

Probabilistic compositional semantics, purely

Julian Grove and Jean-Philippe Bernardy

Centre for Linguistic Theory and Studies in Probability
Department of Philosophy, Linguistics and Theory of Science
University of Gothenburg
`firstname.lastname@gu.se`

Abstract. We provide a general framework for the integration of formal semantics with probabilistic reasoning. This framework is conservative, in the sense that it relies only on typed λ -calculus and is thus compatible with logical systems already in use. The framework is also presented modularly, in that it regards probabilistic effects (i.e., sampling and marginalization) as *side effects*, using continuations. We show how our framework may be used to build probabilistic programs compositionally within higher-order logic and then illustrate its use on two applications: semantic learning and pragmatic inference within the Rational Speech Act framework.

1 Introduction

Formal semantics in the tradition of Montague characterizes linguistic meaning in terms of either a logic or a model, constructed set-theoretically. By exploiting an antecedently well understood formalism, a logical characterization of meaning allows one to reason about it in terms of notions like entailment. Indeed, while the formal description such a characterization provides is necessarily abstract, it can be assembled compositionally, in terms of rules that combine the meanings of syntactic constituents. It is this feature of formal semantics that makes it such an attractive approach to meaning, and one which has persisted throughout the development of the field.

There has been much effort in the last decade to connect formal semantics to mathematically explicit models of pragmatic reasoning, with Rational Speech Act (RSA) models providing a paradigmatic case. RSA models consider utterance interpretation to be a process of updating probability distributions over logically characterized meanings (Goodman and Stuhlmüller, 2013; Lassiter and Goodman, 2013; Goodman and Frank, 2016; Lassiter and Goodman, 2017). In doing so, they aim to capture a central feature of discourse known since the work of Grice 1975; namely, that it is constrained by principles of appropriate social behavior, which, through the reasoning of interlocutors, serve to enrich the very meanings which are communicated.

The present work provides a general approach to the integration of formal semantics with probabilistic reasoning — one which is both modular and conservative. Past efforts to consider linguistic meaning probabilistically (including

the seminal work of Goodman and Lassiter) have drastically modified the underlying logic to express it, typically in a way that freely mixes a logical semantics with probabilities. Goodman and Lassiter (2015), for example, encode meanings using the probabilistic programming language Church (Goodman et al, 2008), a decision which constitutes a radical departure from formal semantics in the style of Montague: while the latter uses a pure λ -calculus, Church programs can invoke probabilistic effects, i.e., by sampling or marginalizing, at any point in a given program.

In contrast to this and similar approaches to probabilistic semantics, ours allows for the usual approach to compositional semantics, in terms of a *pure* logical language. Moreover, we expect that our approach will be quite general: it allows for any simply typed language with products and a finite number of non-logical constants, but it should also be compatible with more expressive systems, e.g., System F (Girard, 1972) and dependent type theory (we defer an investigation of its generality, however). Our trick is to treat probabilistic computation modularly, as a side effect, using continuations. Doing so allows logical meanings to be viewed as values computed by probabilistic programs. Even so, as we shall see, our semantics does not overstep the tight bounds of the λ -calculus: probabilistic programs are *themselves* expressed in higher-order logic; we can thus provide an expressive probabilistic compositional semantics without the use of radically novel tools.

2 Formal semantics

To illustrate, we provide a schematic English fragment, which we translate into a higher-order language with types for individuals (e), truth values (t), and real numbers (r). In addition to function types ($\alpha \rightarrow \beta$), we assume access to products ($\alpha \times \beta$ and unit type \diamond), along with their associated constructors $\langle M, N \rangle : \alpha \times \beta$ (for $M : \alpha$ and $N : \beta$), destructors $M_1 : \alpha$ and $M_2 : \beta$ (for $M : \alpha \times \beta$), and unit $\diamond : \diamond$, as well as n -ary generalizations of these. Notably, we employ an indicator function $\mathbb{1} : t \rightarrow r$ taking \top ('true') and \perp ('false') onto 1 and 0, respectively. We additionally assume the existence of a family d_i of subtypes of r corresponding to degree types. For instance d_{tall} represents degrees of height, d_{happy} degrees of happiness, etc.

In general, we assume the language to have a finite number of non-logical constants. For our example, we employ the following:

$$\text{human} : e \rightarrow t \quad \text{height} : e \rightarrow d_{tall} \quad \theta_{tall} : d_{tall} \quad (\geq) : r \rightarrow r \rightarrow t$$

In terms of these, the following meanings can be given for *someone*, *is*, and *tall*, in order to derive the meaning of *someone is tall* via functional application:

$$\begin{aligned} \llbracket \text{someone} \rrbracket &= \lambda k. \exists x : \text{human}(x) \wedge k(x) \\ \llbracket \text{is} \rrbracket &= \lambda x. x \\ \llbracket \text{tall} \rrbracket &= \lambda x. \text{height}(x) \geq \theta_{tall} \end{aligned}$$

The meaning of *someone is tall*, $\llbracket \text{someone} \rrbracket (\llbracket \text{is} \rrbracket (\llbracket \text{tall} \rrbracket))$, can then be computed to be $\exists x : \text{human}(x) \wedge \text{height}(x) \geq \theta_{\text{tall}}$.

3 The traditional interpretation

We have so far illustrated the use of a language with an intuitively clear meaning. For completeness, we spell out the “traditional” meaning of this language precisely in terms of an interpretation function, $\llbracket \cdot \rrbracket$, given by a λ -homomorphism:¹

$$\begin{aligned}
 \llbracket x \rrbracket &= x && (x \text{ is a variable}) \\
 \llbracket \lambda x.M \rrbracket &= \lambda x. \llbracket M \rrbracket \\
 \llbracket MN \rrbracket &= \llbracket M \rrbracket \llbracket N \rrbracket \\
 \llbracket \langle M, N \rangle \rrbracket &= \langle \llbracket M \rrbracket, \llbracket N \rrbracket \rangle \\
 \llbracket M_i \rrbracket &= \llbracket M \rrbracket_i \\
 \llbracket \theta_{\text{tall}} \rrbracket &= d \\
 \llbracket \text{height} \rrbracket &= \text{height} \\
 \llbracket \text{human} \rrbracket &= \text{human} \\
 \llbracket (\geq) \rrbracket &= (\geq)
 \end{aligned}$$

Here, d is some real number representing the contextual standard of height used by the adjective *tall*, *height* is some function from individuals to real numbers, and *human* is some function from individuals to truth values. The \geq symbol corresponds to the “greater-than-or-equal-to” relation on real numbers. To save space, we have left implicit the interpretation of logical constants, like \top , \perp , and \exists , which is standard. One can now compose $\llbracket \cdot \rrbracket$ with $\llbracket \cdot \rrbracket$ and map *someone is tall* onto $\exists x : \text{human}(x) \wedge \text{height}(x) \geq d$, i.e., a truth value.

4 The probabilistic interpretation

We provide the probabilistic interpretation in two steps. First, we parameterize our interpretation function, $\llbracket \cdot \rrbracket^\kappa$, by a tuple κ of values, which we call a *context*. The idea is to use such a tuple to provide an interpretation for the non-logical constants. In particular, we assume that the n non-logical constants of the language are ordered, such that, if constant c_i has type α_i , then $\kappa : \alpha_1 \times \dots \times \alpha_n$.

¹ The λ -homomorphisms that we employ map one higher-order language into another, preserving variables, abstractions, applications, pairing, and projection. They are accompanied by type-homomorphisms $\bar{\cdot}$ which, for us, preserve implication and products (i.e., $\overline{\alpha \rightarrow \beta} = \bar{\alpha} \rightarrow \bar{\beta}$ and $\overline{\alpha \times \beta} = \bar{\alpha} \times \bar{\beta}$), but which may affect base types. In general, if $M : \alpha$, then $\llbracket M \rrbracket : \bar{\alpha}$. The motivation for these constraints is that they provide meanings to the constants of the source language, leaving the surrounding λ -calculus unaffected (as analogous to a traditional model-theoretic interpretation). In this case, both $\llbracket \cdot \rrbracket$ and its associated type homomorphism are trivial, mapping both constants and base types into (notational variants of) themselves.

For any such κ , $(\cdot)^\kappa$ is the following λ -homomorphism:

$$\begin{aligned}
 (x)^\kappa &= x && (x \text{ is a variable}) \\
 (\lambda x.M)^\kappa &= \lambda x.(M)^\kappa \\
 (MN)^\kappa &= (M)^\kappa (N)^\kappa \\
 (\langle M, N \rangle)^\kappa &= \langle (M)^\kappa, (N)^\kappa \rangle \\
 (M_i)^\kappa &= (M)_i^\kappa \\
 (c_i)^\kappa &= \kappa_i && (c_i \text{ is the } i^{\text{th}} \text{ constant})
 \end{aligned}$$

Thus if c_i is one of θ_{tall} , **height**, or \geq , then its interpretation is determined by the context. (Again, we have omitted the interpretation of logical constants to save space.) Obviously, if c_i is of type α_i , then so is κ_i . We assume that all probabilistic semantic knowledge resides in the interpretation of constants, that is, in the context. It remains to be shown how to evaluate the above expressions when κ is a random variable.

5 Probabilistic programs

In general, we consider something a random variable if it is the value returned by a probabilistic program. In our framework, a probabilistic program returning values of type α is a function of type $(\alpha \rightarrow r) \rightarrow r$; that is, one which consumes a *projection* function (from values of type α to real numbers), in order to return a real number.² The intent is that if p is a probabilistic program and f is a projection function, then $p(f)$ is the sum of $f(x)$, for all possible values of x returned by the program, weighted in proportion to their probabilities.

Given this setup, probabilistic programs form a *monad*. A monad, as stated in Figure 1, is a functor M from types to types, associated with two operators, η (‘unit’) and \star (‘bind’), satisfying certain laws. In general, implementing a monad in a pure setting, such as the λ -calculus, allows one to simulate various notions of side effect, including probabilistic computation, as we shall see. The role of η is to inject an ordinary value into the monad, while that of \star is to compose *computations*. More precisely, \star runs a computation of type $M\alpha$, and then binds the returned value to a variable in the next computation (something of type $\alpha \rightarrow M\beta$). In the case of probabilistic programs, $M\alpha = (\alpha \rightarrow r) \rightarrow r$, and the

² There is some precedent for this representation of probabilistic programs, by Mohammed Ismail and Shan (2016), who describe a small typed probabilistic programming language and provide a denotational semantics for it in terms of continuations. Our formulation is chiefly inspired by the dependently typed language of Bernardy et al (2021). See also Jansson et al 2021.

Operators

$$\begin{aligned} \eta &: \alpha \rightarrow M\alpha \\ (\star) &: M\alpha \rightarrow (\alpha \rightarrow M\beta) \rightarrow M\beta \end{aligned}$$

Laws on terms

$$\begin{aligned} \eta(v) \star k &= k(v) && \text{(Left Identity)} \\ m \star \eta &= m && \text{(Right Identity)} \\ (m \star n) \star o &= m \star (\lambda x. n(x) \star o) && \text{(Associativity)} \end{aligned}$$

Fig. 1. Definition of a monad

unit η and bind operator \star are inherited from the continuation monad:

$$\begin{aligned} \eta &: \alpha \rightarrow (\alpha \rightarrow r) \rightarrow r \\ \eta(a) &= \lambda c. ca \\ (\star) &: ((\alpha \rightarrow r) \rightarrow r) \\ &\quad \rightarrow (\alpha \rightarrow (\beta \rightarrow r) \rightarrow r) \\ &\quad \rightarrow (\beta \rightarrow r) \rightarrow r \\ m \star k &= \lambda c. m(\lambda x. k(x)(c)) \end{aligned}$$

Employing the monadic operators, one may sequence a program $p : (\alpha \rightarrow r) \rightarrow r$ with some projection function $f : \alpha \rightarrow r$, producing:

$$p \star \lambda x. \eta(f(x)) : (r \rightarrow r) \rightarrow r$$

Indeed, feeding the identity function to the result obtains $p(f) : r$.

The encoding of probabilistic programs in terms of continuations may appear somewhat indirect. In general, one can see a continuation (here, of type $\alpha \rightarrow r$) as a question to ask a program. The result type r restricts the sorts of questions one may ask, i.e., to those having real numbers as answers. As if responding with a riddle, moreover, the probabilistic program returns the weighted sum of the answers for its possible values. This means, for instance, that given a probabilistic program p , one may feed it the question $(\lambda x. 1)$, which asks how much mass it assigns in total, in order to get the answer, $p(\lambda x. 1)$, which returns the total mass assigned by p .

If p returns truth values (i.e., is of type $(t \rightarrow r) \rightarrow r$), we can ask for the mass it assigns to \top by passing the indicator function as a continuation: $p(\mathbb{1})$. Consequently, we may compute a *probability* for p as the expected value of $\mathbb{1}$, in terms of a function $P : ((t \rightarrow r) \rightarrow r) \rightarrow r$:

$$P(p) = \frac{p(\mathbb{1})}{p(\lambda b. 1)}$$

The denominator (the total mass assigned by p) normalizes the result.

Now, let K be a probabilistic program representing the distribution of contexts; that is, if $\alpha_1 \times \dots \times \alpha_n$ is the type of contexts, K is of type $(\alpha_1 \times \dots \times \alpha_n \rightarrow r) \rightarrow r$. Given a term ϕ of type t , we encode its interpretation in the context of K as the following probabilistic program:

$$K \star \lambda \kappa. \eta(\llbracket \phi \rrbracket^\kappa)$$

Like all probabilistic programs returning truth values, the above is of type $(t \rightarrow r) \rightarrow r$. Operationally, it reads in the random context returned by K and computes from it a truth value for ϕ in this context. As such, one can determine a probability for it, as outlined above.

To illustrate, consider our running example, *someone is tall*, to which we assigned the interpretation $\exists x : \text{human}(x) \wedge \text{height}(x) \geq \theta_{tall}$. Let us assume a probabilistic program K returning contexts where the interpretation of θ_{tall} is a random variable having a normal distribution with a mean of 72 inches and a standard deviation of 3 inches. Moreover, we assume that the interpretations of the other constants are fixed as the functions *height*, *human*, and \geq , as above. Then, assuming the order of the constants to be *height*, *human*, \geq , θ_{tall} , we have the following definition of K :

$$K = \mathcal{N}(72, 3) \star \lambda d. \eta(\text{height}, \text{human}, (\geq), d)$$

Here, $\mathcal{N}(72, 3)$ is a probabilistic program (of type $(d_{tall} \rightarrow r) \rightarrow r$) representing a normal distribution with the relevant mean and standard deviation. If fed a projection function f of type $d_{tall} \rightarrow r$, it will integrate f over the real line, weighting each $f(d)$ by the probability of d .³ Our strategy allows us to associate a probabilistic program with the sentence *someone is tall*, as follows:

$$\begin{aligned} & K \star \lambda \kappa. \eta(\llbracket \exists x : \text{human}(x) \wedge \text{height}(x) \geq \theta_{tall} \rrbracket^\kappa) \\ &= K \star \lambda \kappa. \eta(\llbracket \exists x : (\llbracket \text{human} \rrbracket^\kappa(x) \wedge \llbracket (\geq) \rrbracket^\kappa(\llbracket \text{height} \rrbracket^\kappa(x)) \rrbracket^\kappa) \rrbracket^\kappa) \\ &= \mathcal{N}(72, 3) \star \lambda d. \eta(\text{height}, \text{human}, (\geq), d) \star \lambda \kappa. \eta(\llbracket \exists x : \kappa_2(x) \wedge \kappa_3(\kappa_1(x)) \rrbracket^\kappa) \\ &= \mathcal{N}(72, 3) \star \lambda d. \eta(\llbracket \exists x : \text{human}(x) \wedge \text{height}(x) \geq d \rrbracket^\kappa) \\ &= \lambda c. \mathcal{N}(72, 3)(\lambda d. c(\llbracket \exists x : \text{human}(x) \wedge \text{height}(x) \geq d \rrbracket^\kappa)) \end{aligned}$$

As expected, we have a program of type $(t \rightarrow r) \rightarrow r$. We may therefore compute a probability for it as:

$$\frac{\mathcal{N}(72, 3)(\lambda d. \mathbf{1}(\llbracket \exists x : \text{human}(x) \wedge \text{height}(x) \geq d \rrbracket^\kappa))}{\mathcal{N}(72, 3)(\lambda d. \mathbf{1})}$$

³ We leave $\mathcal{N} : d_{tall} \times d_{tall} \rightarrow (d_{tall} \rightarrow r) \rightarrow r$ unanalyzed. In general, computing a continuous distribution $\mathcal{D} : p_1 \times \dots \times p_n \rightarrow (d \rightarrow r) \rightarrow r$ over d amounts to computing

$$\lambda \langle p_1, \dots, p_n \rangle, f. \int_{-\infty}^{\infty} \text{PDF}_{\mathcal{D}(p_1, \dots, p_n)}(x) \star f(x) dx$$

where $\text{PDF}_{\mathcal{D}(p_1, \dots, p_n)}$ provides the probability density function associated with \mathcal{D} (given parameters p_1, \dots, p_n). Such integrals don't in general admit closed-form solutions, and so one must resort to approximations. We implement this via Markov chain Monte Carlo sampling in our Haskell implementation, using the library at <https://github.com/jyp/ProbProg>.

Because \mathcal{N} represents a genuine probability distribution, its total mass is 1, and we can simply ignore the denominator. The theoretical value of this expression is determined by computing the truth of the proposition that someone's height exceeds the height threshold d at every possible value of d , and weighting it by the probability associated with d . This model of the uncertainty associated with *someone is tall* locates it in the meaning of *tall*; in particular, how tall one must be to be considered tall.

For example, consider a case in which exactly one person is 72 inches tall and no one is taller. Then the condition imposed by the meaning of *someone is tall* will be met by all $\theta_{tall} \leq 72$, and the sentence will be assigned the probability 0.5. In general, the probability assigned will be equal to the mass of $\mathcal{N}(72, 3)$ that is less than or equal to the height of the tallest human.

6 Bayesian inference

One of the main interests of a probabilistic semantics such as the one proposed is that it can be combined with Bayesian marginalization. For this purpose, we define the following function *observe*:

$$\begin{aligned} \text{observe} &: t \rightarrow (\diamond \rightarrow r) \rightarrow r \\ \text{observe } \phi f &= \mathbf{1}(\phi) * f(\diamond) \end{aligned}$$

Given a proposition ϕ , *observe* either keeps or throws out its continuation, as according to whether ϕ is true or false; hence, the resulting program retains only values from the part of the distribution it represents compatible with ϕ being true.⁴ This function thus allows us to marginalize the truth value of ϕ under a premise ψ as follows, exploiting the monadic structure of probabilistic programs:

$$K \star \lambda\kappa. \text{observe}(\langle \psi \rangle^\kappa) \star \lambda\diamond. \eta(\langle \phi \rangle^\kappa)$$

Such a marginalization process can be used for several purposes: for probabilistic inference (as suggested by our running example), but also to refine the probability distributions associated with constants; that is, for semantic learning. We briefly suggest how each of these tasks can be accomplished in our framework, starting with semantic learning.

6.1 Semantic learning

Semantic learning in our framework is matter of updating (distributions of) contexts. Given a program K_0 returning contexts which represents the initial state of one's semantic knowledge, one may observe a number of propositions to be true or false, thus obtaining a new program, K_1 :

$$K_1 = K_0 \star \lambda\kappa. \text{observe}(\phi_1) \star \lambda\diamond. \dots \text{observe}(\phi_n) \star \lambda\diamond. \eta(\kappa)$$

⁴ Some may recognize it as akin to the *guard* function of Haskell's MonadPlus and Alternative classes.

The effect of sequencing K_0 with such a series of observations is to zero out the portion of its distribution in which ϕ_1, \dots, ϕ_n are false, returning the values that survive.

Let us say that a learner is attempting to learn the meaning of *tall*, and they start out with a distribution of contexts such that the height threshold that the adjective makes use of ranges over a normal distribution with a mean of 68 inches and a standard deviation of 3 inches (we will deal here with the constants *height*, \geq , and θ_{tall} , along with the four names for individuals *c*, *m*, *a*, and *v*):

$$K_0 = \mathcal{N}(68, 3) \star \lambda d. \eta(c, m, a, v, \textit{height}, (\geq), d)$$

In addition, this learner happens to know the following three facts: that Camilla is 65 inches tall, that Matt is 67 inches tall, and that Anna is 72 inches tall:

$$\textit{height}(c) = 65 \quad \textit{height}(m) = 67 \quad \textit{height}(a) = 72$$

One day, someone this learner trusts utters the following three sentences, in sequence: (1) *Camilla isn't tall*, (2) *Matt isn't tall*, (3) *Anna is tall*. Upon hearing these utterances, the learner updates K_0 , in order to obtain K_1 :

$$\begin{aligned} K_1 &= K_0 \\ &\quad \star \lambda \kappa. \textit{observe}(\lnot \textit{height}(c) \geq \theta_{tall})^\kappa \\ &\quad \star \lambda \diamond. \textit{observe}(\lnot \textit{height}(m) \geq \theta_{tall})^\kappa \\ &\quad \star \lambda \diamond. \textit{observe}(\textit{height}(a) \geq \theta_{tall})^\kappa \\ &\quad \star \lambda \diamond. \eta(\kappa) \\ &= \mathcal{N}(68, 3) \quad (\text{by Associativity and Left Identity}) \\ &\quad \star \lambda d. \textit{observe}(-65 \geq d) \\ &\quad \star \lambda \diamond. \textit{observe}(-67 \geq d) \\ &\quad \star \lambda \diamond. \textit{observe}(72 \geq d) \\ &\quad \star \lambda \diamond. \eta(c, m, a, v, \textit{height}, (\geq), d) \end{aligned}$$

This may in turn be simplified to:

$$\mathcal{N}(68, 3) \star \lambda d. \textit{observe}(72 \geq d \wedge d > 67) \star \lambda \diamond. \eta(c, m, a, v, \textit{height}, (\geq), d)$$

Thus K_1 is just like K_0 , but for the fact that the distribution associated with θ_{tall} has been pared down to only include the mass of $\mathcal{N}(68, 3)$ in the interval $(67, 72]$. If Vlad is 68 inches tall ($\textit{height}(v) = 68$), then the sentence *Vlad is tall* would have been associated with the probability 0.5 in K_0 , while it is associated with a probability of around 0.24 in K_1 :

$$\begin{aligned} \frac{K_0(\lambda \kappa. \mathbf{1}(\textit{height}(v) \geq \theta_{tall})^\kappa)}{K_0(\lambda \kappa. \mathbf{1})} &= 0.5 \\ \frac{K_1(\lambda \kappa. \mathbf{1}(\textit{height}(v) \geq \theta_{tall})^\kappa)}{K_1(\lambda \kappa. \mathbf{1})} &\approx 0.24 \end{aligned}$$

6.2 RSA: background

In the case of probabilistic inference, our framework can serve as the basis for complex pragmatic reasoning, as in RSA models. For example, Lassiter and Goodman (2013) present an RSA model of the inference made when someone utters a sentence such as *Vlad is tall*. This model consists of a pragmatic listener (L_1), who reasons about probable meanings based on the expected behavior of a pragmatic speaker (S_1), who, in turn, reasons about a literal listener (L_0). These agents’ behaviors are modeled in terms of the following equations (adapted to the current example):

$$P_{L_1}(h, d_{tall} \mid \text{‘Vlad is tall’}) \propto P_{S_1}(\text{‘Vlad is tall’} \mid h, d_{tall}) * P_{L_1}(h) \quad (L_1)$$

$$P_{S_1}(u \mid h, d_{tall}) \propto (P_{L_0}(h \mid u, d_{tall}) * e^{-C(u)})^\alpha \quad (S_1)$$

$$P_{L_0}(h \mid u, d_{tall}) = P_{L_0}(h \mid \llbracket u \rrbracket^{d_{tall}} = \top) \quad (L_0)$$

Each of these statements defines a probability distribution for the random variables of interest. h and d_{tall} are values of the random variables representing Vlad’s height and the height threshold for the adjective *tall*, respectively. The function C in the S_1 model is utterance cost. α is the “temperature” of the S_1 model: it controls the extent to which the speaker behaves rationally, i.e., by taking the expected behavior of the literal listener L_0 , as well as utterance cost, into account in designing their distribution over utterances.

Given the more general notions of a world state w and a parameter θ , (h and d_{tall} , respectively, in the above), these equations may be presented more perspicuously as follows, given some utterance u_0 :

$$P_{L_1}(w, \theta \mid u_0) = \frac{P_{S_1}(u_0 \mid w, \theta) * P_{L_1}(w, \theta)}{\int_{w' \in W} \int_{\theta' \in \Theta} P_{S_1}(u_0 \mid w', \theta') * P_{L_1}(w', \theta') d\theta' dw'} \quad (L_1)$$

$$P_{S_1}(u \mid w, \theta) = \frac{(P_{L_0}(w \mid u, \theta) * e^{-C(u)})^\alpha}{\sum_{u' \in U} (P_{L_0}(w \mid u', \theta) * e^{-C(u')})^\alpha} \quad (S_1)$$

$$P_{L_0}(w \mid u, \theta) = P_{L_0}(w \mid \llbracket u \rrbracket^\theta = \top) \quad (L_0)$$

Thus abstractly, pragmatic listeners provide a joint posterior distribution over world states w and parameters θ , given an utterance u_0 .⁵ Pragmatic speakers provide a distribution of utterances, given the particular world state w (and parameter θ) they wish to communicate. These utterances, moreover, are taken from an antecedently chosen set U of possible utterances, which is generally assumed to be finite, thus justifying the use of summation in the normalizing factor for S_1 . Finally, linguistic uncertainty is represented by the parameter θ , which is passed from the pragmatic listener L_1 down to the literal listener L_0 , through the speaker model S_1 . Note, therefore, that L_1 differs from L_0 in a

⁵ Note that we define this posterior in terms of a prior joint distribution. Lassiter and Goodman (2013) assume the prior distributions over world states and linguistic parameters to be independent, with an effectively uniform prior over parameters.

crucial respect: while L_1 samples both world states and parameters, L_0 samples only world states, relying on a parameter which has been fixed by L_1 .

6.3 RSA: implementation

Our purpose is to illustrate how the RSA framework may be realized in the vocabulary of probabilistic programs. Thus taking u to be the type of utterances, σ the type of world states, and π the type of linguistic parameters, we aim to find a program L_1 of type $u \rightarrow (\sigma \times \pi \rightarrow r) \rightarrow r$, which, given an utterance, provides a joint distribution over world states and parameters, and which satisfies the desiderata laid out above. In order to do so, it is useful to introduce the following generalization of *observe* to fuzzy conditions:

$$\begin{aligned} \mathit{factor} &: r \rightarrow (\diamond \rightarrow r) \rightarrow r \\ \mathit{factor} \ x \ f &= x * f(\diamond) \end{aligned}$$

Instead of a truth value, *factor* takes a real number and applies it as a weight to the result of its continuation. Thus *observe* may be viewed as the specific case of *factor* in which the relevant weight is either 1 or 0.

Now, we may formulate L_1 as follows. Say that S_1 provides a probabilistic program returning *utterances*, given a world state and a parameter; i.e., it is of type $\sigma \times \pi \rightarrow (u \rightarrow r) \rightarrow r$. Then given some w and θ , we would like access to the probability density function corresponding to $S_1(w, \theta)$ — PDF $_{S_1(w, \theta)}$ — of type $u \rightarrow r$, so that we may appropriately factor the probability of $\langle w, \theta \rangle$ in L_1 , given an utterance. (We will come back to how we obtain the PDFs of probabilistic programs shortly. For now, we simply take them for granted.) Moreover, let us assume that world states and parameters take prior distributions $W : (\sigma \rightarrow r) \rightarrow r$ and $\Theta : (\pi \rightarrow r) \rightarrow r$, respectively. These assumptions leave us with the following definition of L_1 :

$$L_1(u_0) = W \star \lambda w. \Theta \star \lambda \theta. \mathit{factor}(\text{PDF}_{S_1(w, \theta)}(u_0)) \star \lambda \diamond. \eta(w, \theta)$$

Now, given some prior distribution U over utterances (i.e., of type $(u \rightarrow r) \rightarrow r$), we may similarly provide definitions of S_1 and L_0 :

$$\begin{aligned} S_1 &: \sigma \times \pi \rightarrow (u \rightarrow r) \rightarrow r \\ S_1(w, \theta) &= U \star \lambda u. \mathit{factor}(\text{PDF}_{L_0(u, \theta)}(w) * e^{-C(u)}^\alpha) \star \lambda \diamond. \eta(u) \\ L_0 &: u \times \pi \rightarrow (\sigma \rightarrow r) \rightarrow r \\ L_0(u, \theta) &= W \star \lambda w. \mathit{observe}(\langle u \rangle^{(w, \theta)}) \star \lambda \diamond. \eta(w) \end{aligned}$$

Note our use of notation in the definition of L_0 : $\langle w, \theta \rangle$ is assumed to provide a context in which we can interpret the utterance u , which we encode as a formula. Thus the relevant λ -homomorphism maps the type u onto t ($\bar{u} = t$).

Having stated our formulation of RSA somewhat abstractly, let us now turn to the problem of PDFs; that is, of obtaining a PDF of type $\alpha \rightarrow r$ from a

probabilistic program of type $(\alpha \rightarrow r) \rightarrow r$. If α is finite, we may construct this function as follows (recall that P takes a probabilistic program of type $(t \rightarrow r) \rightarrow r$ into a probability):⁶

$$\begin{aligned} \text{PDF}_{(\cdot)} &: ((\alpha \rightarrow r) \rightarrow r) \rightarrow \alpha \rightarrow r \\ \text{PDF}_p &= \lambda x. P(p \star \lambda y. \eta(y = x)) \end{aligned}$$

That is, for every $x : \alpha$, $\text{PDF}_p(x)$ evaluates the probability that p returns x .

If α is continuous, however, we have a problem: the probability that any two values x and y are equal is zero, and the above definition of a PDF would have it return zero everywhere! Fortunately, there are sound remedies which we may adopt. For instance, we may take the derivative of the cumulative mass of a given distribution p with respect to the argument:

$$\text{PDF}_p = \lambda x. \frac{d}{dx} [P(p \star \lambda y. \eta(y \leq x))]$$

Indeed, these two definitions of PDF may be plugged into the definitions above, depending on whether the type of the relevant argument is finite (e.g., an utterance) or continuous, in order to obtain a fuller specification of L_1 .⁷ One need only determine what the distributions U , W , and Θ are. To realize the model of Lassiter and Goodman (2013), we would take U to be a small finite set, W to be a normal distribution, and Θ to be, effectively, uniform.⁸ The resulting probabilistic program can be computed approximately using Monte Carlo methods; in this case, one will typically evaluate a probabilistic program to an approximate, finite PDF.

We close out this section by observing a noteworthy feature of the foregoing formulation of RSA: it highlights an odd lack of symmetry between the L_1 model and the L_0 model. Why does L_1 sample both world states from W and linguistic parameters from Θ , while L_0 samples only the former? Indeed, this fact is reflected in their types! L_1 is of type $u \rightarrow (\sigma \times \pi \rightarrow r) \rightarrow r$: it takes an utterance and returns a distribution over pairs of world states and parameters. Meanwhile, L_0 is of type $u \times \pi \rightarrow (\sigma \rightarrow r) \rightarrow r$: it takes a pair of an utterance *and* a parameter and returns a distribution over world states. Thus L_0 considers a linguistic parameter which has been *fixed* by L_1 and S_1 . Put differently, S_1 reasons about an L_0 that knows θ ahead of time, when determining what to say.

⁶ When α is finite, one generally speaks of a probability mass function. We use the term ‘PDF’ in a generic way, however.

⁷ An alternative, syntactically closer to the discrete case, relies on the Dirac δ distribution, whose value is zero everywhere except when its argument is zero, and whose total mass sums to one. Thus we recover a non-zero result after integration:

$$\text{PDF}_p = \lambda x. p(\lambda y. \delta(x - y))$$

⁸ More accurately, we would take U to be uniform over a finite set, S_U . Thus we would define it as $U = \lambda k. \sum_{u \in S_U} k(u)$.

Yet more vividly, the pragmatic listener assumes that the speaker is under the impression that the two have already (telepathically, perhaps) coordinated on linguistic parameters.

Maybe, it is more realistic not to assume that S_1 imagines such an omniscient L_0 . In fact, relaxing this assumption restores the symmetry of the model. At the same time, it conveniently allows us not to explicitly split the context κ into two parts w and θ . As in previous sections, we assume that the context has some type $\kappa = \alpha_1 \times \dots \times \alpha_n$:

$$\begin{aligned} L_1 &: u \rightarrow (\kappa \rightarrow r) \rightarrow r \\ L_1(u) &= K \star \lambda\kappa.\mathit{factor}(\text{PDF}_{S_1(\kappa)}(u)) \star \lambda\circ.\eta(\kappa) \\ S_1 &: \kappa \rightarrow (u \rightarrow r) \rightarrow r \\ S_1(\kappa) &= U^* \star \lambda u.\mathit{factor}(\text{PDF}_{L_0(u)}(\kappa)^\alpha) \star \lambda\circ.\eta(u) \\ L_0 &: u \rightarrow (\kappa \rightarrow r) \rightarrow r \\ L_0(u) &= K \star \lambda\kappa.\mathit{observe}((u)^\kappa) \star \lambda\circ.\eta(\kappa) \end{aligned}$$

To simplify the presentation, we have used the notation U^* to stand for a distribution over utterances which has already incorporated a notion of cost.⁹ In our final formulation, both L_1 and L_0 have the same type. There is thus a more general notion of “listener”, corresponding to a family of maps from utterances to distributions over contexts (or, equivalently, joint distributions over world states and linguistic parameters). Such listeners differ only in how they marginalize the prior: the literal listener uses a literal interpretation, while the pragmatic listener uses a pragmatic interpretation. Such pragmatic interpretations arise from the speaker model, which chooses utterances which best fit the state of the world that it wishes to communicate.¹⁰

In summary, we have a realization of RSA that is highly compositional, in two senses. First, the models themselves are assembled compositionally in terms of probabilistic programs and monadic combinators. Second, utterances, represented by logical formulae, are interpreted compositionally, and such formulae may be obtained from natural language sentences using standard compositional techniques. At the same time, the mathematical vocabulary for describing RSA models is one and the same as that for describing linguistic meanings.

7 Conclusion

Our aim has been to lay a strong foundation for compositional probabilistic semantics. Many details have been left out, including about how one might represent prior knowledge, concretely. Many possibilities arise here. For instance,

⁹ To implement the definition of cost employed by RSA models, for example, U^* could be $U \star \lambda u.\mathit{factor}(e^{-\alpha \cdot C(u)}) \star \lambda\circ.\eta(u)$, given some uniform distribution U .

¹⁰ Systematically, if α tends to ∞ ; probabilistically, otherwise.

one may follow machine-learning methods and use vectors to represent individuals (Bernardy et al, 2021), while predicates are represented by hyperplanes in the relevant space (Bernardy et al, 2019). An alternative would encode prior knowledge in terms of the same logic used to represent meanings, i.e., as sets of formulae. One may then constrain distributions over contexts in terms of such formulae (Grove et al, 2021). Following this route, one may obtain a seamless integration of Bayesian and logical representations of knowledge.

We should note that, while the logical fragments provided here are rudimentary, they are also merely expository: there is no deep reason that we did not provide a richer semantics, e.g., incorporating dynamism (following the tradition of combining dynamic semantics with the simply typed λ -calculus). Indeed, one could combine the framework we have illustrated with a logical semantics that *itself* uses continuations (Barker and Shan, 2014; Lebedeva, 2012) or monads (Charlow, 2014).

Finally, while our contribution is chiefly a theoretical one, the system described in this paper has been implemented using the Haskell programming language (available at <https://github.com/juliangrove/grove-bernardy-lens18>). The mathematical vocabulary that we have employed here to assemble expressions of type r is closely mirrored by the implementation in terms of a domain-specific language for characterizing Markov chain Monte Carlo sampling procedures. Thus while many of the probabilistic programs provided above cannot be evaluated to closed-form solutions, they may all be finitely approximated, given sufficiently many samples. Most important, however, the modular division of labor between logical expressions and probabilistic side effects is straightforward to simulate in Haskell, given the pure functional setting it provides.

We have shown that a probabilistic semantics of natural language is amenable to a fully formal treatment — one which remains squarely within the realm of pure typed λ -calculi. The key idea is to use an effect system to capture probabilistic operations (i.e., sampling and marginalization). Our approach fits the general framework of monadic semantics, and, as such, augments a literature that has grown in many exciting ways since the work of Shan (2002).

Bibliography

- Barker C, Shan Cc (2014) Continuations and natural language, vol 53. Oxford studies in theoretical linguistics
- Bernardy JP, Blanck R, Chatzikyriakidis S, Lappin S, Maskharashvili A (2019) Predicates as Boxes in Bayesian Semantics for Natural Language. In: Proceedings of the 22nd Nordic Conference on Computational Linguistics, Linköping University Electronic Press, Turku, Finland, pp 333–337
- Bernardy JP, Blanck R, Chatzikyriakidis S, Maskharashvili A (2021) Bayesian Natural Language Semantics and Pragmatics. In: Bernardy JP, Blanck R, Chatzikyriakidis S, Lappin S, Maskharashvili A (eds) Probabilistic Approaches to Linguistic Theory, CSLI Publications
- Charlow S (2014) On the semantics of exceptional scope. PhD Thesis, NYU, New York, URL <https://semanticsarchive.net/Archive/2JmMWRjY>

- Girard JY (1972) Interprétation fonctionnelle et élimination des coupures de l'arithmétique d'ordre supérieur. PhD thesis, Université Paris 7
- Goodman ND, Frank MC (2016) Pragmatic Language Interpretation as Probabilistic Inference. *Trends in Cognitive Sciences* 20(11):818–829
- Goodman ND, Lassiter D (2015) Probabilistic Semantics and Pragmatics Uncertainty in Language and Thought. In: Lappin S, Fox C (eds) *The Handbook of Contemporary Semantic Theory*, John Wiley & Sons, Ltd, pp 655–686
- Goodman ND, Stuhlmüller A (2013) Knowledge and Implicature: Modeling Language Understanding as Social Cognition. *Topics in Cognitive Science* 5(1):173–184
- Goodman ND, Mansinghka VK, Roy D, Bonawitz K, Tenenbaum JB (2008) Church: a language for generative models. In: *Proceedings of the Twenty-Fourth Conference on Uncertainty in Artificial Intelligence*, AUAI Press, Arlington, Virginia, USA, UAI'08, pp 220–229
- Grice HP (1975) Logic and Conversation. In: Cole P, Morgan JL (eds) *Syntax and Semantics*, vol 3, *Speech Acts*, Academic Press, New York, pp 41–58
- Grove J, Bernardy JP, Chatzikiyiakidis S (2021) From compositional semantics to Bayesian pragmatics via logical inference. In: *Proceedings of Natural Logic meets Machine Learning 2021*
- Jansson P, Ionescu C, Bernardy JP (2021) Probability Theory. In: *Domain Specific Languages of Mathematics*, pp 196–215, URL https://github.com/DSLsofMath/DSLsofMath/blob/master/L/snapshots/DSLsofMathBook_snapshot_2021-03-06.pdf
- Lassiter D, Goodman ND (2013) Context, scale structure, and statistics in the interpretation of positive-form adjectives. *Semantics and Linguistic Theory* 23(0):587–610, number: 0
- Lassiter D, Goodman ND (2017) Adjectival vagueness in a Bayesian model of interpretation. *Synthese* 194(10):3801–3836
- Lebedeva E (2012) Expressing discourse dynamics through continuations. PhD thesis, Université de Lorraine, Lorraine
- Mohammed Ismail W, Shan Cc (2016) Deriving a probability density calculator (functional pearl). In: *Proceedings of the 21st ACM SIGPLAN International Conference on Functional Programming*, Association for Computing Machinery, New York, NY, USA, ICFP 2016, pp 47–59
- Shan Cc (2002) Monads for natural language semantics. In: Striegnitz K (ed) *Proceedings of the ESSLLI 2001 Student Session*, 13th European Summer School in Logic, Language, and Information, Helsinki