# 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:

```
dim = 2
nclusters = 3
PointType = Vector(dim, Double)
ClusterType = Categorical(nclusters)

get_point = rand_function(ClusterType, PointType)

def sample():
  cluster = uniform_categorical(nclusters)
  return (cluster, get_point(cluster))
```

Notice that we declared two types (`PointType` and `ClusterType`) and one random function `get_point`). 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 `get_point` 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 result is distributed from some exponential family whose natural parameters are determined from the arguments:

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

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

To determine $\eta$ as a function of $x$, 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 entries not in the first row corresponding to sufficient statistics of $y$ that are not features. This ensures that only the natural parameters for features 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 EXAMPLES

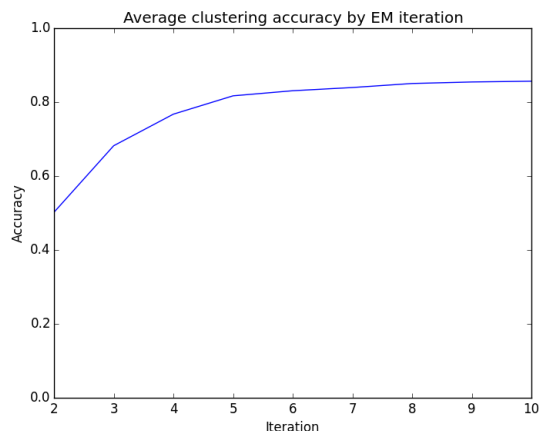### 5.1 CLUSTERING

```
dim = 2
nclusters = 3

PointType = Vector(dim, Double)
ClusterType = Categorical(nclusters)

get_cluster = rand_function(ClusterType)
get_point = rand_function(ClusterType, PointType)

def sample():
  cluster = get_cluster()
  return (cluster, get_point(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 it this graph:



3

## 5.2 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)

## 5.3 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 (comenents, 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.

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

## 5.4 HIDDEN MARKOV MODEL

```
num_states = 5
num_observations = 3
sequence_length = 10

StateType = Categorical(num_states)
ObsType = Categorical(num_observations)

trans_fun = rand_function(StateType, StateType)
obs_fun = rand_function(StateType, ObsType)

def sample():
  states = [uniform(num_states)]
  for i in range(sequence_length - 1):
    states.append(trans_fun(states[-1]))
  return (states, [obs_fun(s) for s in states])
```

4

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

## 5.5  LATENT DIRICHLET ALLOCATION

```
n_classes = 10
n_words = 100
n_words_per_document = 1000

ClassType = Categorical(n_classes)
WordType = Categorical(n_words)

class_to_word = rand_function(ClassType, WordType)

def sample():
    which_class = uniform_categorical(n_classes)
    words = [class_to_word(which_class)
             for i in range(n_words_per_document)]
    return (which_class, words)
```

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

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

## 5.6  NEURAL NETWORK

```
input_dim = 100
hidden_dim = 20
output_dim = 2

InputType = Categorical(input_dim)
HiddenType = Categorical(hidden_dim)
OutputType = Categorical(output_dim)

input_to_hidden = rand_function(InputType, HiddenType)
hidden_to_output = rand_function(HiddenType, OutputType)

def sample(input_layer):
  hidden_layer = input_to_hidden(input_layer)
  output_layer = hidden_to_output(hidden_layer)
  return ((), output_layer)
```

## 5.7  A MORE COMPLEX MODEL

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

## 6  DISCUSSION

We have found that it is possible to write many useful machine learning models as Quipp programs and then use generic algorithms for inference. Furthermore, performance is ¡???¿. This should make it much easier for non-experts to write useful machine learning models.

In the future, it will be useful to support unbounded recursion in models (as in Church). Furthermore, it will also be useful to expand the set of types supported. It is possible to define reasonable default distributions for non-recursive algebraic data types, and it may also be possible to define them for recursive algebraic data types using catamorphisms. Also, it will be useful to create a more usable interface to infer parameters and perform additional data pracessing given these parameters.