

# AUXILIARY VARIABLE MARKOV CHAIN MONTE CARLO METHODS

MATTHEW MCKENZIE GRAHAM



Doctor of Philosophy  
University of Edinburgh

2017

Matthew Mckenzie Graham: *Auxiliary variable Markov chain Monte Carlo methods*, 2017.  
Submitted for the degree of Doctor of Philosophy, University of Edinburgh.

SUPERVISOR:

Dr. Amos J. Storkey

## ABSTRACT

Inference, the process of drawing conclusions from evidence, is at the heart of the scientific method. Probability theory offers a consistent framework for performing inference in the presence of uncertainty, posing the inference problem as the task of computing expectations of functions of interest with respect to a probability distribution.

In complex models with a large number of unknown variables to be inferred, these expectations can be intractable to compute exactly. This has motivated the development of a large class of approximate inference methods which tradeoff a lack of exactness for greater computational tractability. Monte Carlo methods are one class of such techniques, with the integrals or summations across the whole state space approximated by summations over a finite number of randomly sampled points. This maps the inference problem to that of drawing samples from, often complex, high-dimensional, probability distributions.

## LAY SUMMARY

A lay summary is intended to facilitate knowledge exchange, public engagement and outreach. It should be in simple, non-technical terms that are easily understandable by a lay audience, who may be non-professional, non-scientific and outside the research area.

Abstracts, particularly in science, engineering, medicine and veterinary medicine, may be highly technical or contain scientific language that is not easily understandable to readers outside the research area. Therefore, the lay summary is intended as supplementary to the abstract.

## ACKNOWLEDGEMENTS

Put acknowledgments here.

Lorem ipsum at nusquam appellantur his, ut eos erant homero concluda-  
turque. Albucius appellantur deterruisset id eam, vivendum partiendo  
dissentiet ei ius. Vis melius facilisis ea, sea id convenire referrentur,  
takimata adolescens ex duo. Ei harum argumentum per. Eam vidit ex-  
erci appetere ad, ut vel zzril intellegam interpretaris.

Errem omnium ea per, pro Markov chain Monte Carlo ([MCMC](#)) con pop-  
ulo ornatus cu, ex qui dicant nemore melius. No pri diam iriure euismod.  
Graecis eleifend appellantur quo id. Id corpora inimicus nam, facer  
nonummy ne pro, kasd repudiandae ei mei. Mea menandri mediocrem  
dissentiet cu, ex nominati imperdiet nec, sea odio duis vocent ei. Tem-  
por everti appareat cu ius, ridens audiam an qui, aliquid admodum con-  
ceptam ne qui. Vis ea melius nostrum, mel alienum euripidis eu.

## DECLARATION

I declare that this thesis was composed by myself, that the work contained herein is my own except where explicitly stated otherwise in the text, and that this work has not been submitted for any other degree or professional qualification except as specified.

*Edinburgh, June 2017*

---

Matthew Mckenzie Graham

# CONTENTS

## FRONT MATTER

Abstract	1
Lay summary	2
Acknowledgements	3
Declaration	4
List of Figures	7
List of Tables	7
List of Abbreviations	7
1	PROBABILISTIC INFERENCE 9
1.1	Probability theory 9
1.1.1	Random variables 10
1.1.2	Joint and conditional probability 10
1.1.3	Probability densities 12
1.1.4	Transforms of random variables 14
1.1.5	Expectations 17
1.1.6	Conditional expectations and densities 18
1.2	Graphical models 19
1.2.1	Directed and undirected graphical models 20
1.2.2	Factor graphs 22
1.2.3	Computation graphs 27
1.3	Inference 32
1.3.1	Posterior expectations 33
1.3.2	Model evidence 33
2	APPROXIMATE INFERENCE 35
2.1	Deterministic approaches 35
2.1.1	Laplace's method 35
2.1.2	Variational inference 35
2.1.3	Expectation propagation 35
2.2	Stochastic approaches 35
2.2.1	Monte Carlo method 35
2.2.2	Rejection sampling 35
2.2.3	Importance sampling 35
2.2.4	Markov chain Monte Carlo 35

3	MARKOV CHAIN MONTE CARLO	37
3.1	Metropolis–Hastings	37
3.2	Gibbs sampling	37
3.3	Slice sampling	37
3.4	Hamiltonian Monte Carlo	37
4	THERMODYNAMIC METHODS	39
4.1	Simulated tempering	39
4.2	Parallel tempering	39
4.3	Tempered transitions	39
4.4	Annealed importance sampling	39
4.5	Path sampling	39
4.6	Adiabatic Monte Carlo	39
	BIBLIOGRAPHY	41



## LIST OF FIGURES

Figure 1.1	Directed and undirected graphical models.	20
Figure 1.2	Factor graph examples.	23
Figure 1.3	Boltzmann machine factor graph.	24
Figure 1.4	Hierarchical linear regression model factor graph.	25
Figure 1.5	Example computation graph.	28
Figure 1.6	Visualisation of reverse-mode automatic differentiation.	29
Figure 1.7	Factor node reparameterisation.	32
Figure 1.8	Hierarchical linear regression model stochastic computation graph.	33

## LIST OF TABLES

## LIST OF ABBREVIATIONS

**MCMC** Markov chain Monte Carlo

**wrt** with respect to



# 1

## PROBABILISTIC INFERENCE

Inference is the process of drawing conclusions from evidence. Much of our lives are spent making inferences about the world given our observations of it; in particular inference is a central aspect of the scientific process. Although deductive logic offers a framework for inferring conclusions from absolute statements of truth, it does not apply to the more typical real-world setting where the information we receive is subject to uncertainty.

To make inferences under conditions of uncertainty, we must instead turn to probability theory. Probabilities offer a consistent framework for quantifying the uncertainty in our beliefs about the world and making inferences given these beliefs. The output of the inference process is itself probabilistic, reflecting that the conclusions we make given uncertain information will themselves be subject to uncertainty.

In this chapter we will first introduce the probability notation we will use in the rest of this work, and state some basic results which will be important in the later chapters. We will introduce graphical models as a compact way of visualising structure in probabilistic models. Finally we will give a concrete definition of the probabilistic inference tasks that the methods presented in the rest of this thesis are aimed at computing (approximate) solutions to, and motivate why such approximate computational methods are needed.

### 1.1 PROBABILITY THEORY

A *probability space* is defined as a triplet  $(S, \mathcal{E}, P)$  where

- $S$  is the *sample space*, the set of all possible outcomes,
- $\mathcal{E}$  is the *event space*, a  $\sigma$ -algebra on  $S$ , defining all possible events (measurable subsets of  $S$ ),
- $P$  is the *probability measure*, a finite measure satisfying  $P(S) = 1$ , which specifies the probabilities of events in  $\mathcal{E}$ .

*The actual science of logic is conversant at present only with things either certain, impossible, or entirely doubtful, none of which (fortunately) we have to reason on. Therefore the true logic for this world is the calculus of probabilities*  
—James Clerk Maxwell

*Probability theory is nothing but common sense reduced to calculation.*  
— Pierre-Simon Laplace

*A  $\sigma$ -algebra,  $\mathcal{E}$ , on a set  $S$  is set of subsets of  $S$  with  $S \in \mathcal{E}$ ,  $\emptyset \in \mathcal{E}$  and which is closed under complement and countable unions and intersections.*

Kolmogorov's axioms:

1. *Non-negativity:*  
 $P(E) \geq 0 \forall E \in \mathcal{E}$ ,
2. *Normalisation:*  
 $P(S) = 1$ ,
3. *Countable additivity:*  
*for any countable set of disjoint events*  
 $\{E_i\}_i : E_i \in \mathcal{F} \forall i$ ,  
 $E_i \cap E_j = \emptyset \forall i \neq j$ ,  
 $P(\cup_i E_i) = \sum_i P(E_i)$ .

Given this definition of a probability space, Kolmogorov's axioms [12] can be used to derive a measure-theoretic formulation of probability theory. The probability of an event  $E \in \mathcal{E}$  is defined as the measure of that event  $P(E)$ . Two events  $A, B \in \mathcal{E}$  are said to be *independent* if  $P(A \cap B) = P(A)P(B)$ .

A measure-theoretic approach has the advantage of providing a unified treatment for describing probabilities on both finite and infinite sample spaces. Although alternative derivations of the laws of probability from different premises such as Cox's theorem [5, 6] have been proposed, modern extensions of this work result in a calculus of probabilities that is equivalent to Kolmogorov's [17], with the differences mainly being in the philosophical interpretations of probabilities.

### 1.1.1 Random variables

If  $(X, \mathcal{F})$  and  $(Y, \mathcal{G})$  are two measurable spaces, a function  $f : X \rightarrow Y$  is measurable if  $f^{-1}(E) \in \mathcal{F} \forall E \in \mathcal{G}$ .

When modelling real-world processes, rather than considering events as subsets of an abstract sample space, it is usually more helpful to consider *random variables* which represent quantities in the model of interest. A random variable  $x : S \rightarrow X$  is defined as a measurable function from the sample space to a measurable space  $(X, \mathcal{F})$ .

The Borel  $\sigma$ -algebra  $\mathcal{B}(\mathbb{R})$  is the smallest  $\sigma$ -algebra on  $\mathbb{R}$  which contains all open real intervals.

Often  $X$  is the reals,  $\mathbb{R}$ , and  $\mathcal{F}$  is the Borel  $\sigma$ -algebra on the reals,  $\mathcal{B}(\mathbb{R})$ , in which case we will refer to a *real random variable*. It is also common to consider cases where  $X$  is a real vector space,  $\mathbb{R}^D$ , and  $\mathcal{F} = \mathcal{B}(\mathbb{R}^D)$  - in this case we will term the resulting random variable a *random vector* and use the notation  $\mathbf{x} : S \rightarrow X$ . A final special case is when  $X$  is countable and  $\mathcal{F}$  is the power set  $\mathcal{P}(X)$  in which case we will refer to  $x$  as a *discrete random variable*.

If  $(X, \mathcal{F})$  and  $(Y, \mathcal{G})$  are two measurable spaces,  $\mu$  a measure on these spaces and  $f : X \rightarrow Y$  a measurable function, the pushforward measure  $\mu_f$  satisfies  $\mu_f(A) = \mu \circ f^{-1}(A) \forall A \in \mathcal{G}$ .

Due to the definition of a random variable as a measurable function, we can define a pushforward measure on a random variable  $x$

$$P_x(A) = P \circ x^{-1}(A) = P(\{s \in S : x(s) \in A\}) \quad \forall A \in \mathcal{F}. \quad (1.1)$$

The measure  $P_x$  specifies that the probability of the event that the random variable  $x$  takes a value in a measurable set  $A \in \mathcal{F}$  is  $P_x(A)$ .

### 1.1.2 Joint and conditional probability

Often we will jointly define multiple random variables on the same probability space. Let  $(S, \mathcal{E}, P)$  be a probability space and  $x : S \rightarrow X$ ,

$y : S \rightarrow Y$  be two random variables with corresponding  $\sigma$ -algebras  $\mathcal{F}$  and  $\mathcal{G}$ . Then the *joint probability* of  $x$  and  $y$  is defined as

$$P_{x,y}(A, B) = P(x^{-1}(A) \cap y^{-1}(B)) \quad \forall A \in \mathcal{F}, B \in \mathcal{G}. \quad (1.2)$$

The joint probability is related to the probabilities  $P_x$  and  $P_y$  by

$$P_{x,y}(A, Y) = P_x(A), P_{x,y}(X, B) = P_y(B) \quad \forall A \in \mathcal{F}, B \in \mathcal{G}. \quad (1.3)$$

In this context  $P_x$  and  $P_y$  are referred to as *marginals* of the joint.

The two random variables are said to be independent if and only if

$$P_{x,y}(A, B) = P_x(A)P_y(B) \quad \forall A \in \mathcal{F}, B \in \mathcal{G}. \quad (1.4)$$

Also useful is the definition of *conditional probability*

$$P(A | B) = \frac{P(A \cap B)}{P(B)} \quad \forall A \in \mathcal{E}, B \in \mathcal{E} : P(B) \neq 0. \quad (1.5)$$

Correspondingly, the conditional probabilities of random variables  $P_{x|y}$  and  $P_{y|x}$  can likewise be defined as satisfying

$$P_{x,y}(A, B) = P_{x|y}(A | B) P_y(B) = P_{y|x}(B | A) P_x(A) \quad (1.6)$$

$$\forall A \in \mathcal{F}, B \in \mathcal{G} : P_{x,y}(A, B) \neq 0,$$

*In Kolmogorov's probability theory, (1.5) is given as an additional definition distinct from the basic axioms. In alternatives such as the work of Cox [5, 6] and de Finetti [8], conditional probabilities are instead viewed as a primitive.*

which is sometimes referred to as the product rule.

An implication of (1.6) is what is often termed *Bayes' theorem*

$$P_{x|y}(A | B) = \frac{P_{y|x}(B | A) P_x(A)}{P_y(B)} \quad \forall A \in \mathcal{F}, B \in \mathcal{G} : P_y(B) \neq 0, \quad (1.7)$$

which will be of key importance in the later discussion of inference.

The definition in (1.2) of the joint probability of a pair of random variables can be extended to arbitrarily large collections of random variables. Similarly conditional probabilities can be defined for collections of multiple jointly dependent random variables, with the product rule given in (1.6) generalising to a combinatorial number of possible factorisations of the joint probability. Graphical models offer a convenient way of representing the dependencies between large collections of random variables and any resulting factorisation structure in their joint probability, and will be discussed later in this chapter in section 1.2.

## 1.1.3 Probability densities

So far we have ignored how the probability measure  $P$  is defined and by consequence the probability of a random variable.

*A measure on  $X$  is  $\sigma$ -finite if  $X$  is a countable union of finite measure sets.*

The Radon–Nikodym theorem guarantees that for a pair of  $\sigma$ -finite measures  $\mu$  and  $\nu$  on a measurable space  $(X, \mathcal{F})$  where  $\nu$  is absolutely continuous with respect to  $\mu$ , then there is a unique (up to  $\mu$ -null sets) measurable function  $f : X \rightarrow [0, \infty)$  termed a *density* such that

$$\nu(A) = \int_A f \, d\mu \quad \forall A \in \mathcal{F}. \quad (1.8)$$

*If  $\mu$  and  $\nu$  are measures on a measurable space  $(X, \mathcal{F})$  then  $\nu$  has absolute continuity wrt to  $\mu$  if  $\forall A \in \mathcal{F}$ ,  $\mu(A) = 0 \Rightarrow \nu(A) = 0$ .*

The density function  $f$  is also termed the *Radon-Nikodym derivative* of  $\nu$  with respect to  $\mu$ , denoted  $\frac{d\nu}{d\mu}$ . Density functions therefore represent a convenient way to define a probability measure with respect to an appropriate base measure. It can also be shown that if  $f = \frac{d\nu}{d\mu}$  and  $g$  is a measurable function that

$$\int_X g \, d\nu = \int_X g f \, d\mu, \quad (1.9)$$

which we will use later when discussing calculation of expectations.

For real random variables, an appropriate base measure is usually the *Lebesgue measure*,  $\lambda$ , on  $\mathbb{R}$ . The probability  $P_x$  of a real random variable  $x$  can then be defined via a *probability density*  $p_x : \mathbb{R} \rightarrow [0, \infty)$  by

$$P_x(A) = \int_A p_x \, d\lambda = \int_A p_x(x) \, dx \quad \forall A \in \mathcal{B}(\mathbb{R}). \quad (1.10)$$

Analogously for a random vector  $\mathbf{x}$  with density  $p_{\mathbf{x}} : \mathbb{R}^D \rightarrow [0, \infty)$  with respect to the  $D$ -dimensional Lebesgue measure  $\lambda^D$

$$P_{\mathbf{x}}(A) = \int_A p_{\mathbf{x}} \, d\lambda^D = \int_A p_{\mathbf{x}}(\mathbf{x}) \, d\mathbf{x} \quad \forall A \in \mathcal{B}(\mathbb{R}^D). \quad (1.11)$$

The notation in the second equalities in (1.10) and (1.11) uses a convention that will be used throughout this thesis that integrals without an explicit measure are with respect to the Lebesgue measure.

*The counting measure  $\#$  is defined as  $\#(A) = |A|$  for all finite  $A$  and  $\#(A) = +\infty$  otherwise.*

For discrete random variables, an appropriate base measure is instead the *counting measure*,  $\#$ . The probability of a discrete random variable

is then defined via a probability density  $p_x : X \rightarrow [0, 1]$  by

$$P_x(A) = \int_A p_x d\# = \sum_{x \in A} p_x(x) \quad \forall A \in \mathcal{P}(X). \quad (1.12)$$

The co-domain of a probability density  $p_x$  for a discrete random variable is restricted to  $[0, 1]$  due to the non-negativity and normalisation requirements for the probability measure  $P_x$ , with  $\sum_{x \in X} p_x(x) = 1$ . Commonly for the case of a discrete random variable, the density  $p_x$  is instead referred to as a *probability mass function*, with density reserved for real random variables. We will however use *probability density* in both cases in keeping with the earlier definition of a density with respect to a base measure, this avoiding difficulties when defining joint probabilities on a mixture of real and discrete random variables.

The joint probability  $P_{x,y}$  of a pair of random variables  $x$  and  $y$  with co-domains the measurable spaces  $(X, \mathcal{F})$  and  $(Y, \mathcal{G})$  respectively, can be defined via a joint probability density  $p_{x,y} : X \times Y \rightarrow [0, \infty)$  by

$$P_{x,y}(A, B) = \int_{A \times B} p_{x,y} d(\mu_x \times \mu_y) \quad \forall A \in \mathcal{F}, B \in \mathcal{G}, \quad (1.13)$$

where  $\mu_x \times \mu_y$  represents the product measure of two appropriate base measures  $\mu_x$  and  $\mu_y$ , e.g.  $\mu_x = \lambda$  and  $\mu_y = \#$  if  $x$  is a real random variable and  $y$  is a discrete random variable.

When dealing with random variables, we will usually only specify the co-domain of the random variable(s) and a (joint) probability density, with the base measure being implicitly defined as the Lebesgue measure for real random variables (or vectors), counting measure for discrete random variables and an appropriate product measure for a mix of random variables. Similarly we will usually neglect to explicitly define the probability space  $(S, \mathcal{E}, P)$  which the random variable(s) map from. In this case we will typically use the loose notation  $x \in X$  to mean a random variable  $x$  with co-domain  $X$ .

This less explicit but more succinct probability notation in terms of random variables and densities is common in the machine learning and computational statistics literature and will generally be preferred to improve readability. The underlying measure-theoretic basis of these concepts will however be important for some of the upcoming definitions in this chapter and some of the derivations later in the thesis.

*If  $(X_1, \mathcal{F}_1, \mu_1)$  and  $(X_2, \mathcal{F}_2, \mu_2)$  are two measure spaces, the product measure  $\mu_1 \times \mu_2$  on a measurable space  $(X_1 \times X_2, \mathcal{F}_1 \otimes \mathcal{F}_2)$  is defined as satisfying  $(\mu_1 \times \mu_2)(A_1 \times A_2) = \mu_1(A_1)\mu_2(A_2)$   $\forall A_1 \in \mathcal{F}_1, A_2 \in \mathcal{F}_2$ .*

## 1.1.4 Transforms of random variables

It is common to define a random variable via a transform of another. Let  $x$  be a random variable with co-domain the measurable space  $(X, \mathcal{F})$ . Further let  $(Y, \mathcal{G})$  be a second measurable space and  $\phi : X \rightarrow Y$  a measurable function between the two spaces. If we define  $y = \phi \circ x$  then analogously to our original definition of  $P_x$  as the pushforward measure of  $P$  under the measurable function defining  $x$ , we can define  $P_y$  in terms of  $P_x$  as

$$P_y(A) = P_x \circ \phi^{-1}(A) = P_x(\{x \in X : \phi(x) \in A\}) \quad \forall A \in \mathcal{G}, \quad (1.14)$$

i.e. the probability of the event  $y \in A$  is equal to the probability of  $x$  being in the pre-image under  $\phi$  of  $A$ . To calculate probabilities of transformed random variables therefore we will therefore need to be able to find the pre-images of values of the transformed variable.

If the probability  $P_x$  is defined by a probability density  $p_x$  with respect to a measure  $\mu_x$ , we can also in some cases find a density  $p_y$  on the transformed variable  $y = \phi(x)$  with respect to a (potentially different) measure  $\mu_y$  which can be used to calculate the probability  $P_y$ ,

$$P_y(A) = \int_{\phi^{-1}(A)} p_x d\mu_x = \int_A p_y d\mu_y \quad \forall A \in \mathcal{G}. \quad (1.15)$$

For random variables with countable co-domains where the integral in (1.15) corresponds to a sum, a  $p_y$  satisfying (1.15) is simple to identify. If  $x$  is a discrete random variable with probability density  $p_x$  with respect to the counting measure, then  $y = \phi(x)$  will necessarily also be a discrete random variable. Applying (1.15) for  $p_x = \frac{dP_x}{d\#}$  we have that

$$\begin{aligned} \int_{\phi^{-1}(A)} p_x(x) d\#(x) &= \sum_{x \in \phi^{-1}(A)} p_x(x) = \sum_{y \in A} \sum_{x \in \phi^{-1}(y)} p_x(x) \\ &= \int_A \sum_{x \in \phi^{-1}(y)} p_x(x) d\#(y) \quad \forall A \in \mathcal{G}. \end{aligned} \quad (1.16)$$

We can therefore define  $p_y = \frac{dP_y}{d\#}$  in terms of  $p_x$  as

$$p_y(y) = \sum_{x \in \phi^{-1}(y)} p_x(x) \quad \forall y \in Y. \quad (1.17)$$



In the special case that  $\phi$  is bijective we have that

$$p_Y(y) = p_X \circ \phi^{-1}(y) \quad \forall y \in Y. \quad (1.18)$$

For transformations of real random variables and vectors, the situation is more complicated as we need to account for any local contraction or expansion of space by the map  $\phi$ . Let  $X = \mathbb{R}^M$  and  $Y = \mathbb{R}^N$  with  $N \leq M$ ,  $N, M \in \mathbb{N}$ . We will need a result from geometric measure theory, the *co-area formula* [7]. Let  $g$  be an  $L^1$  integrable function and  $\phi : X \rightarrow Y$  a Lipschitz map. Then the co-area formula states that

$$\int_X g(\mathbf{x}) J_\phi(\mathbf{x}) d\lambda^M(\mathbf{x}) = \int_Y \int_{\phi^{-1}(\mathbf{y})} g(\mathbf{x}) d\mathcal{H}^{M-N}(\mathbf{x}) d\lambda^N(\mathbf{y}) \quad (1.19)$$

where  $\mathcal{H}^D$  is the  $D$ -dimensional Hausdorff measure and  $J_\phi : X \rightarrow [0, \infty)$  is the Jacobian determinant defined as

$$J_\phi(\mathbf{x}) = \left| \frac{\partial \phi}{\partial \mathbf{x}} \frac{\partial \phi}{\partial \mathbf{x}}^\top \right|^{\frac{1}{2}} \quad \forall \mathbf{x} \in X. \quad (1.20)$$

Now let  $\mathbf{x}$  be a random vector with co-domain the measurable space  $(X, \mathcal{B}(\mathbb{R}^M))$  and define  $\mathbf{y} = \phi \circ \mathbf{x}$  as a random vector with co-domain the measurable space  $(Y, \mathcal{B}(\mathbb{R}^N))$  with  $\phi : X \rightarrow Y$  a Lipschitz map as above. Let  $Z = \{\mathbf{x} \in X : J_\phi(\mathbf{x}) = 0\}$  and require that  $P_X(Z) = 0$ . Then for  $A \in \mathcal{B}(\mathbb{R}^N)$  define an  $L^1$  integrable function  $g$  as

$$g(\mathbf{x}) = \begin{cases} \mathbb{1}_A \circ \phi(\mathbf{x}) p_X(\mathbf{x}) J_\phi(\mathbf{x})^{-1} & \forall \mathbf{x} \in X \setminus Z \\ 0 & \forall \mathbf{x} \in Z \end{cases}. \quad (1.21)$$

Integrating  $g(\mathbf{x}) J_\phi(\mathbf{x})$  over  $\mathbf{x} \in X$  we have that

$$\int_X g(\mathbf{x}) J_\phi(\mathbf{x}) d\lambda^M(\mathbf{x}) = \int_{X \setminus Z} \mathbb{1}_A \circ \phi(\mathbf{x}) p_X(\mathbf{x}) d\lambda^M(\mathbf{x}) \quad (1.22)$$

$$= \int_X \mathbb{1}_A \circ \phi(\mathbf{x}) dP_X(\mathbf{x}) \quad (1.23)$$

$$= \int_{\phi^{-1}(A)} dP_X(\mathbf{x}) = P_Y(A). \quad (1.24)$$

The equality between first and second lines comes from the requirement  $P_X(Z) = 0$ , with the Lebesgue integrals of a function over two

*The  $D$ -dimensional Hausdorff measure  $\mathcal{H}^D$  on  $\mathbb{R}^N$  for  $D \in \mathbb{N}$ ,  $0 < D < N$  formalises a measure of the ‘volume’ of  $D$ -dimensional submanifolds of  $\mathbb{R}^N$  - e.g. for  $D = 1$  it corresponds to the length of a curve in  $\mathbb{R}^N$ . Additionally  $\mathcal{H}^N = \lambda^N$  and  $\mathcal{H}^0 = \#$ .*

sets which differ by only a zero-measure set equal. Now applying the co-area formula (1.19) to the left-hand side gives

$$\int_Y \int_{\phi^{-1}(\mathbf{y})} g(\mathbf{x}) d\mathcal{H}^{M-N}(\mathbf{x}) d\lambda^N(\mathbf{y}) = P_Y(A). \quad (1.25)$$

Therefore we can define a density  $p_Y = \frac{dP_Y}{d\lambda^N}$  satisfying (1.15) as

$$p_Y(\mathbf{y}) = \int_{\phi^{-1}(\mathbf{y})} p_X(\mathbf{x}) J_\phi(\mathbf{x})^{-1} d\mathcal{H}^{M-N}(\mathbf{x}) \quad \forall \mathbf{y} \notin \phi(Z). \quad (1.26)$$

For the special case of a dimension-preserving map  $\phi$  with  $N = M$  the integral in (1.26) is with respect to  $\mathcal{H}^0$  which is equivalent to the counting measure  $\#$ . In this case  $J_\phi(\mathbf{x}) = \left| \frac{\partial \phi}{\partial \mathbf{x}} \right|$  and we therefore get

$$p_Y(\mathbf{y}) = \sum_{\mathbf{x} \in \phi^{-1}(\mathbf{y})} p_X(\mathbf{x}) \left| \frac{\partial \phi}{\partial \mathbf{x}} \right|^{-1} \quad \forall \mathbf{y} \notin \phi(Z). \quad (1.27)$$

Under the further restriction that  $\phi$  is bi-Lipschitz, i.e. it is bijective and Lipschitz in both directions, we recover the more commonly presented multidimensional change of variables formula

$$p_Y(\mathbf{y}) = p_X \circ \phi^{-1}(\mathbf{y}) \left| \frac{\partial \phi^{-1}}{\partial \mathbf{y}} \right| \quad \forall \mathbf{y} \in Y. \quad (1.28)$$

In both of the cases considered, we have seen that if the function  $\phi$  the random variable  $\mathbf{x}$  is mapped through is bijective, the resulting expression for the density on the mapped random variable  $\mathbf{y}$  is simpler in the sense that the pre-image  $\phi^{-1}(\mathbf{y})$  of a point  $\mathbf{y} \in Y$  is itself a point and so we do not need to integrate or sum over points in the pre-image which will often be difficult to do analytically.

Bijectivity is a very limiting condition however, with many models involving non-bijective transformations of random variables. Later in this thesis we will see that methods used for defining the more general forms for calculating the density of a transformed variable are key to proposed methods for performing inference in generative models defined by complex, non-dimension preserving and non-bijective transformations of random variables.

## 1.1.5 Expectations

A fundamental operation when working with probabilistic models is computing expectations of random variables. Let  $(S, \mathcal{E}, P)$  be a probability space, and  $x : S \rightarrow X$  a random variable on this space. Then the *expected value of  $x$*  is defined as

$$\mathbb{E}[x] = \int_S x(s) dP(s). \quad (1.29)$$

Often it will be more convenient to express expectations in terms of the probability  $P_x$  instead. If  $f : S \rightarrow X$  is a measurable function and  $\mu$  a measure on  $S$  then the integral with respect to the pushforward measure  $\mu_f$  of an integrable function  $g$  satisfies

$$\int_X g(x) d\mu_f(x) = \int_S g \circ f(s) d\mu(s). \quad (1.30)$$

If we take  $g$  as the identity map we therefore have that

$$\mathbb{E}[x] = \int_X x dP_x(x). \quad (1.31)$$

If  $P_x$  is given by a density  $p_x = \frac{dP_x}{d\mu}$  then using (1.9) we also have

$$\mathbb{E}[x] = \int_X x p_x(x) d\mu(x), \quad (1.32)$$

which is often the form used for computation.

A further useful implication of (1.30) is what is sometimes termed the *Law of the unconscious statistician*. Let  $x : S \rightarrow X$  be a random variable,  $\phi : X \rightarrow Y$  a measurable function and define  $y = \phi \circ x$ . Then the expected value of  $y$  is

$$\mathbb{E}[y] = \int_S y(s) dP(s) = \int_S \phi \circ x(s) dP(s) = \int_X \phi(x) dP_x(x), \quad (1.33)$$

i.e. it can be calculated by integrating  $\phi$  with respect to  $P_x$ . This means we can calculate expectations of a transformed random variable  $y = \phi(x)$  without needing to use the change of variables formulae from Section 1.1.4 to explicitly calculate the probability  $P_y$  (or density  $p_y$ ) and with a relatively weak condition of measurability on  $\phi$ .

## 1.1.6 Conditional expectations and densities

A related concept, and one which will be key in our discussion of inference, is conditional expectation. Let  $(S, \mathcal{E}, P)$  be a probability space,  $(X, \mathcal{F})$  and  $(Y, \mathcal{G})$  two measurable spaces and  $x : S \rightarrow X$  and  $y : S \rightarrow Y$  two random variables. Then the *conditional expectation of  $x$  given  $y$* , is defined as a measurable function  $\mathbb{E}[x | y] : Y \rightarrow X$  satisfying

$$\int_{y^{-1}(A)} x(s) dP(s) = \int_A \mathbb{E}[x | y](y) dP_y(y) \quad \forall A \in \mathcal{G}. \quad (1.34)$$

$\mathbb{E}[x | y]$  is guaranteed to be uniquely defined almost everywhere in  $Y$  by (1.34), i.e. up to  $P_y$ -null sets. As a particular case where  $A = Y$  we recover what is sometimes termed the *Law of total expectation*

$$\int_S x dP = \int_S \mathbb{E}[x | y] \circ y dP \implies \mathbb{E}[x] = \mathbb{E}[\mathbb{E}[x | y] \circ y]. \quad (1.35)$$

We can also motivate a definition of conditional density in terms of conditional expectation. Assume a joint density  $p_{x,y} = \frac{dP_{x,y}}{d(\mu_x \times \mu_y)}$  exists and has marginal density  $p_y = \frac{dP_y}{d\mu_y}$ . Then for all  $A \in \mathcal{G}$

$$\int_{y^{-1}(A)} x(s) dP(s) = \int_S x(s) \mathbb{1}_A \circ y(s) dP(s) \quad (1.36)$$

$$= \int_{X \times Y} x \mathbb{1}_A(y) dP_{x,y}(x, y) \quad (1.37)$$

$$= \int_A \int_X x p_{x,y}(x, y) d\mu_x(x) d\mu_y(y). \quad (1.38)$$

Define  $g : Y \rightarrow X$  as

$$g(y) = \begin{cases} \int_X x \frac{p_{x,y}(x, y)}{p_y(y)} d\mu_x(x) & \forall y \in Y : p_y(y) > 0 \\ 0 & \forall y \in Y : p_y(y) = 0. \end{cases} \quad (1.39)$$

Then from (1.38) we have that for all  $A \in \mathcal{G}$

$$\int_{y^{-1}(A)} x(s) dP(s) = \int_A g(y) p_y(y) d\mu_y(y) = \int_A g(y) dP_y(y). \quad (1.40)$$

The definition of  $g$  in (1.39) therefore satisfies the definition of conditional expectation in (1.34) and is uniquely defined up to a  $P_Y$ -null set. Therefore if  $p_{x,y}$  and  $p_Y$  can be defined we have that

$$\mathbb{E}[x | y](y) = \int_X x p_{x|y}(x | y) d\mu_x(x) \quad \forall y \in Y : p_Y(y) > 0 \quad (1.41)$$

where the *conditional density of  $x$  given  $y$* ,  $p_{x|y}$ , is defined as

$$p_{x|y}(x | y) = \frac{p_{x,y}(x, y)}{p_Y(y)} \quad \forall x \in X, y \in Y : p_Y(y) > 0 \quad (1.42)$$

which can be seen to be analogous to the definition of conditional probability in (1.5). Note the definition of conditional expectation in (1.34) was not dependent on a joint density  $p_{x,y}$  being defined and so is more general than (1.41).

## 1.2 GRAPHICAL MODELS

When working with probabilistic models involving large numbers of random variables, it will often be the case that not all the variables are jointly dependent on each other but that instead there are more local conditional relationships between them. Graphical models, which use graphs to describe the dependencies between random variables, are a useful framework for visualising the structure in complex probabilistic models and for giving a graph-theoretic basis for establishing the dependence between sets of random variables.

*Graphical models =  
statistics × graph  
theory × computer  
science  
—Zoubin Ghahramani*

Central to all graphical models is the concept of conditional independence. Let  $(S, \mathcal{E}, P)$  be a probability space and  $x : S \rightarrow X$ ,  $y : S \rightarrow Y$  and  $z : S \rightarrow Z$  be three random variables with corresponding  $\sigma$ -algebras,  $\mathcal{F}_x$ ,  $\mathcal{F}_y$  and  $\mathcal{F}_z$  respectively. Following from our earlier definition of (unconditional) independence of random variables in (1.4), we say that  $x$  and  $y$  are *conditionally independent given  $z$* , denoted  $x \perp y | z$ , if

$$\mathbb{E}[\mathbb{1}_A \circ x \mathbb{1}_B \circ y | z] = \mathbb{E}[\mathbb{1}_A \circ x | z] \mathbb{E}[\mathbb{1}_B \circ y | z] \quad \forall A \in \mathcal{F}_x, B \in \mathcal{F}_y, \quad (1.43)$$

holds almost everywhere with respect to  $P_z$ . If a joint density on the random variables exists, a sufficient condition for  $x \perp y | z$  is that the conditional density  $p_{x,y|z}$  factorises as

$$p_{x,y|z}(x, y | z) = p_{x|z}(x | z) p_{y|z}(y | z) \quad \forall x \in X, y \in Y, z \in Z. \quad (1.44)$$

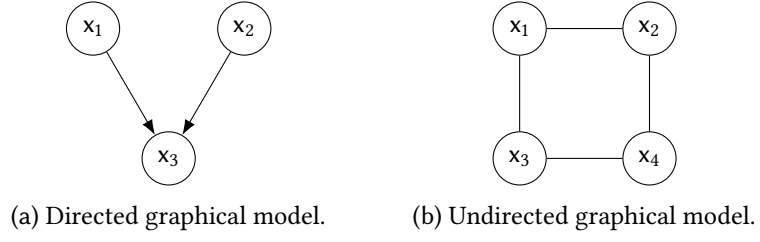


Figure 1.1: Examples of directed and undirected graphical models. Circular nodes represent random variables in the model, with edges between them indicating dependencies between variables.

This definition can be naturally extended to conditional independence when conditioning on more than one random variable, for example

$$v \perp x \mid y, z \implies p_{v,x \mid y,z}(v, x \mid y, z) = p_{v \mid y,z}(v \mid y, z) p_{x \mid y,z}(x \mid y, z) \quad (1.45)$$

### 1.2.1 Directed and undirected graphical models

Several different graphical frameworks have been proposed for representing conditional independency relationships (and other information) in probabilistic models.

*Directed graphical models* [14], also known as *Bayesian networks*, represent probabilistic models as *directed acyclic graphs* (i.e. a directed graph in which there are no directed cycles), with the nodes in the graph representing random variables in the model and the edges of the graph defining a factorisation of the joint density over these variables into a product of conditional and marginal densities. In particular a conditional density factor is included for each node with parents (on the node random variable value given the parent variable values) and a marginal density factor for each root node without any parents.

An example directed graphical model for three random variables,  $x_1$ ,  $x_2$  and  $x_3$ , is shown in Figure 1.1a. The graph implies that the joint density can be factorised as

$$p_{x_1, x_2, x_3}(x_1, x_2, x_3) = p_{x_3 \mid x_1, x_2}(x_3 \mid x_1, x_2) p_{x_1}(x_1) p_{x_2}(x_2). \quad (1.46)$$

Note that this factorisation would not be valid for all joint densities on the three variables; in particular we have that  $x_1$  and  $x_2$  are (unconditionally) independent and so that the joint density  $p_{x_1, x_2}$  can be written as the product of the two marginals  $p_{x_1}$  and  $p_{x_2}$ .

Directed graphical models are a natural way of specifying *generative models* - i.e. probabilistic models which can be used to generate simulated observable quantities. Typically the factorisation specified by a directed graphical model gives a straightforward method to generate values from the joint density via *ancestral sampling*.

An alternative formalism for graphically representing probabilistic models is that of *undirected graphical models* [11], which are also known as *Markov random fields*. As with directed graphical models, each node in the graph represents a random variable, but here the edges connecting nodes are undirected. Rather than describing a factorisation of a joint density into conditional and marginal densities, an undirected graphical model indicates the factorisation of a joint density into a product of *clique potentials* on each of the *maximal cliques* in the graph.

A *clique* is a fully connected component of the graph - i.e. a subset of nodes in the graph such that all pairs of nodes in the subset are connected by an edge. A *maximal clique* is a clique which is not a strict subset of any other clique. A *clique potential* is a non-negative function of the values of the random variables in the clique; it does not necessarily correspond to any conditional or marginal probability density.

An example undirected graphical model on four random variables,  $x_1$ ,  $x_2$ ,  $x_3$  and  $x_4$ , is shown in Figure 1.1b. Here the (maximal) cliques correspond to all the connected pairs of nodes. If  $\psi_{a,b}$  denotes the clique potential on the pair  $(a, b)$  then the graphical model implies the joint density can be factorised as

$$p_{x_1, x_2, x_3, x_4}(x_1, x_2, x_3, x_4) = \frac{1}{Z} \psi_{x_1, x_2}(x_1, x_2) \psi_{x_1, x_3}(x_1, x_3) \psi_{x_2, x_4}(x_2, x_4) \psi_{x_3, x_4}(x_3, x_4), \quad (1.47)$$

with  $Z$  a normalising constant such that the density integrates to 1 and so defines a valid probability measure.

Undirected graphical models are a natural representation for models of systems of mutually interacting components. For example they are commonly used in models of images to represent dependencies between pixel values and models of ferromagnetism to represent interactions between lattices of particles.

Unlike directed models, generating joint configurations of the random variables in an undirected graphical model from the implied joint distri-

*Ancestral sampling in a directed graphical model corresponds to first sampling values from all the root nodes from their marginal densities, then iteratively sampling from the conditional densities on each node for which all the parents nodes already have sampled values to condition on.*

bution is typically a non-trivial task, with no general equivalent to ancestral sampling. Further the joint density can typically only be evaluated up to an unknown normalising constant, with the integral needed to evaluate this constant often intractable for models involving a large number of variables or complex potentials. These properties mean that inference in distributions defined by undirected graphical models is often particularly challenging.

*D-separation:*  
 $x \perp y | C \iff$  all  
 paths in the graph  
 between  $x$  and  $y$  are  
 blocked. A path is  
 blocked if at least one  
 of the following holds:  
 1. The path includes a  
 $\rightarrow \bigcirc \rightarrow$  node or a  
 $\leftarrow \bigcirc \leftarrow$  node in  $C$ .  
 2. The path includes a  
 $\rightarrow \bigcirc \leftarrow$  node and  
 neither the node or its  
 descendants are in  $C$ .

*U-separation:*  
 $x$  and  $y$  in the model  
 and a conditioning set  
 of random variables  
 $C$ ,  $x \perp y | C \iff$  at  
 least one random  
 variable node on  
 every path between  $x$   
 and  $y$  is in  $C$ .

As suggested at the start of this section, both directed and undirected graphical models encode conditional independence properties of probabilistic models. In particular the rules of *D-separation* for directed graphical models and *U-separation* for undirected model give graph-based algorithmic descriptions of how to determine whether a pair of random variables are conditionally independent for a given conditioning set of random variables.

For example the directed graphical model in Figure 1.1a encodes the (un)conditional independence property  $x_1 \perp x_2 | \emptyset = x_1 \perp x_2$  i.e. that  $x_1$  and  $x_2$  are independent if the value of  $x_3$  is *not* conditioned on. The undirected graphical model in Figure 1.1b encodes the conditional independence properties  $x_1 \perp x_4 | x_2, x_3$  and  $x_2 \perp x_3 | x_1, x_4$ .

Although there are methods to convert a directed graphical model to an undirected one and vice versa, in general these transformations are lossy - not all of the conditional independence relationships encoded in the original graph will necessarily be maintained in the transformed graph. For example there is no undirected graphical model which will represent the exact set of conditional independence properties represented by the directed graphical model in Figure 1.1a. Likewise there is no directed graphical model which will represent the exact set of conditional independence properties represented by the undirected graphical model in Figure 1.1b. Further there are distributions with dependency structures and factorisations which cannot be uniquely represented by either directed or undirected graphical models [9].

### 1.2.2 Factor graphs

An alternative graphical model formalism which overcomes some of the limitations of directed and undirected graphical models is that of factor graphs [9, 10]. In factor graphs, in addition to nodes representing random variables, represented as in directed and undirected graphical



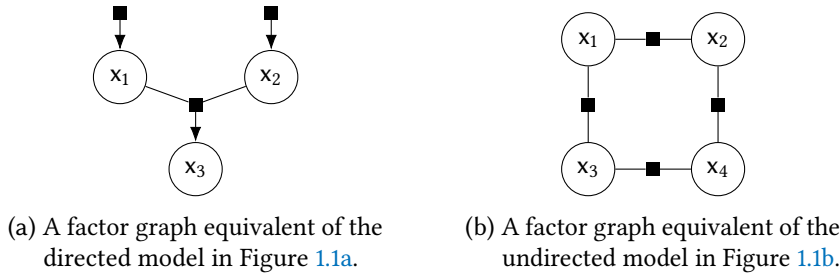


Figure 1.2: Examples of factor graphs corresponding to the directed and undirected graphical models in Figure 1.1. Square black nodes correspond to individual factors depending on the connected variables (represented by circular nodes) in the joint density.

models by circular nodes, a second class of nodes, denoted by filled squares ( $\blacksquare$ ), are introduced which represent individual factors in the joint density across the random variables represented in the model.

Factors may be either directed or undirected. Undirected factors, denoted by factor nodes in which all edges connecting to variable nodes are undirected, correspond to a factor in the joint density which depends on all of the variables with nodes connected to the factor, but without any requirement that the factor corresponds to a conditional or marginal probability density. Directed factors, denoted by factor nodes in which at least one edge from the factor node to a variable node is directed, correspond to a conditional density on the variables pointed to by directed edges given the values of the variables connected to the factor node by undirected edges (if there are no such variables then the factor instead corresponds to a marginal density).

Edges between nodes in a factor graph are always between nodes of disparate types i.e. between factor and variable nodes, but never between factor and factor or variable and variable nodes. As with directed graphical models, factor graphs with directed factors must not contain any directed cycles (i.e. a connected loop of edges in which one of every pair of edges connected to a factor on the loop is directed and all of the directed edges point in the same sense around the loop).

In the original extension of undirected factor graphs [10] to include directedness [9], it was proposed to allow multiple directed factors to connect via directed edges to the same variable node, representing multiple factors in a conditional density on that node. This generalisation introduces extra normalisation requirements and loses the interpretation

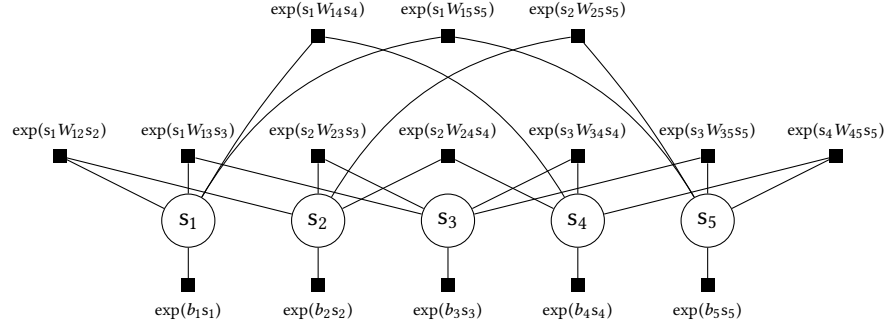


Figure 1.3: Five unit Boltzmann machine factor graph showing explicit factorisation of distribution into pairwise and single variable potentials.

of a directed factor as directly representing a conditional density, and so we will here only use directed factor graphs in which there is at most one directed edge connecting from a factor to a node.

Whether two variables are conditionally independent given a set of other variables can be checked from a factor graph by checking if all paths (i.e. connected series of edges and nodes) between the two corresponding variables nodes in the factor graph are *blocked*. A path is blocked if at least one of the following conditions is satisfied [9]

1. One of the variable nodes in the path is in the conditioning set.
2. One of the directed factor nodes in the path has two connected undirected edges in the path and there is no second directed path from the node to a variable node in the conditioning set.

Both directed and undirected graphical models can always be losslessly converted to a factor graph, i.e. such that by applying the above blocking rules after the transformation we obtain exactly the same set of conditional independency properties as present in the original graph, and thus they have a superset of the capacity to represent conditional independence properties as either of these two alternative frameworks. For example, factor graph equivalents of the directed and undirected graphical model examples in Figure 1.1 are shown in Figure 1.2.

As well as allowing representations of mixed graphs with both directed and undirected factors which cannot be represented with either directed or undirected graphical models, factor graphs are also able to include finer-grained information about the factorisation of the joint density than either of the other two model types by explicitly indicating the presence of individual factors. For instance Figure 1.3 shows the

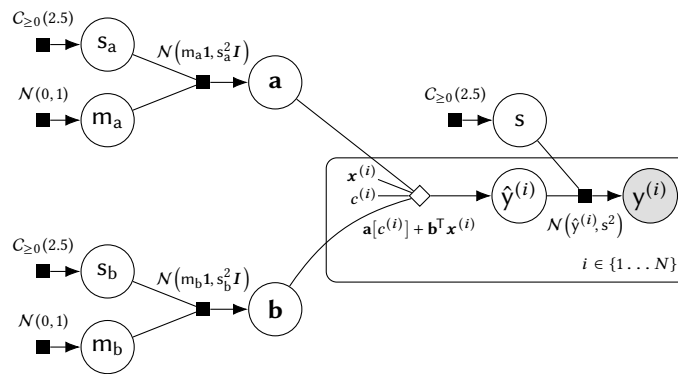


Figure 1.4: Hierarchical linear regression model factor graph showing examples of extended factor graph notation.

factor graph for a *Boltzmann machine* distribution, sometimes called a *pairwise binary Markov random field* or *Ising model*, on five binary random variables  $\{s_i\}_{i=1}^5$ . A Boltzmann machine distribution can be factored in to a product of pairwise weighted interactions  $\exp(s_i W_{ij} s_j)$  and single variable bias potentials  $\exp(b_i s_i)$ , each of which are explicitly represented by labelled factors in Figure 1.3. A corresponding undirected graphical model representation would have a single clique involving all five variables, and so would not indicate any information about the factorisation of the joint density.

In Figure 1.4 we illustrate some additional useful factor graph notation we will use in this thesis. We use a factor graph corresponding to a hierarchical linear regression model which will be discussed in more detail later in the thesis as a motivating example. The exact meaning of the model and its various factors are unimportant to the discussion of notation here so will be skipped for now.

It will often be useful to be able to explicitly represent deterministic functions applied to the random variables in a factor graph. For this purpose we introduce an additional node type denoted by an unfilled diamond ( $\diamond$ ). The semantics of this node type are very similar to standard directed factor nodes. Variables acting as inputs to the function are connected to the node by undirected edges and the variable corresponding to the function output indicated by a directed edge from the node to the relevant variable. Like standard factor nodes, the deterministic factor nodes only ever connect to variable nodes. The operations performed by the function on the inputs will usually be included as a label adjacent to the node as illustrated by the example in Figure 1.4.

The deterministic factor nodes are directly equivalent to a standard directed factor node corresponding to a degenerate Dirac delta conditional density on the output variable which concentrates all the probability mass at the output value of the function applied to inputs (conditioning variables). From this perspective it can be seen that the previously discussed rules for evaluating conditional independency properties in factor graphs can be directly extended to account for the new node type by just considering it as a (specialised) directed factor node.

Optionally constant values used in a model may be included in a factor graph as plain nodes indicated only by a label. The  $\mathbf{x}^{(i)}$  and  $c^{(i)}$  nodes in Figure 1.4 are an example of this notation.

A commonly used convention in both factor graphs and directed graphical models is *plate notation* [4], with an example of a plate shown by the rounded rectangle bounding some of the nodes in Figure 1.4. Plates are used to indicate a subgraph in the model which is replicated multiple times (with the number of replications and an index over the replications typically indicated somewhere on the plate, for example in the lower right corner in the plate in Figure 1.4). The subgraph entirely contained on the plate is assumed to be replicated the relevant number of times, with any edges crossing into the plate from variable nodes outside of the plate being repeated once for each subgraph replication. Plates are commonly used to represent a model component repeated across multiple data items.

A half-Cauchy density is defined as

$$C_{\geq 0}(x | \gamma) = \frac{2}{\pi \gamma \left(1 + \frac{x^2}{\gamma^2}\right)}$$

for  $x \geq 0$  and 0 otherwise.

Each of the factors in Figure 1.4 is labelled with a shorthand for a probability density function corresponding to the conditional or marginal density factor associated with the node. In the examples in Figure 1.4,  $\mathcal{N}(\mu, \sigma^2)$  corresponds to a Gaussian density with mean  $\mu$  and variance  $\sigma^2$ . Similarly  $\mathcal{N}(\mu, \Sigma)$  corresponds to a multivariate Gaussian with mean  $\mu$  and covariance  $\Sigma$  and  $C_{\geq 0}(\gamma)$  corresponds to a half-Cauchy density with scale parameter  $\gamma$ . The dependence of the factors on the value of the random variable the density is defined on is omitted in the factor labels for brevity.

A final additional notation used in Figure 1.4 is the use of a shaded variable node (corresponding to  $y^{(i)}$ ) to indicate a random variable corresponding to an observable quantity in the model.

### 1.2.3 Computation graphs

A final graph based tool we will make use of in this thesis is that of *computation graphs* [1]. In particular computation graphs (via associated software frameworks [16]) will be used to allow automatic differentiation of complex probabilistic models used in later chapters. Computation graphs are not typically considered in the context of probabilistic graphical models, but they share many of the same features and as we will see are closely related to directed factor graphs.

A *computation graph*, sometimes instead termed a *computational graph* or *data flow graph*, represents the computations involved in evaluating a mathematical expression. In this thesis we will distinguish between two types of nodes in a computation graph. *Variable nodes* correspond to variables which hold either inputs to the computation or intermediate results corresponding to the outputs of sub-expressions. *Operation nodes* describe how non-input variable nodes are computed as functions of other variable nodes. In other presentations of computation graphs, often the operation nodes are instead implicitly represented by directed edges between variable nodes. However analogously to the more explicit factorisation afforded by directed factor graphs compared to directed graphical models, directly representing operations as nodes allows finer grained information about the decomposition of the operations associated with a computation graph to be included.

As with directed graphical models and directed factor graphs, computation graphs cannot contain directed cycles. This does not preclude recursive and recurrent computations however as these can always be unrolled to form a directed acyclic graph. The ‘mathematical expressions’ a computation graph is constructed to evaluate can be arbitrarily complex - a computation graph corresponding to the evaluation of any numerical algorithm can always be constructed including use of arbitrary nested flow control and branching statements.

An example of a computation graph representing the calculation of the negative log density of a univariate Gaussian distribution, i.e.

$$c = \frac{1}{2} \left( \frac{x - m}{s} \right)^2 + \log s + \frac{1}{2} \log(2\pi) \quad (1.48)$$

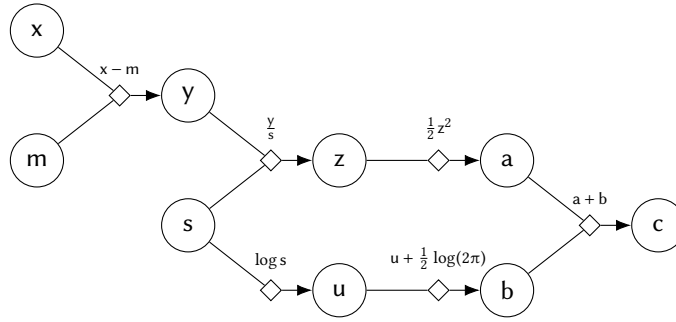


Figure 1.5: Example computation graph corresponding to calculation of the negative log density of a univariate Gaussian distribution.

is shown in Figure 1.5. The graph inputs have chosen to be the value of the random variable ( $x$ ) to evaluate the density at and the mean ( $m$ ) and the standard deviation ( $s$ ) parameters of the density.

Variable nodes in the computation graph have been represented by labelled circles and operation nodes with labelled diamonds. Undirected edges connecting from a variable node to an operation node correspond to the inputs to the operation, and directed edges from an operation node to variable nodes to the outputs of the operation.

The computation graph associated with a given expression is not uniquely defined. There will usually be multiple possible orderings in which operations can be applied to achieve the same result (up to differences due to non-exact floating point computation). Similarly what should be considered a single operation to be represented by a node in the computation graph as opposed to being split up into a sub-graph of multiple operations is a matter of choice. For example in Figure 1.5 the addition of the constant  $\frac{1}{2} \log(2\pi)$  could have been included at various other points in the graph and the operation  $\frac{1}{2} z^2$  could have been split in to separate multiplication and exponentiation operations.

The main motivation for representing expressions as computation graphs is to formalise an efficient general procedure for automatically calculating derivatives of the output of an expression with respect to its inputs termed automatic differentiation [3, 13]. The key ideas in automatic differentiation are to use the chain rule to decompose the derivatives into products and sums of the partial derivatives of the output of each individual operation in the expression with respect to its input, and to use an efficient recursive accumulation of these partial derivative sum-

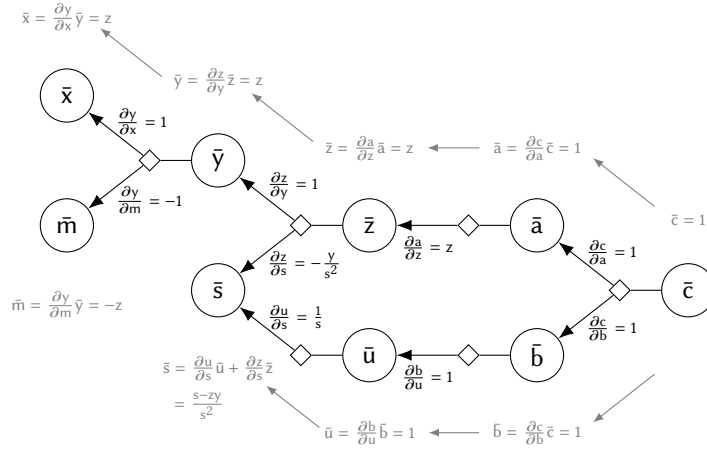


Figure 1.6: Visualisation of applying reverse-mode automatic differentiation to the computation graph in Figure 1.5 to calculate the derivatives of the negative log density of a univariate Gaussian distribution.

products corresponding to a traversal of the computation graph such that multiple derivatives can be efficiently calculated together.

Depending on how the computation graph is traversed to accumulate the derivative terms, different modes of automatic differentiation are possible. Of most direct use in most inference problems is *reverse-mode accumulation* [15], in which the derivatives of an output node with respect to all input nodes are accumulated by a reverse pass through the computation graph from the output node to inputs.

As an example the partial derivatives of the expression for univariate Gaussian log density given in (1.48) with respect to  $x$ ,  $m$  and  $s$  can be decomposed using the chain rule in terms of the intermediate variables in the computation graph shown in Figure 1.5 as

$$\frac{\partial c}{\partial x} = \frac{\partial c}{\partial a} \frac{\partial a}{\partial z} \frac{\partial z}{\partial y} \frac{\partial y}{\partial x}, \quad (1.49)$$

$$\frac{\partial c}{\partial m} = \frac{\partial c}{\partial a} \frac{\partial a}{\partial z} \frac{\partial z}{\partial y} \frac{\partial y}{\partial m}, \quad (1.50)$$

$$\frac{\partial c}{\partial s} = \frac{\partial c}{\partial a} \frac{\partial a}{\partial z} \frac{\partial z}{\partial s} + \frac{\partial c}{\partial b} \frac{\partial b}{\partial u} \frac{\partial u}{\partial s}. \quad (1.51)$$

We can immediately see that some of the chains of products of partial derivatives are repeated in the different derivative expressions - for example  $\frac{\partial c}{\partial a} \frac{\partial a}{\partial z}$  appears in the expressions for all three derivatives. Reverse-mode accumulation is effectively an automatic way of exploiting these possibilities for reusing calculations.

Figure 1.6 shows a visualisation of reverse-mode accumulation applied to the computation graph in Figure 1.5. The first step is for a *forward pass* through the graph to be performed, i.e. values are provided for each of the input variables and then each of the intermediate and output variables calculated from the incoming operation applied to their parent values. Importantly the values of all variables in the graph calculated during the forward pass must be maintained in memory.

The *reverse pass* recursively calculates the values of the partial derivatives of the relevant output node with respect to each variable node in the graph - we will term these intermediate derivatives *accumulators* denoted with barred symbols in Figure 1.6 e.g.  $\bar{a} = \frac{\partial c}{\partial a}$ . The reverse pass begins by seeding an accumulator for the output node to one (i.e.  $\bar{c} = \frac{\partial c}{\partial c} = 1$  in Figure 1.6).

Accumulators for the input variables of an operation are calculated by multiplying the accumulator for the operation output by the partial derivatives of the operation output with respect to each input variable. For non-linear operations multiplying by the operator partial derivatives will require access to the value of the input variables calculated in the forward pass. If a variable is an input to multiple operations, the derivative terms from each operation are added together in the relevant accumulator, as for example shown for  $\bar{s}$  in Figure 1.6. By recursively applying these product and sum operations, the derivatives of the output with respect to all variables in the graph can be calculated.

This reverse accumulation method allows computation of numerically exact (up to floating point error) derivatives of a single output variable in a computation graph with respect to *all input variables* with a computational cost, in terms of the number of atomic operations which need to be performed, that is a constant factor of the cost of the evaluation of the original expression represented by the computation graph in the forward pass. The constant factor is typically two to three and at most six [2]). This efficient computational cost is balanced by the requirement that the values of all intermediate variables in the computation graph evaluated in the forward pass through the graph must be stored in memory for the derivative accumulation in a reverse pass, which for large computational graphs can become a bottleneck.

To calculate the full Jacobian from a computation graph representing a function with  $M$  inputs  $\{x_i\}_{i=1}^M$  and  $N$  outputs  $\{y_i\}_{i=1}^N$ , i.e. the  $N \times M$



matrix  $mtxJ$  with entries  $J_{i,j} = \frac{\partial y_i}{\partial x_j}$ , we can do a single forward pass and  $N$  reverse passes each time accumulating the derivatives of one output variable with respect to all inputs. This leads to an overall computational cost that is  $O(N)$  times the cost of a single (forward) function evaluation to evaluate the full Jacobian. Note that as each of the reverse passes can trivially be run in parallel (in addition to any parallelisation of the operations in the forward and reverse passes themselves), this  $O(N)$  factor in the operation count need not corresponds to an equivalent increase in compute time.

An alternative to reverse-mode accumulation is *forward-mode accumulation* [18], which insteads accumulates partial derivatives with respect to a single input variable alongside a forward pass through the graph. In contrast to reverse-mode, this allows calculation of the partial derivatives of all output variables with respect to a single input variable at a computational cost that is a constant factor of the cost of the evaluation of the original expression in the forward pass. Forward-mode accumulation therefore allows evaluation of the Jacobian of a function with  $M$  inputs and  $N$  outputs at an overall computational cost that is  $O(M)$  times the cost of a single function evaluation. For functions with  $M \gg N$ , e.g. scalar valued functions of multiple inputs, reverse-mode accumulation is generally therefore significantly more efficient at computing the Jacobian. Forward-mode accumulation is however useful for functions with  $N \gg M$ , and also has the advantage over reverse-mode accumulation of avoiding the requirement to store the values of intermediate variables from the forward pass for the reverse pass(es).

The direct overlap in our notation to represent variable and operation nodes in computation graphs and that used to represent (random) variable nodes and deterministic factor nodes in our above presentation of factor graphs is intentional. Although often the operations associated with a deterministic factor node in a factor graph will be more complex than the operations usually represented by nodes in a computation graph, this is only a matter of granularity of representation - fundamentally they perform the same role. Similarly while we often consider random variables as stochastic quantities, formally they are defined as measurable functions on the sample space, and likewise the variable nodes in a computation graph represent functional transforms of their ancestor nodes.

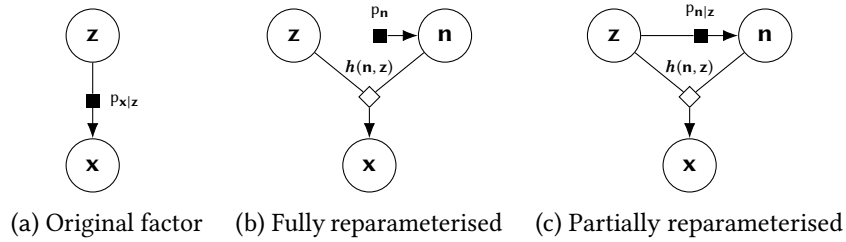


Figure 1.7: Full and partial reparameterisation transforms of a factor node representing a conditional density factor  $p_{x|z}$ . In full reparameterisation an additional random variable  $\mathbf{n}$  is introduced which (unconditionally) independent of  $\mathbf{x}$  such that the output of a deterministic transformation  $\mathbf{x} = \mathbf{h}(\mathbf{n}, \mathbf{z})$  has the same conditional density given  $\mathbf{z}$  as the original factor. In partial reparameterisation the introduced random variable  $\mathbf{n}$  is instead conditionally dependent on  $\mathbf{z}$ .

Like directed graphical models, a directed factor graph naturally specifies a process for sampling from a generative model via ancestral sampling, with values for the random variables in the graph successively calculated in a forward pass consisting of a combination of deterministic and stochastic transformations of the values of parent variables. A computation graph likewise specifies a generative process (how to compute the expression outputs given inputs) computed via a forward pass through the graph with the main differences being here that the inputs to that process are assumed to be given rather than sampled from marginal densities and the intermediate operations are all deterministic.

Although we originally motivated computation graphs as giving a graphical representation of ‘mathematical expressions’, it should be noted that these expressions need not be composed of the standard arithmetic type operations we might typically think. Any computation, involving arbitrarily complex

### 1.3 INFERENCE

Having now introduced the underlying theoretical tools we use to construct probabilistic models, we will now describe more concretely the problem of probabilistic inference.

The starting point of any inference problem is a model (or models). The model codifies our assumptions about the world. In particular we will assume the model proposes a probabilistic relationship between observ-

*You cannot do  
inference without  
making assumptions  
—David Mackay*

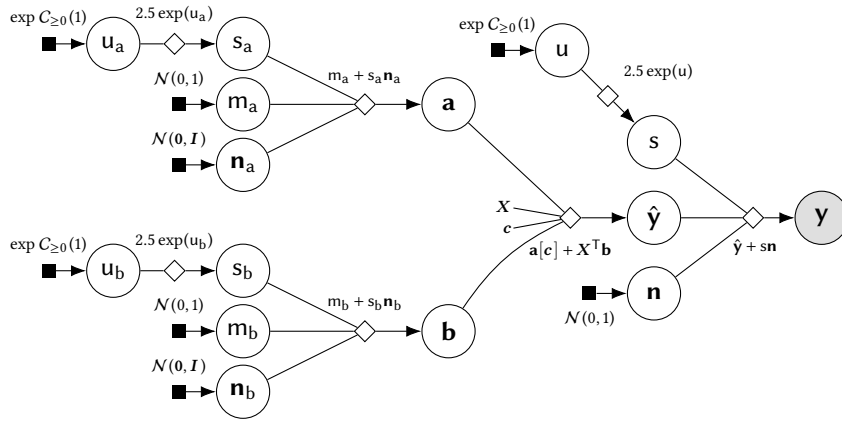


Figure 1.8: Hierarchical linear regression model stochastic computation graph.

able quantities, which we will denote with the random vector  $\mathbf{y}$ , and unobserved quantities which we wish to infer, which we will denote as the random vector  $\mathbf{z}$ .

Implicitly or explicitly the model defines the joint probability  $P_{\mathbf{y}, \mathbf{z}}$ ; this might be by defining a joint probability density  $p_{\mathbf{y}, \mathbf{z}}$  which specifies  $P_{\mathbf{y}, \mathbf{z}}$ , or by construction of a *simulator* or *generative model* which allows to generate values for  $\mathbf{y}$  and  $\mathbf{z}$  without the resulting joint probability  $P_{\mathbf{y}, \mathbf{z}}$  being directly specified.

#### 1.3.1 Posterior expectations

#### 1.3.2 Model evidence



# 2 | APPROXIMATE INFERENCE

## 2.1 DETERMINISTIC APPROACHES

2.1.1 Laplace's method

2.1.2 Variational inference

2.1.3 Expectation propagation

## 2.2 STOCHASTIC APPROACHES

2.2.1 Monte Carlo method

2.2.2 Rejection sampling

2.2.3 Importance sampling

2.2.4 Markov chain Monte Carlo



# 3 | MARKOV CHAIN MONTE CARLO

## 3.1 METROPOLIS-HASTINGS

## 3.2 GIBBS SAMPLING

## 3.3 SLICE SAMPLING

## 3.4 HAMILTONIAN MONTE CARLO





# 4 | THERMODYNAMIC METHODS

4.1 SIMULATED TEMPERING

4.2 PARALLEL TEMPERING

4.3 TEMPERED TRANSITIONS

4.4 ANNEALED IMPORTANCE SAMPLING

4.5 PATH SAMPLING

4.6 ADIABATIC MONTE CARLO



## BIBLIOGRAPHY

- [1] Friedrich L Bauer. ‘Computational graphs and rounding error’. In: *SIAM Journal on Numerical Analysis* 11.1 (1974), pp. 87–96.
- [2] Atilim Gunes Baydin, Barak A Pearlmutter, Alexey Andreyevich Radul and Jeffrey Mark Siskind. ‘Automatic differentiation in machine learning: a survey’. In: *arXiv preprint arXiv:1502.05767* (2015).
- [3] L. M. Beda, L. N. Korolev, N. V. Sukkikh and T. S. Frolova. *Programs for automatic differentiation for the machine BESM*. Technical Report. (In Russian). Moscow, USSR: Institute for Precise Mechanics and Computation Techniques, Academy of Science, 1959.
- [4] Wray L Buntine. ‘Operations for learning with graphical models’. In: *Journal of artificial intelligence research* (1994).
- [5] Richard T Cox. ‘Probability, frequency and reasonable expectation’. In: *American Journal of Physics* 14.1 (1946), pp. 1–13. URL: <http://dx.doi.org/10.1119/1.1990764>.
- [6] Richard T Cox. ‘The algebra of probable inference’. In: *American Journal of Physics* 31.1 (1963), pp. 66–67. URL: <http://dx.doi.org/10.1119/1.1969248>.
- [7] Herbert Federer. *Geometric measure theory*. Springer, 1969.
- [8] Bruno de Finetti. ‘Foresight: its logical laws, its subjective sources’. In: *Studies in Subjective Probability*. Ed. by H. E. Kyburg. English translation of original 1937 French article *La Prévision: ses lois logiques, ses sources subjectives*. Springer, 1992, pp. 134–174.
- [9] Brendan J Frey. ‘Extending factor graphs so as to unify directed and undirected graphical models’. In: *Proceedings of the Nineteenth conference on Uncertainty in Artificial Intelligence*. Morgan Kaufmann Publishers Inc. 2002, pp. 257–264. URL: <https://arxiv.org/abs/1212.2486>.
- [10] Brendan J Frey, Frank R Kschischang, Hans-Andrea Loeliger and Niclas Wiberg. ‘Factor graphs and algorithms’. In: *Proceedings of the 35th Annual Allerton Conference on Communication Control and Computing*. 1997.

- [11] Ross Kindermann and Laurie Snell. *Markov random fields and their applications*. American Mathematical Society, 1980.
- [12] Andreï Nikolaevich Kolmogorov. *Foundations of the Theory of Probability*. Ed. by Nathan Morrison. 2nd English Edition. English translation of original 1933 German monograph, *Grundbegriffe der Wahrscheinlichkeitrechnung*. Chelsea Publishing Company, 1956. URL: <https://pdfs.semanticscholar.org/c3e1/51f71168a5f348bdebfde11752ca603fa6d0.pdf>.
- [13] John F Nolan. ‘Analytical differentiation on a digital computer’. PhD thesis. Massachusetts Institute of Technology, 1953.
- [14] Judea Pearl. *Probabilistic reasoning in intelligent systems: Networks of plausible inference*. Morgan Kaufmann, 1988.
- [15] Bert Speelpenning. ‘Compiling Fast Partial Derivatives of Functions Given by Algorithms’. PhD thesis. University of Illinois at Urbana-Champaign, 1980.
- [16] Theano Development Team et al. ‘Theano: A Python framework for fast computation of mathematical expressions’. In: *arXiv e-prints* abs/1605.02688 (May 2016). URL: <http://arxiv.org/abs/1605.02688>.
- [17] Alexander Terenin and David Draper. ‘Cox’s Theorem and the Jaynesian Interpretation of Probability’. arXiv preprint. 2015. URL: <https://arxiv.org/abs/1507.06597v2>.
- [18] Robert Edwin Wengert. ‘A simple automatic derivative evaluation program’. In: *Communications of the ACM* 7.8 (1964), pp. 463–464.