

---

# 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;

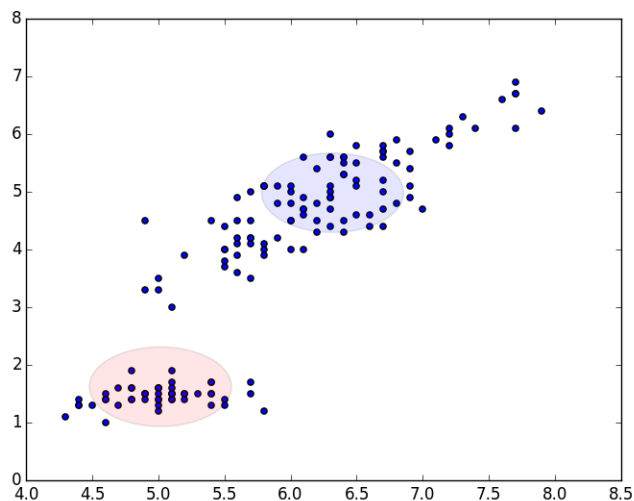
var n = 100;

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

```

It is written as a generative model producing the observations. 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 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. The `observe` function allows us to observe data; here it says that the observation named `point3` is equal to the result of calling `getPoint(cluster)`, where `cluster` is cluster number 3.

To demonstrate, let us run this example on a dataset consisting of 150 points (TODO cite). When we run the program on this data, we infer the parameters to the random function `getPoint`. In this case, `getPoint` is a linear function with Gaussian noise, so it will naturally split the data into 2 clusters with equal variance:



The first cluster is at (6.3, 5.0) and the second is at (5.0, 1.6). They both have a standard deviation of 0.54 in the x direction and 0.69 in the y direction. We could use these parameters to fill in the generative model:

```

return repeat(n, function(i) {

```

```

108     var cluster = randomInteger(nclusters);
109     return [gaussian(cluster == 0 ? 6.3 : 5.0, 0.54),
110             gaussian(cluster == 0 ? 5.0 : 1.6, 0.69)];
111   });
112
113
114

```

This model is estimated to assign probability density  $e^{-393}$  to the data, or 0.0728 per point.

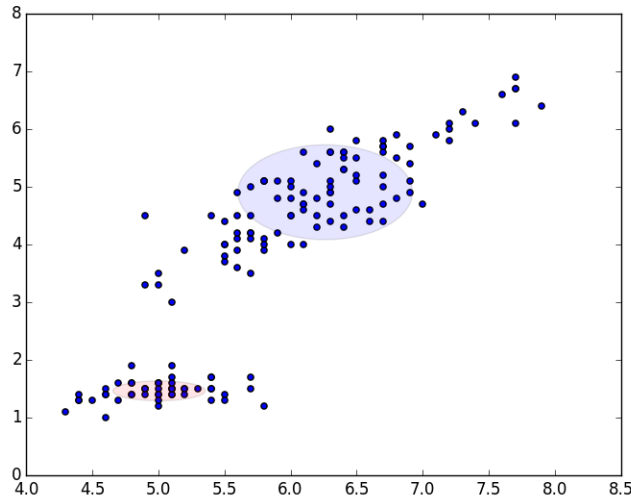
Note that, because the two clusters are forced to have the same variance, they do not fit the data well, since the data has a different shape in each location. To fix this problem, we can substitute the following model:

```

118
119 testParamInference(function() {
120   var pointType = Vector(dim, Double);
121   var getPointFunctions = repeat(nclusters, function(i) {
122     return randFunction(pointType);
123   });
124   return function() {
125     return repeat(n, function(i) {
126       var cluster = randomInteger(nclusters);
127       return observe('point' + i, getPointFunctions[cluster]);
128     });
129   });
130
131
132

```

Using this model, we get the following clusters:



This model (with the parameters filled in) is estimated to assign probability density  $e^{-337}$  to the data, or 0.1058 per point, which means it fits the data better.

### 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  $\theta$
- Generate datasets  $x_{train}, x_{test}$  using  $\theta$
- 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. Between observations,

## 6 Examples

### 6.1 Clustering

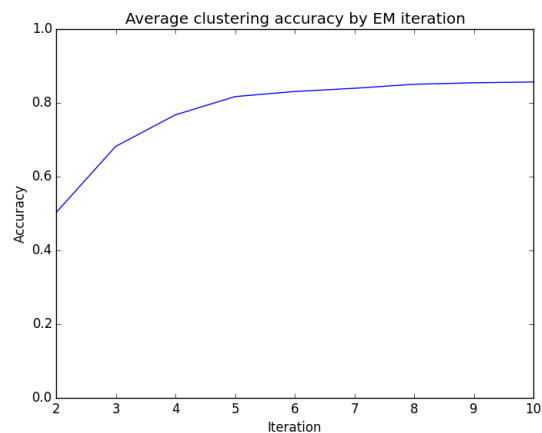
```
var nclusters = 3;
var dim = 2;

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

```
var nfeatures = 10;
var featuresType = Vector(nfeatures, Bool);
var nclasses = 2;
var n = 50;

testParamInference(function(randFunction) {
  var class0Features = randFunction(featuresType);
  var class1Features = randFunction(featuresType);

  var getFeatures = randFunction(Categorical(nclusters), pointType);
  return function() {
    repeat(n, function(i) {
      var whichClass = randomInteger(2);
      if (whichClass == 0) {
```

```

270         observe('features' + i, class0Features);
271     } else {
272         observe('features' + i, class1Features);
273     }
274     });
275 };
276 });

```

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

```

285 var nfactors = 2;
286 var dim = 4;
287
288 var n = 50;
289
290 testParamInference(function(randFunction) {
291     var factorType = Vector(nfactors, Double);
292     var pointType = Vector(dim, Double);
293     var getPoint = randFunction(factorType, pointType);
294     var getFactors = function() {
295         return repeat(nfactors, function() { return gaussian(0, 1); });
296     };
297     return function() {
298         return repeat(n, function(i) {
299             var factors = getFactors();
300             return observe('point' + i, getPoint, factors);
301         });
302     });
303 });

```

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.

### 6.4 Hidden Markov model

```

308 var nstates = 2;
309 var nobs = 3;
310
311 var chainLength = 8;
312
313 var n = 30;
314
315 testParamInference(function(randFunction) {
316     var stateType = Categorical(nstates);
317     var obsType = Categorical(nobs);
318     var transFun = randFunction(stateType, stateType);
319     var obsFun = randFunction(stateType, obsType);
320
321     var observeStates = function(sampIndex, startState, i) {
322         if (i == chainLength) {
323             return [];
324         } else {
325             observe('obs' + sampIndex + '_' + i, obsFun, startState);
326             return [startState].concat(observeStates(sampIndex, transFun(startState), i+1));
327         }
328     };
329 }

```

```

324     }
325   };
326   return function() {
327     return repeat(n, function(sampIndex) {
328       return observeStates(sampIndex, randomInteger(nstates), 0);
329     });
330   });
331 }

```

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

## 6.5 Latent Dirichlet allocation

```

336 var nClasses = 2;
337 var nWords = 100;
338 var maxWordsPerDocument = 1000;
339 var nDocuments = 20;
340
341 var classType = Categorical(nClasses);
342 var wordType = Categorical(nWords);
343
344 testParamInference(function(randFunction) {
345   var classToWord = randFunction(ClassType, WordType);
346   return function() {
347     return repeat(nDocuments, function(docIndex) {
348       var whichClass = randomInteger(nClasses);
349       var nWords = observe('len' + docIndex, randomInteger, maxWordsPerDocument);
350       repeat(nWordsPerDocument, function(wordIndex) {
351         observe('word' + doc + '_' + wordIndex, classToWord, whichClass);
352       });
353       return whichClass;
354     });
355   });
356 }

```

In this example, we use the unknown function `classToWord` to map classes to word distributions. Note that each column of the matrix of parameters for `classToWord` 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)

## 6.6 Neural network

```

366 var inputDim = 100;
367 var hiddenDim = 20;
368 var outputDim = 2;
369
370 var inputType = Categorical(inputDim);
371 var hiddenType = Categorical(hiddenDim);
372 var outputType = Categorical(outputDim);
373
374 var inputs = [...];
375
376 testParamInference(function(randFunction) {
377   var inputToHidden = randFunction(inputType, hiddenType);
378   var hiddenToOutput = randFunction(hiddenType, outputType);
379
380   return function() {
381     return repeat(inputs.length, function(sampIndex) {
382       var hiddenLayer = inputToHidden(inputs[sampIndex]);
383       observe('output' + sampIndex, hiddenToOutput, hiddenLayer);
384     });
385   });
386 }

```

```
378         });
379     };
380 });
```

381

## 382 6.7 A more complex model

383

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

386

## 387 7 Discussion

388

389 We have found that it is possible to write many useful machine learning models as Quipp programs  
390 and then use generic algorithms for inference. Furthermore, performance is  $\approx 10^6$ . This should make  
391 it much easier for non-experts to write useful machine learning models.

392

393 In the future, it will be useful to expand the set of types supported. It is possible to define reasonable  
394 default distributions for non-recursive algebraic data types, and it may also be possible to define  
395 them for recursive algebraic data types using catamorphisms. Also, it will be useful to create a more  
396 usable interface to infer parameters and perform additional data processing given these parameters.

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431