

# Discrete Probabilistic Programming Languages II<sup>1</sup>

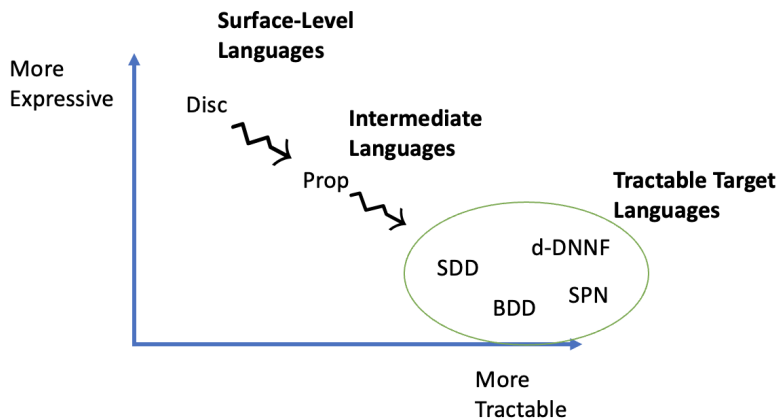
Steven Holtzen

s.holtzen@northeastern.edu

October 3, 2023

<sup>1</sup> CS7470 Fall 2023: Foundations of Probabilistic Programming.

- We have been filling in the PPL compilation picture which organizes different kinds of probabilistic programming languages:



- We now understand every component. Recall:
  - A language consists of syntax and semantics. We customarily give two kinds of semantics: a denotational semantics that describes the mathematical object associated with the program, and an operational semantics that tells us precisely how to compute that object. We typically prove something called “adequacy” that relates the denotational and operational semantics.
  - A compilation  $L_1 \rightsquigarrow L_2$  is a procedure for translating programs written in language  $L_1$  into  $L_2$  that preserves semantics.
  - Language  $L_1$  is more expressive than  $L_2$  if we can efficiently compile any program from  $L_2$  into  $L_1$  but not vice versa;
  - We can compare tractability of languages based on their worst-case runtime complexity.
- Some bibliographic notes:
  - The process of compiling a Boolean formula into a tractable target is often called *knowledge compilation* in the literature [Darwiche and Marquis, 2002]. There are many tools for compiling formulae. Some knowledge compilation tools:
    - \* c2d compiles from CNF to d-DNNF: <http://reasoning.cs.ucla.edu/c2d/>

- \* minic2d compiles from CNF to SDD: <http://reasoning.cs.ucla.edu/minic2d/>
- \* rsdd compiles from Prop to SDD and BDD: <https://github.com/neuppl/rsdd>
- \* sharpSAT compiles from CNF to d-DNNF [Muise et al., 2012]
- Disc is based on Dice [Holtzen et al., 2020]. There are other PPLs that work by knowledge compilation, such as ProbLog [Fierens et al., 2015, De Raedt et al., 2007], SPPL [Saad et al., 2021]. Historically, inference via compilation has been very effectively used in inference for graphical models [Chavira and Darwiche, 2008, Sang et al., 2005]
- What are some research questions in here?
  - Optimizing compilers for each stage.
  - New tractable target languages.
  - Compilation with performance guarantees.
  - Tools for making building these compilers.

## 1 DISC: A simple discrete PPL

- Syntax:

```

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

```

- Disc looks very similar to a standard functional programming language, but has two some interesting new keywords: `flip`, `observe`, and `return`
- `flip  $\theta$`  allocates a new random quantity that is true with probability  $\theta$  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  $[\text{true} \mapsto 0.25, \text{false} \mapsto 0.75]$ .

- **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  $\Gamma$  is a map from identifiers to types.

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

- This is often called a *monadic* type system [Moggi, 1991]; it was introduced as a way of handling effects in purely functional languages. We won't get too into monads yet. If you're curious, see Ramsey and Pfeffer [2002].

### Denotational semantics of DISC

- We will define our denotation on *closed terms*, i.e. terms with no free variables in them.
- There are two kinds of terms: *probabilistic* and *pure* terms. Probabilistic terms have type  $\text{Dist}(\mathbb{B})$ , and pure terms have type  $\mathbb{B}$ .
- The denotational semantics of pure terms is standard: they map closed terms to Boolean values. Examples of closed terms:

$$\begin{aligned} \llbracket \text{true} \rrbracket &= \text{true} \\ \llbracket \text{false} \rrbracket &= \text{false} \\ \llbracket e_1 \wedge e_2 \rrbracket &= \llbracket e_1 \rrbracket \wedge \llbracket e_2 \rrbracket \end{aligned}$$

- The denotational semantics of probabilistic terms maps closed terms to probability distributions,  $\llbracket e \rrbracket : \text{Bool} \rightarrow [0, 1]$ , and has the following inductive definition:

$$\begin{aligned}
\llbracket \text{flip } \theta \rrbracket (v) &= \begin{cases} \theta & \text{if } v = T \\ 1 - \theta & \text{if } v = F \end{cases} \\
\llbracket \text{return } e \rrbracket (v) &= \begin{cases} 1 & \text{if } v = \llbracket e \rrbracket \\ 0 & \text{otherwise} \end{cases} \\
\llbracket x \leftarrow e_1; e_2 \rrbracket (v) &= \sum_{v'} \llbracket e_1 \rrbracket (v') \times \llbracket e_2[v'/x] \rrbracket (v) \\
\llbracket \text{if true then } e_1 \text{ else } e_2 \rrbracket &= \llbracket e_1 \rrbracket \\
\llbracket \text{if false then } e_1 \text{ else } e_2 \rrbracket &= \llbracket e_2 \rrbracket
\end{aligned}$$

- Recall our notation for substitution  $e[v/x]$ , which substitutes the variable  $x$  for value  $v$ . Here we define substitution to be *capture-avoiding*. This concept is best illustrated by example. Considering the following program:

```

1 x ← true;
2 x ← false;
3 return true

```

What should this program return? Intuitively, it should return the distribution  $[\text{false} \mapsto 1, \text{true} \mapsto 0]$ .

- A substitution  $e[v/x]$  is *capture-avoiding* if it only substitutes *free occurrences* of  $x$  with  $v$  in  $e$ . A variable is **free** if it does not occur inside of a binding. We can accomplish this with the following inductive description of substitution (here we show only a subset of the rules):

$$\begin{aligned}
\text{true}[x/v] &= \text{true} & \text{false}[x/v] &= \text{false} & x[x/v] &= v \\
\frac{x \neq y}{y[x/v] = y} & & (x \leftarrow e_1; e_2)[v/x] &= x \leftarrow e_1[v/x]; e_2 \\
\frac{x \neq y}{(x \leftarrow e_1; e_2)[v/x] = y \leftarrow e_1[v/x]; e_2[v/x]}
\end{aligned}$$

- In your project, we have implemented capture-avoiding substitution for terms in Disc for you.
- A simple example of the let-rule:

$$\llbracket x \leftarrow \text{flip } \theta; \text{return } x \rrbracket (\text{true}) = \sum_{v \in \mathbb{B}} \llbracket \text{flip } \theta \rrbracket (v) \times \llbracket \text{return } v \rrbracket (\text{true}) = \theta$$

## 2 Compiling DISC to PROP

- **Goal:** give a semantics-preserving compilation  $\rightsquigarrow$  that compiles DISC to PROP.
- We want to define this relation to satisfy the following **adequacy condition**, assuming the program  $e \rightsquigarrow p$ :

$$\llbracket e \rrbracket (\text{true}) = \llbracket p \rrbracket. \quad (2)$$

- Compilation relation:

$$\begin{aligned} \text{true} &\rightsquigarrow (\text{true}, \emptyset) & \text{false} &\rightsquigarrow (\text{false}, \emptyset) \\ x &\rightsquigarrow (x, [x \mapsto 1/2, \bar{x} \mapsto 1/2]) & \frac{\text{fresh } f}{\text{flip } \theta &\rightsquigarrow (f, [f \mapsto \theta, \bar{f} \mapsto 1 - \theta])} \\ \frac{e_1 &\rightsquigarrow (\varphi, w) \quad e_2 &\rightsquigarrow (\varphi', w')}{x \leftarrow e_1; e_2 &\rightsquigarrow (\varphi'[\varphi/x], w_1 \cup w_2)} \\ \frac{e_1 &\rightsquigarrow (\varphi_1, w_1) \quad e_2 &\rightsquigarrow (\varphi_2, w_2)}{e_1 \wedge e_2 &\rightsquigarrow (\varphi_1 \wedge \varphi_2, w_1 \cup w_2)} \\ \frac{e_1 &\rightsquigarrow (\varphi_1, w_1) \quad e_2 &\rightsquigarrow (\varphi_2, w_2) \quad e_3 &\rightsquigarrow (\varphi_3, w_3)}{\text{if } e_1 \text{ then } e_2 \text{ else } e_3 &\rightsquigarrow ((\varphi_1 \wedge \varphi_2) \vee (\neg \varphi_1 \wedge \varphi_3), w_1 \cup w_2 \cup w_3)} \end{aligned}$$

- One thing to notice above: we have a new substitution rule  $\varphi'[\varphi/x]$  that substitutes  $x$  with an *arbitrary propositional formula*  $\varphi$  in  $\varphi'$ . For example:

$$(x \vee y)[z \wedge w/x] = (x \wedge y) \vee y.$$

You can perform this operation directly on BDDs; it is called *BDD composition*.

- Now we are almost ready to state the semantics-preservation property of this compilation. We need one extra piece: *environment substitution*.
- Suppose  $\Gamma \vdash e$ . A substitution  $\gamma \in \llbracket \Gamma \rrbracket$  gives an assignment to every free variable in  $e$ . We use the notation  $e[\gamma]$  to denote performing the substituting of all variables in  $\gamma$  in  $e$ , i.e.:

$$e[\{[v_1/x_1, v_2/x_2, \dots, v_n/x_n]\}] = e[v_1/x_1][v_2/x_2] \dots [v_n/x_n].$$

Additionally, note that we can apply this same substitution notation to propositional formulae, i.e.  $\varphi[\gamma]$  is also well-defined according to the above description.

**Theorem 1.** Assume  $\Gamma \vdash e$  and  $e \rightsquigarrow p$ . Then, for any substitution  $\gamma \in \llbracket \Gamma \rrbracket$ , is is the case that  $\llbracket e[\gamma] \rrbracket(\text{true}) = \llbracket p[\gamma] \rrbracket$ .

*Proof.* The proof is on structural induction on terms in  $e$ .

The base case for `flip` is simple but a good place to start. Assume  $\text{flip } \theta \rightsquigarrow (f, [f \mapsto \theta, \bar{f} \mapsto 1 - \theta])$ . Then,

$$\llbracket \text{flip } \theta \rrbracket(\text{true}) = \theta = \llbracket (f, [f \mapsto \theta, \bar{f} \mapsto 1 - \theta]) \rrbracket.$$

The critical case is the inductive step for binding. Inductive hypothesis: assume that  $e_1 \rightsquigarrow (\varphi_1, w_1)$  and  $e_2 \rightsquigarrow (\varphi_2, w_2)$ , and that for any well-typed substitution  $\gamma$  it is the case that:

- $\llbracket e_1[\gamma] \rrbracket(\text{true}) = \llbracket \varphi_1[\gamma] \rrbracket$
- $\llbracket e_2[\gamma] \rrbracket(\text{true}) = \llbracket \varphi_2[\gamma] \rrbracket$

We need to show that  $\llbracket x \leftarrow e_1; e_2 \rrbracket(\text{true}) = \llbracket \varphi_2[\varphi_1/x], (w_1 \cup w_2) \rrbracket$ .

**TODO;** if you want, you can take a peek at Holtzen et al. [2020]

□

- Question: Is this an efficient compilation? *Yes.*

### 3 Observation

- `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:

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

- 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.

- 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  $e \rightsquigarrow (p_1, p_2)$ .
- Compilation relation:

$$\begin{array}{c}
\text{true} \rightsquigarrow (\text{true}, \text{true}, \emptyset) \qquad \text{false} \rightsquigarrow (\text{false}, \text{true}, \emptyset) \\
\\
x \rightsquigarrow (x, \text{true}, \emptyset) \qquad \frac{\text{fresh } x}{\text{flip } \theta \rightsquigarrow (x, \text{true}, [x \mapsto \theta, \bar{x} \mapsto 1 - \theta])} \\
\\
\frac{e_1 \rightsquigarrow (\varphi, \varphi_A, w) \quad e_2 \rightsquigarrow (\varphi', \varphi'_A, w')}{x \leftarrow e_1; e_2 \rightsquigarrow (\varphi'[\varphi/x], \varphi'_A[\varphi/x] \wedge \varphi_A, w_1 \cup w_2)} \\
\\
\frac{e_1 \rightsquigarrow (\varphi, \varphi_A, w) \quad e_2 \rightsquigarrow (\varphi', \varphi'_A, w')}{\text{observe } e_1; e_2 \rightsquigarrow (\varphi', \varphi'_A \wedge \varphi_A, w_1 \cup w_2)} \\
\\
\llbracket \text{observe } e_1; e_2 \rrbracket(v) = \sum_{\{v' \mid \llbracket e_1(v') \rrbracket = T\}} \llbracket e_2 \rrbracket(v)
\end{array}$$

**Theorem 2** (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

- Mark Chavira and Adnan Darwiche. On probabilistic inference by weighted model counting. *Artificial Intelligence*, 172(6-7):772–799, 2008.
- Adnan Darwiche and Pierre Marquis. A knowledge compilation map. *Journal of Artificial Intelligence Research*, 17:229–264, 2002.
- Luc De Raedt, Angelika Kimmig, Hannu Toivonen, and M Veloso. Problog: A probabilistic prolog and its application in link discovery. In *IJCAI 2007, Proceedings of the 20th international joint conference on artificial intelligence*, pages 2462–2467. IJCAI-INT JOINT CONF ARTIF INTELL, 2007.
- Daan Fierens, Guy Van den Broeck, Joris Renkens, Dimitar Shterionov, Bernd Gutmann, Ingo Thon, Gerda Janssens, and Luc De Raedt. Inference and learning in probabilistic logic programs using weighted boolean formulas. *Theory and Practice of Logic Programming*, 15(3):358–401, 2015.

- Steven Holtzen, Guy Van den Broeck, and Todd Millstein. Scaling exact inference for discrete probabilistic programs. *Proceedings of the ACM on Programming Languages*, 4(OOPSLA):1–31, 2020.
- Eugenio Moggi. Notions of computation and monads. *Information and computation*, 93(1):55–92, 1991.
- Christian Muise, Sheila A McIlraith, J Christopher Beck, and Eric I Hsu. D sharp: fast d-dnnf compilation with sharpsat. In *Advances in Artificial Intelligence: 25th Canadian Conference on Artificial Intelligence, Canadian AI 2012, Toronto, ON, Canada, May 28-30, 2012. Proceedings 25*, pages 356–361. Springer, 2012.
- Norman Ramsey and Avi Pfeffer. Stochastic lambda calculus and monads of probability distributions. In *Proceedings of the 29th ACM SIGPLAN-SIGACT symposium on Principles of programming languages*, pages 154–165, 2002.
- Feras A Saad, Martin C Rinard, and Vikash K Mansinghka. Sppl: probabilistic programming with fast exact symbolic inference. In *Proceedings of the 42nd ACM SIGPLAN International Conference on Programming Language Design and Implementation*, pages 804–819, 2021.
- Tian Sang, Paul Beame, and Henry A Kautz. Performing bayesian inference by weighted model counting. In *AAAI*, volume 5, pages 475–481, 2005.