and complex PM so we encode
30 our models via a probabilistic programming language. Probabilistic programming uses concepts from
31 programming languages to compactly encode complex PM [2]. Specifically, we use an open source
32 probabilistic programming language, Figaro [6].

33 In computing inference queries, there is no information embedded in the answer that provides
34 insight as to how *stable* the solution is. E.g., if a parameter that defines the network is changed
35 slightly, will the answer to the query change substantially? Questions of this type can be classified as

36 *sensitivity analysis*, which can be roughly described as the study of how the variation in the output of
 37 a mathematical model or system can be affected by variation in its inputs [8]. Suppose it is known to
 38 the analyst that there is significant uncertainty associated with a parameter x for this a probabilistic
 39 model; it reasonable to ask: how far off must the true value of x be from our estimated value, to
 40 change the query output by 5%? Referencing the running example, the probabilities that describe the
 41 relationship between variables are encoded via a conditional probability table. These probabilities
 42 *parametrize* the model, giving numerical values which are used in performing queries on the it. For
 43 example, suppose the parameter of interest is the probability that alarm is true, given earthquake
 44 is true, and burglary is false: $x = P(\text{alarm} = \text{true} | \text{earthquake} = \text{true}, \text{burglary} = \text{false})$.
 45 Traditional analysis of changes in the output of a model with respect to a specific parameter is
 46 possible [5], but manually repeating this analysis for all parameters is slow and laborious. More
 47 recent efforts explore means to compute node-to-node derivatives [1], but do not scale to more general
 48 inference tasks.

49 At the core of sensitivity analysis is the question of how much a function is changing with respect to
 50 changes to its input, captured by the mathematical notion of a gradient. Once gradients are obtained,
 51 they can be used to search for optimal parameter values which answer the sensitivity queries posed
 52 by the user. In our example the inference queries act as the function and the input are the parameters
 53 that define the PM. Computing these queries is often computationally expensive, and can be subject
 54 to variation due to approximations made to render calculations computationally feasible, even for
 55 parameters with constant value. A variety of methods exist to compute gradients, and we consider
 56 dual number enabled Automatic Differentiation (AD) [7] for its ability to compute exact derivatives
 57 in a computationally efficient manner [3]. There are several different mechanisms for AD, which
 58 compute gradients distinctly (e.g., forward accumulation, reverse mode), and we adopted a pure dual
 59 number approach for our prototypes. The semi-ring they form is analogous to the semi-ring used in
 60 the Variable Elimination solver used Figaro, allowing for a more straightforward implementation.

61 2 Approach

62 2.1 Sensitivity Query Example Problem

63 The motivation for developing a tool for performing sensitivity analysis, was answering questions
 64 such as: what is the minimum amount we can change parameter x by, such that the output of a query
 65 changes by ϵ . This can be expressed in the minimization problem

$$\begin{aligned} & \underset{\delta}{\operatorname{argmin}} \quad \delta \\ & \text{subject to} \quad |f(x_0 + \delta) - f(x_0)| > \epsilon. \end{aligned} \tag{1}$$

66 where $f(x)$ is the query, and x is the parameter of interest. We now extend our example with a
 67 sensitive analysis query and pose it as the minimization problem in Eq. 1. Let the query of interest
 68 be the probability of the alarm being triggered given the neighbor is calling, and the parameter of
 69 interest be the prior on an earthquake occurring: $x = P(\text{earthquake} = \text{true})$, $f(x) = P(\text{alarm} =$
 70 $\text{true} | \text{call} = \text{true})$. Note that for this analysis we consider all other parameters constant, so that f is
 71 only a function of x . In order to solve the minimization problem, we will a Newton's line search to
 72 update x :

$$x_{i+1} = x_i - \eta \frac{f(x_i)}{f'(x_i)} \tag{2}$$

73 where η is the learning rate for the search. The challenge now becomes computing the derivative
 74 $f'(x)$ at each step. Symbolic methods are appealing in that they are exact, but they suffer from
 75 expression bulge as models get complex. This quickly leads to intractable calculations as models
 76 become complex. Numerical derivatives are unappealing because computing the query $f(x)$ may be
 77 expensive. Worse, approximations and sampling methods entail that the results of successive queries
 78 of $f(x)$ may vary on the same scale as the true derivative, which creates "noisy" gradient estimates,

79 e.g., when evaluating $f(x + \delta) - f(x)$. Therefore, we use dual numbers to perform automatic
 80 differentiation to yield an exact derivative without computing the query value multiple times.¹

81 2.2 Extending Figaro’s Variable Elimination Algorithm with AD

82 Dual numbers form a semiring, extending real numbers by adjoining a new element d with the
 83 property $d^2 = 0$ (e.g., a dual number may be written as $a + bd$, where $a, b \in \mathbb{R}$). Dual numbers have
 84 the interesting property that when a dual number is passed into a function, the output contains the
 85 gradient value in its dual component: I.e.,

$$f(a + bd) = c + ed \implies f'(a) = e \quad (3)$$

86 This result depends on the property $d^2 = 0$ and the arithmetic associated with the dual number
 87 semiring. To perform inference on this PM we use Figaro’s Variable Elimination (VE) algorithm,
 88 but instead of standard arithmetic we compute over a dual number semiring. With the parameter of
 89 interest expressed as a dual number, the coefficient of the dual number in the output is the derivative
 90 of the query with respect to the parameter of interest.

91 To see this, consider our example (the query is the probability of the probability of the alarm being
 92 triggered given the neighbor is calling, and the parameter of interest is the prior on an earthquake
 93 occurring). Using the chain rule, we can write out the analytic expression the probability:

$$P(A^+|C^+) = \frac{1}{Z} \sum_{E,B} P(E)P(B)P(A|B,C)P(C|A) \quad (4)$$

94 where Z denotes $P(C^+)$, X^+ denotes X is true, and X^- denotes X is false. We can then now plug
 95 in a dual number $x + d$ for the parameter of interest.

$$P(A^+|C^+) = \frac{1}{Z} \sum_{E,B} (x + d)_E P(B)P(A|B,C)P(C|A) \quad (5)$$

96 where $(x + d)_E$ takes on the value of $P(E = true)$ or $P(E = false)$ depending on which value
 97 of E is being used in the summation. After the summation is performed we group terms by real
 98 components and dual components to get:

$$P(N^+) = \alpha + \beta d \quad (6)$$

99 Where α is the numerical answer to the inference query, and β with derivative of the query with
 100 respect to the parameter of interest.

101 We implemented this in Figaro by extending the initial *factors* produced by the VE algorithm with
 102 dual numbers (i.e., the factors produced from VE are similar to the individual terms produced in
 103 Eq. 4). These factors will have numerical values associated with them; for the factors relevant to
 104 the parameter of interest we give the dual coefficients a value of 1, and all others 0. Once the terms
 105 are assigned the correct dual numbers the VE algorithm runs as usual, but with arithmetic defined
 106 by the dual number semiring. The output contains a dual number (such as in Eq 6), which will
 107 contain both the query output in the natural number component and the gradient information in the
 108 dual component. We refer to this algorithm as Variable Elimination with Automatic Differentiation
 109 (VEAD).

110 2.3 Results

111 For the example in 2.1 we used the gradients obtained by the VEAD algorithm to execute a Newton’s
 112 line search as in Eq. 2. The results are depicted in Fig. 1(b), where one can see the parameter value

¹The computational cost is roughly a small constant factor more than the cost of computing the query

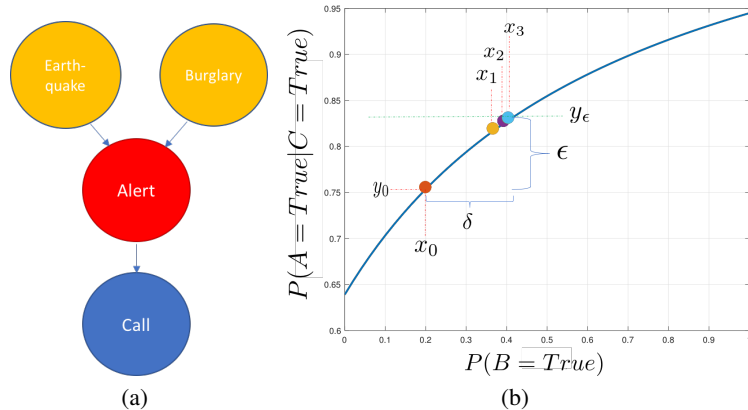


Figure 1: a) An example four node Bayesian model describing a scenario where a burglary or an earthquake influences whether an alarm goes off, which influences whether a neighbor calls. Each circle is a Boolean random variable characterized by conditional probability tables. b) An iterative search over the parameter space leads to the optimum value, solving the sensitivity query.

113 quickly converging to optimum value which caused the query to change by a target 5%. The method
 114 was also tested with a variety of PMs of varying complexity and results were verified by manually
 115 and numerically checking the gradients.

116 3 Conclusion and Future Work

117 We explored the usage of Automatic Differentiation to extend the probabilistic programming language
 118 Figaro, with a tool for efficiently calculating gradients of probabilistic inference queries. These
 119 gradients can be used to perform sensitivity analysis on these queries in order for an analyst to answer
 120 such questions as: how far off must the true value of a parameter of the system be from our estimated
 121 value, to change the query output by 5%? We have shown questions such as this can be answered
 122 with our framework utilizing a newtons line search to solve a for the optimum parameter value.

123 There are ample directions for future work. First, to validate the efficiency of our method, we would
 124 like to conduct a series of "wall clock tests" versus purely numerical means (even though these
 125 numerical derivatives may suffer from numerical instabilities, which the dual number approach for
 126 calculating does not). Secondly, we would like to explore augmenting powerful, approximate solvers
 127 such as Markov Chain Monte Carlo, Gibbs Sampling, or Importance Sampling, with the ability to
 128 ingest dual numbers in order to automatically compute gradients when computing queries, in the
 129 same we extended Variable Elimination to produce gradients.

130 References

- 131 [1] Chan, H., Darwiche, A.: Sensitivity analysis in bayesian networks: From single to multiple parameters. In:
 132 Proceedings of the 20th conference on Uncertainty in artificial intelligence. pp. 67–75. AUAI Press (2004)
- 133 [2] Goodman, N.D.: The principles and practice of probabilistic programming. In: Proceedings of the 40th
 134 annual ACM SIGPLAN-SIGACT symposium on Principles of programming languages. pp. 399–402. ACM
 135 (2013)
- 136 [3] Griewank, A., Walther, A.: Evaluating derivatives: principles and techniques of algorithmic differentiation.
 137 SIAM (2008)
- 138 [4] Koller, D., Friedman, N.: Probabilistic graphical models: principles and techniques. MIT press (2009)
- 139 [5] Laskey, K.B.: Sensitivity analysis for probability assessments in bayesian networks. IEEE Transactions on
 140 Systems, Man, and Cybernetics 25(6), 901–909 (1995)
- 141 [6] Pfeffer, A.: Creating and manipulating probabilistic programs with Figaro. In: 2nd International Workshop
 142 on Statistical Relational AI (2012)

- 143 [7] Rall, L.B.: Automatic differentiation: Techniques and applications (1981)
- 144 [8] Saltelli, A., Ratto, M., Andres, T., Campolongo, F., Cariboni, J., Gatelli, D., Saisana, M., Tarantola, S.:
145 Global sensitivity analysis: the primer. John Wiley & Sons (2008)