
Qualitative Probabilistic Programming

Anonymous Author(s)

Affiliation

Address

email

Abstract

In probabilistic programs, sometimes it is difficult to specify the correct parameterized family of distributions. We explore an extension to probabilistic programming languages that allows programmers to mark some distributions as unspecified. Then, we can fill in the distribution with some family and infer parameters.

1 Introduction

By separating model specification and inference, probabilistic programming has made it easier for non-experts to implement and use probabilistic models. Practitioners frequently have strong intuitions about the *structure* of their domain knowledge, such as which latent variables exist and what their causal relations are, and probabilistic programming allows them to encode this knowledge. However, it also requires them to specify the specific parametric shape and parameterization of any distributions used, and intuitions tend to be much less precise there. We present Quipp, a system that does *not* require such specification; instead, random variables and random functions can be left undefined and will automatically be filled in under maximum entropy assumptions based on their types and available datasets.

Our formalism can concisely express a wide variety of models that machine learning practitioners care about, and we provide an expectation maximization algorithm that can learn the parameters for many of these models with reasonable efficiency. This system makes it easy for non-experts to encode their beliefs about the data and to get predictions based on as few additional assumptions as possible.

In an ordinary probabilistic programming language (such as Church), it is possible to treat parameters as random variables. This would allow ordinary inference algorithms to infer parameters. However, there are advantages of having unknown functions as a feature in the language. First, it is easier to write programs without knowing the details of different parameterized distributions. Second, the system can use specialized algorithms to infer parameters faster.

In the following, we first specify the syntax used to write Quipp programs, including the notation for unknown variables and functions. We describe the class of exponential family variables and functions that our system can learn, and present the expectation maximization algorithm used to learn them. We then demonstrate the expressiveness of our language, and the broad applicability of our algorithm, by writing some of the most common machine learning models in Quipp: clustering, naive Bayes, factor analysis, a Hidden Markov model, Latent Dirichlet Allocation, and a neural net.

2 Syntax

Quipp is implemented as a library for webppl programs. Webppl [TODO cite] is a probabilistic programming language implemented in Javascript. Quipp programs are written as webppl programs that have access to special functions.

Here is an example of a Quipp program to cluster 2d points into 3 clusters:

```

var nclusters = 3;
var dim = 2;

var n = 100;

testParamInference(function(randFunction) {
  var pointType = Vector(dim, Double);
  var getPoint = randFunction(Categorical(nclusters), pointType);
  return function() {
    return repeat(n, function(i) {
      var cluster = randomInteger(nclusters);
      return observe('point' + i, getPoint, cluster);
    });
  };
});

```

Notice that we declared two types (`PointType` and `ClusterType`) and one random function (`getPoint`). Type annotations are necessary for random functions. The type `Vector(2, Double)` expresses the fact that the points are 2-dimensional, and the type `Categorical(3)` expresses the fact that there are 3 possible clusters (so a cluster may be either 0, 1, or 2). We assume that the distribution over clusters is uniform, but we do not know the distribution of points in each cluster. We will fill in the `getPoint` function with a learned function that will take a random sample from the given cluster.

3 Family of distributions

For unknown functions, the family of random functions used is a kind of generalized linear model. We assume that the distribution of the function's return value is some exponential family whose natural parameters are determined from the arguments:

$$p_{\eta}(y|x) = \exp(\eta(x)^T \phi(y) - g(\eta(x)))$$

Here, $\eta(x)$ is the natural parameter, $\phi(y)$ is a vector of y 's sufficient statistics, and g is the log partition function.

To determine $\eta(x)$, we label some subset of the sufficient statistics of both x and y as *features*. The natural parameters corresponding to non-features are constant, while natural parameters corresponding to features are determined as an affine function of the features of the arguments. Features are the same as sufficient statistics for the categorical distribution, but while both X and X^2 are sufficient statistics for the Gaussian distribution, only X is a feature.

Let $\psi(x)$ be the features of x . Then

$$p_{\mathbf{N}}(y|x) = \exp\left(\begin{bmatrix} 1 \\ \psi(x) \end{bmatrix}^T \mathbf{N} \phi(y) - g\left(\mathbf{N}^T \begin{bmatrix} 1 \\ \psi(x) \end{bmatrix}\right)\right)$$

where \mathbf{N} is a matrix containing our parameters. It must have 0 for each entry whose row corresponds to a sufficient statistics of y that is not a feature and whose column is not 1. This ensures that only the natural parameters that are features of y are affected by x .

4 Inference

To infer both latent variables and parameters, we run the expectation maximization algorithm on the probabilistic model, iterating stages of estimating latent variables using Metropolis Hastings and inferring parameters using gradient descent.

For the expectation step, we must estimate latent variable distributions given fixed values for the parameters. To do this, we can run the Quipp program with random functions set to use these fixed parameter values to generate their results. We use the Metropolis Hastings algorithm to perform inference in this program, yielding traces. Next, for each random function, we can find all calls to it in the trace to get the training data.

For the maximization step, given samples from each random function, we set the parameters of the function to maximize the likelihood of the samples. To do this we, we use gradient descent.

Given (x, y) samples, parameter estimation to maximize log probability is a convex problem because the log probability function is concave:

$$\log p_{\mathbf{N}}(y|x) = \begin{bmatrix} 1 \\ \psi(x) \end{bmatrix}^T \mathbf{N}\phi(y) - g\left(\mathbf{N}^T \begin{bmatrix} 1 \\ \psi(x) \end{bmatrix}\right)$$

This relies on the fact that g is convex, but this is true in general for any exponential family distribution. Since the problem is convex, it is possible to use gradient descent to optimize the parameters. Although the only exponential family distributions we use in this paper are the categorical and Gaussian distributions, we can use the same algorithms for other exponential families, such as the Poisson and gamma distributions.

5 Evaluation

To evaluate performance, for each model, we:

- Randomly generate parameters θ
- Generate datasets x_{train}, x_{test} using θ
- Estimate $\log P(x_{test}|\theta)$
- Use the EM algorithm to infer approximate parameters $\hat{\theta}$ from x_{train}
- Estimate $\log P(x_{test}|\hat{\theta})$ and compare to $\log P(x_{test}|\theta)$

Estimating $\log P(x_{test}|\theta)$ is nontrivial, given that the model contains latent variables. We use the Sequential Monte Carlo algorithm for this.

6 Examples

6.1 Clustering

```
var nclusters = 3;
var dim = 2;

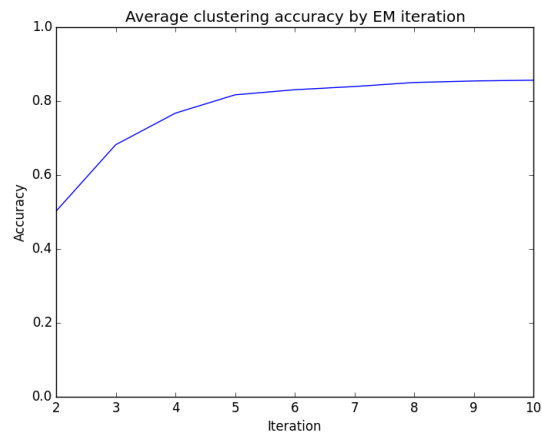
var n = 100;

testParamInference(function(randFunction) {
  var pointType = Vector(dim, Double);
  var getPoint = randFunction(Categorical(nclusters), pointType);
  return function() {
    return repeat(n, function(i) {
      var cluster = randomInteger(nclusters);
      return observe('point' + i, getPoint, cluster);
    });
  };
});
```

In this example, we cluster 2d points into 3 different clusters. Given a cluster, the distribution for a point is some independent Gaussian distribution. This is similar to fuzzy c-means clustering.

We randomly generated parameters for this example 100 times, and each time took 100 samples and then ran 10 EM iterations. The accuracy is defined as the maximum percentage of points assigned

to the correct cluster, for any permutation of clusters. On average, accuracy increased in each EM iteration, as shown in this graph:



6.2 Naive Bayes

```
nfeatures = 10
FeaturesType = Vector(nfeatures, Bool)

get_class = rand_function(Unit, Bool)

class_1_features = rand_function(Unit, FeaturesType)
class_2_features = rand_function(Unit, FeaturesType)

def sample():
    which_class = get_class()
    if which_class:
        features = class_1_features()
    else:
        features = class_2_features()
    return (which_class, features)
```

The naive Bayes model is similar to the clustering model. We have two classes and a feature distribution for each. Since each feature is boolean, we will learn a different categorical distribution for each class.

(figure should show average classification accuracy)

6.3 Factor analysis

```
num_components = 3
point_dim = 5

ComponentsType = Vector(num_components, Double)
PointType = Vector(point_dim, Double)

get_point = rand_function(ComponentsType, PointType)

def sample():
    components = [normal(0, 1)
                  for i in range(num_components)]
    return (components, get_point(components))
```

The factor analysis model is very similar to the clustering model. The main difference is that we replace the categorical `ClusterType` type with a vector type. This results in the model attempting to find each point as an affine function of a vector of standard normal values.

216 (figure should show accuracy by EM iteration. Accuracy can be measured by distances between
217 different factor coefficients, when ordered according to some permutation)
218

219 **6.4 Hidden Markov model**

```
220
221 num_states = 5
222 num_observations = 3
223 sequence_length = 10
224
225 StateType = Categorical(num_states)
226 ObsType = Categorical(num_observations)
227
228 trans_fun = rand_function(StateType, StateType)
229 obs_fun = rand_function(StateType, ObsType)
230
231 def sample():
232     states = [uniform(num_states)]
233     for i in range(sequence_length - 1):
234         states.append(trans_fun(states[-1]))
235     return (states, [obs_fun(s) for s in states])
```

234 In this example, we use the unknown function `trans_fun` for state transitions and `obs_fun` for
235 observations. This means that we will learn both the state transitions and the observation distribution.
236

237 **6.5 Latent Dirichlet allocation**

```
238
239 n_classes = 10
240 n_words = 100
241 n_words_per_document = 1000
242
243 ClassType = Categorical(n_classes)
244 WordType = Categorical(n_words)
245
246 class_to_word = rand_function(ClassType, WordType)
247
248 def sample():
249     which_class = uniform_categorical(n_classes)
250     words = [class_to_word(which_class)
251              for i in range(n_words_per_document)]
252     return (which_class, words)
```

252 In this example, we use the unknown function `class_to_word` to map classes to word distribu-
253 tions. Note that each column of the matrix of parameters for `class_to_word` will represent a
254 categorical distribution over words, and there will be one column for each class.

255 (figure should show accuracy over time. Accuracy can be measured as distance between the learned
256 categorical distributions, for some permutation of classes)
257

258 **6.6 Neural network**

```
259
260 input_dim = 100
261 hidden_dim = 20
262 output_dim = 2
263
264 InputType = Categorical(input_dim)
265 HiddenType = Categorical(hidden_dim)
266 OutputType = Categorical(output_dim)
267
268 input_to_hidden = rand_function(InputType, HiddenType)
269 hidden_to_output = rand_function(HiddenType, OutputType)
270
271 def sample(input_layer):
272     hidden_layer = input_to_hidden(input_layer)
```

```
270     output_layer = hidden_to_output(hidden_layer)
271     return ((), output_layer)
272
```

273 6.7 A more complex model

274
275 TODO: if there is time, we should put a more complex data science type example, where we add
276 dependencies and show change in accuracy as we add/remove assumptions.

277 7 Discussion

278
279 We have found that it is possible to write many useful machine learning models as Quipp programs
280 and then use generic algorithms for inference. Furthermore, performance is $\Theta(1)$. This should make
281 it much easier for non-experts to write useful machine learning models.

282
283 In the future, it will be useful to expand the set of types supported. It is possible to define reasonable
284 default distributions for non-recursive algebraic data types, and it may also be possible to define
285 them for recursive algebraic data types using catamorphisms. Also, it will be useful to create a more
286 usable interface to infer parameters and perform additional data processing given these parameters.
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323