

Discrete Probabilistic Programming Languages¹

Steven Holtzen

s.holtzen@northeastern.edu

September 29, 2023

¹ CS7470 Fall 2023: Foundations of Probabilistic Programming.

- Goal for today: compile a simple discrete PPL to BDD

1 Compositional compilation of BDDs

- So far we have compiled by inducting on variables
- **Problem:** this is not *compositional*! A compilation process is compositional if it works by compiling big programs out of smaller sub-programs. I.e., it would have a rule that looks something like:

$$\frac{\alpha \rightsquigarrow \alpha' \quad \beta \rightsquigarrow \beta'}{\alpha \wedge \beta \rightsquigarrow \alpha' \times \beta'}$$

- Compositional compilation is great: gives us modular reasoning about performance, ... (other reasons?)
- **Goal:** Design a compositional process for compiling PROP_S to BDD
- How can we build big BDDs out of smaller ones? Define a way to compose together BDDs, again using a type-directed step-relation $\Gamma \vdash \alpha_1 \wedge \alpha_2 \Downarrow \beta$, shown in Figure 1.

Theorem 1 (Correctness). *If $\Gamma \vdash \alpha_1$, $\Gamma \vdash \alpha_2$, and $\Gamma \vdash \alpha_1 \wedge \alpha_2 \Downarrow \beta$, then $\llbracket \alpha_1 \rrbracket \cap \llbracket \alpha_2 \rrbracket = \llbracket \beta \rrbracket$.*

Proof. By simultaneous structural induction on syntax of BDDs (note that we need to perform simultaneous induction since there are two BDDs at play here). The rules in Figure 1 are exhaustive (i.e., every pair of syntactic BDDs matches exactly one of these structural rules), so we can proceed by case analysis on each of the compilation rules. The base cases are quite simple and we elide them here. The interesting cases are the inductive cases.

We will show the case for (SAMEVARNE). Assume that $\Gamma \vdash \alpha_1 \wedge \alpha_3 \Downarrow \alpha_{13}$ and $\Gamma \vdash \alpha_2 \wedge \alpha_4 \Downarrow \alpha_{24}$. As usual, our inductive hypothesis assumes that the theorem holds for compiled subterms:

- $\llbracket \Gamma \vdash \alpha_1 \rrbracket \cap \llbracket \Gamma \vdash \alpha_3 \rrbracket = \llbracket \Gamma \vdash \alpha_{13} \rrbracket$
- $\llbracket \Gamma \vdash \alpha_2 \rrbracket \cap \llbracket \Gamma \vdash \alpha_4 \rrbracket = \llbracket \Gamma \vdash \alpha_{24} \rrbracket$

We want to show that:

$$\left\| x :: \Gamma \vdash \begin{array}{c} \alpha_1 \quad \alpha_2 \\ \nearrow \quad \nwarrow \\ x \end{array} \right\| \cap \left\| x :: \Gamma \vdash \begin{array}{c} \alpha_3 \quad \alpha_4 \\ \nearrow \quad \nwarrow \\ x \end{array} \right\| = \left\| x :: \Gamma \vdash \begin{array}{c} \alpha_{13} \quad \alpha_{24} \\ \nearrow \quad \nwarrow \\ x \end{array} \right\|$$

We reason forward:

$$\begin{aligned} \left\| x :: \Gamma \vdash \begin{array}{c} \alpha_{13} \quad \alpha_{24} \\ \nearrow \quad \nwarrow \\ x \end{array} \right\| &= [x \mapsto \text{true}] \otimes [\Gamma \vdash \alpha_{13}] \cup [x \mapsto \text{false}] \otimes [\Gamma \vdash \alpha_{24}] \\ &= [x \mapsto \text{true}] \otimes ([\Gamma \vdash \alpha_1] \cap [\Gamma \vdash \alpha_3]) \cup [x \mapsto \text{false}] \otimes ([\Gamma \vdash \alpha_2] \cap [\Gamma \vdash \alpha_4]) && \text{By I.H.} \\ &= ([x \mapsto \text{true}] [\Gamma \vdash \alpha_1] \cap [x \mapsto \text{true}] [\Gamma \vdash \alpha_3]) \cup ([x \mapsto \text{false}] [\Gamma \vdash \alpha_2] \cap [x \mapsto \text{false}] [\Gamma \vdash \alpha_4]) && (*) \\ &= ([x \mapsto \text{true}] [\Gamma \vdash \alpha_1] \cup [x \mapsto \text{false}] [\Gamma \vdash \alpha_2]) \cap ([x \mapsto \text{true}] [\Gamma \vdash \alpha_1] \cup [x \mapsto \text{false}] [\Gamma \vdash \alpha_4]) && (†) \\ &= \left\| x :: \Gamma \vdash \begin{array}{c} \alpha_1 \quad \alpha_2 \\ \nearrow \quad \nwarrow \\ x \end{array} \right\| \cap \left\| x :: \Gamma \vdash \begin{array}{c} \alpha_3 \quad \alpha_4 \\ \nearrow \quad \nwarrow \\ x \end{array} \right\| \end{aligned}$$

where $(*)$ follows from a simple lemma that \otimes distributes over intersection and $(†)$ follows from distributivity properties of union and intersection, in particular the fact that for any sets A, B, C, D it is the case that $(A \cup B) \cap (C \cup D) = (A \cap C) \cup (B \cap D)$.²

□

² This set theory property has a nice “proof by Venn-diagram”; draw the Venn-diagram of these sets to see this fact clearly.

- We would also like to ensure that our compilation rules produce reduced and ordered BDDs. We can formalize this with a type-preservation theorem:

Theorem 2 (Type preservation). *If $\Gamma \vdash \alpha_1$, $\Gamma \vdash \alpha_2$, and $\Gamma \vdash \alpha_1 \wedge \alpha_2 \Downarrow \alpha$, then $\Gamma \vdash \alpha$.*

- These rules are called *bottom-up BDD compilation* [Darwiche and Marquis, 2002, Oztok and Darwiche, 2015].
- There are several other compositional BDD operations we won’t have time to explain in lecture, but you can use in your lab, like disjunction, negation, substitution, and existential quantification.
- Why might one prefer one mode of compilation over the other? Do they have different runtime cost? *Yes!*

- **Exercise:** Give example families of formulae where top-down compilation is faster than bottom-up and vice-versa.

2 *Disc: A simple discrete PPL*

- Syntax:

```

1  e ::=
2  | x ← e; e
3  | observe e; e
4  | flip q                                // q is a rational value
5  | if e then e else e
6  | return e
7  | true | false
8  | e ∧ e | e ∨ e | ¬ e |
9  | ( e )
10 p ::= e

```

- Disc looks very similar to a standard functional programming language, but has two some interesting new keywords: `flip`, `observe`, and `return`
- `flip θ` allocates a new random quantity that is true with probability θ and false with probability $1 - \theta$
- `return e` turns a non-probabilistic quantity into a probabilistic one, i.e. `return true` is a *probabilistic quantity* that is true with probability 1 and false with probability 0
- Example program and its interpretation:

```

1 x ← flip 0.5;
2 y ← flip 0.5;
3 return x ∧ y

```

This program outputs the probability distribution $[\text{true} \mapsto 0.25, \text{false} \mapsto 0.75]$.

- `observe` is a powerful keyword that lets us *condition* the program. For instance, suppose I want to model the following scenario: “flip two coins and observe that at least one of them is heads. What is the probability that the first coin was heads?”

We can encode this scenario as a Disc program:

```

1 x ← flip 0.5;
2 y ← flip 0.5;
3 observe x ∨ y;
4 return x

```

This program outputs the probability distribution:

$[\text{true} \mapsto (0.25 + 0.25)/0.75, \text{false} \mapsto 0.25/0.75]$

- **Type system:** terms can either be pure Booleans of type \mathbb{B} or distributions on Booleans of type $\text{Dist}(\mathbb{B})$. So, we have the following type definition:

$$\tau ::= \mathbb{B} \mid \text{Dist}(\mathbb{B}). \quad (1)$$

- We define a typing judgment $\Gamma \vdash e : \tau$ that associates each term with a type. The typing context Γ is a map from identifiers to types.

$$\begin{array}{c} \Gamma \vdash \text{true} : \mathbb{B} \quad \Gamma \vdash \text{false} : \mathbb{B} \quad \Gamma \vdash \text{flip } \theta : \text{Dist}(\mathbb{B}) \\[10pt] \frac{\Gamma \vdash e : \mathbb{B}}{\Gamma \vdash \text{return } e : \text{Dist}(\mathbb{B})} \quad \frac{\Gamma \vdash e_1 : \mathbb{B} \quad \Gamma \vdash e_2 : \tau}{\Gamma \vdash \text{observe } e_1; e_2 : \tau} \\[10pt] \frac{\Gamma \vdash e_1 : \text{Dist}(\mathbb{B}) \quad \Gamma \cup [x \mapsto \mathbb{B}] \vdash e_2 : \tau}{\Gamma \vdash x \leftarrow e_1; e_2 : \tau} \\[10pt] \frac{\Gamma \vdash e_1 : \mathbb{B} \quad \Gamma \vdash e_2 : \tau \quad \Gamma \vdash e_3 : \tau}{\Gamma \vdash \text{if } e_1 \text{ then } e_2 \text{ else } e_3 : \tau} \quad \frac{\Gamma \vdash e_1 : \mathbb{B} \quad \Gamma \vdash e_2 : \mathbb{B}}{\Gamma \vdash e_1 \wedge e_2 : \mathbb{B}} \end{array}$$

Denotational semantics of DISC

- Associates each term with an *unnormalized probability distribution* (i.e., the total probability mass may be less than 1).
- Has the type $\llbracket e \rrbracket : \text{Bool} \rightarrow [0, 1]$, and has the following inductive definition:

$$\begin{aligned} \llbracket \text{flip } \theta \rrbracket(v) &= \begin{cases} \theta & \text{if } v = T \\ 1 - \theta & \text{if } v = F \end{cases} \\[10pt] \llbracket \text{return } e \rrbracket(v) &= \begin{cases} 1 & \text{if } v = \llbracket e \rrbracket \\ 0 & \text{otherwise} \end{cases} \\[10pt] \llbracket x \leftarrow e_1; e_2 \rrbracket(v) &= \sum_{v'} \llbracket e_1 \rrbracket(v') \times \llbracket e_2[x \mapsto v'] \rrbracket(v) \\[10pt] \llbracket \text{observe } e_1; e_2 \rrbracket(v) &= \sum_{\{v' \mid \llbracket e_1(v') \rrbracket = T\}} \llbracket e_2 \rrbracket(v) \end{aligned}$$

- The semantics for non-probabilistic terms is standard. The semantic evaluation has type $\llbracket e \rrbracket : \text{Bool}$ for closed terms.

- These semantics give an unnormalized distribution. The main semantic object of interest is the normalized distribution, which is given by the **normalized semantics**:

$$\llbracket e \rrbracket_D(T) = \frac{\llbracket e \rrbracket(T)}{\llbracket e \rrbracket(T) + \llbracket e \rrbracket(F)},$$

defined analogously for the false case.

Compiling DISC to PROP

- **Goal:** give a semantics-preserving compilation \rightsquigarrow that compiles DISC to PROP
- In order to handle observations, we will compile DISC programs into *two* PROP programs: one that computes the unnormalized probability of returning true, and one that computes the probability of evidence (i.e. normalizing constant)
- Inductive description has the shape $e \rightsquigarrow (p_1, p_2)$. We want to define this relation to satisfy the following **adequacy condition**:

$$\llbracket e \rrbracket_D(\text{true}) = \frac{\llbracket p_1 \rrbracket}{\llbracket p_2 \rrbracket}. \quad (2)$$

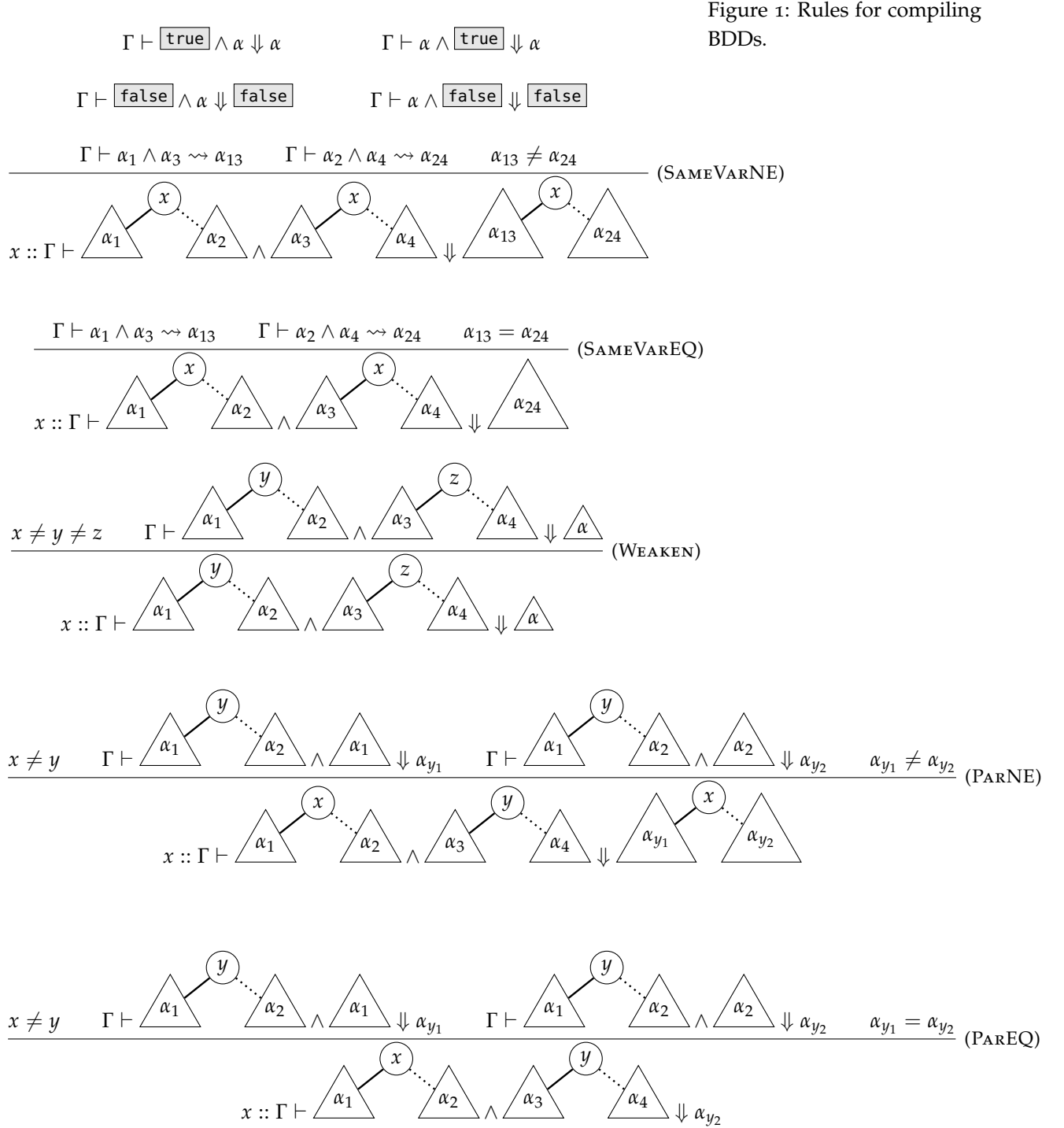
We will shorten this description to $e \rightsquigarrow (\varphi, \varphi_A, w)$ and assume that the two formulae share a common w . The adequacy condition then becomes:

$$\llbracket e \rrbracket_D(\text{true}) = \frac{\llbracket (\varphi, w) \rrbracket}{\llbracket (\varphi_A, w) \rrbracket}. \quad (3)$$

- Compilation relation:

$$\begin{array}{c} \text{true} \rightsquigarrow (\text{true}, \text{true}, \emptyset) \qquad \text{false} \rightsquigarrow (\text{false}, \text{true}, \emptyset) \\[10pt] x \rightsquigarrow (x, \text{true}, \emptyset) \qquad \frac{\text{fresh } x}{\text{flip } \theta \rightsquigarrow (x, \text{true}, [x \mapsto \theta, \bar{x} \mapsto 1 - \theta])} \\[10pt] \frac{e_1 \rightsquigarrow (\varphi, \varphi_A, w) \quad e_2 \rightsquigarrow (\varphi', \varphi'_A, w')}{x \leftarrow e_1; e_2 \rightsquigarrow (\varphi'[\varphi/x], \varphi'_A[\varphi/x] \wedge \varphi_A, w_1 \cup w_2)} \\[10pt] \frac{e_1 \rightsquigarrow (\varphi, \varphi_A, w) \quad e_2 \rightsquigarrow (\varphi', \varphi'_A, w')}{\text{observe } e_1; e_2 \rightsquigarrow (\varphi', \varphi'_A \wedge \varphi_A, w_1 \cup w_2)} \end{array}$$

Theorem 3 (Adequacy). *For well-typed term e , assume $e \rightsquigarrow (\varphi, \varphi_A, w)$. Then, $\llbracket e \rrbracket_D(\text{true}) = \llbracket (\varphi, w) \rrbracket / \llbracket (\varphi_A, w) \rrbracket$.*



References

- Adnan Darwiche and Pierre Marquis. A knowledge compilation map. *Journal of Artificial Intelligence Research*, 17:229–264, 2002.
- Umut Oztok and Adnan Darwiche. A top-down compiler for sentential decision diagrams. In *Twenty-Fourth International Joint Conference on Artificial Intelligence*, 2015.