

---

# OpenPPL Documentation

*Release 0.1.0*

**All contributors**

October 25, 2015



<b>1</b>	<b>Basics</b>	<b>3</b>
1.1	Terminologies . . . . .	3
1.2	Getting Started: Gaussian Mixture Model . . . . .	4
1.3	Generic Specification: Finite Mixture Model . . . . .	5
1.4	Queries . . . . .	5
<b>2</b>	<b>More Examples</b>	<b>9</b>
2.1	Latent Dirichlet Allocation . . . . .	9
2.2	Hidden Markov Model . . . . .	10
2.3	Markov Random Fields . . . . .	12
2.4	Conditional Random Fields . . . . .	14
2.5	Deep Boltzmann Machines . . . . .	15
<b>3</b>	<b>Nonparametric Models</b>	<b>19</b>
3.1	Dirichlet Process Mixture Model . . . . .	19
3.2	Hierarchical Dirichlet Processes . . . . .	20
3.3	Gaussian Processes . . . . .	21
<b>4</b>	<b>The Inference Framework</b>	<b>23</b>
4.1	Gaussian Mixture Model Revisited . . . . .	23
4.2	Inference via Message Passing . . . . .	25



This document proposes the design of *OpenPPL*, an open source probabilistic programming framework implemented as a domain-specific language on top of Julia. This document discusses the design of basic language constructs for model specification, query, as well as directives to control algorithmic choices in inference.

Contents:



*OpenPPL* is a domain-specific language built on top of Julia using macros. The language consists of two parts: *model specification* and *query*. In particular, a *model specification* formalizes a probabilistic model, which involves declaring variables and specifying relations between them; while a *query* specifies *what is given* and *what is to be inferred*.

## 1.1 Terminologies

Here is a list of terminologies that would be involved in the description.

**Variable** A variable generally refers to an entity that can take a value of certain type. It can be a random variable directly associated with a distribution, a deterministic transformation of another variable, or just some value given by the user. The value of a variable can be given or unknown.

**Constant** A value that is fixed when a query is constructed and fixed throughout the inference procedure. A constant is typically used to represent vector dimensions, model sizes, and hyper-parameters etc. Note that model parameters are typically considered as variables instead of constants. For example, a learning process can be formulated as a query that solves the parameter given a set of observations, where the parameter values can be iteratively updated.

**Domain** The domain of a variable refers to the set of possible values that it may take. Any Julia type (e.g. `Int`, `Float64`) can be considered as a domain that contains any value of that type. This package also supports array domains and restrictive domains that contain only a subset of a specific type. Here are some examples:

```

1  1..K           # integers between 1 and K
2  0.0..Inf      # all non-negative real values
3  Float64^n     # n-dimensional real vectors
4  Float64^(m,n) # matrices of size (m, n)
5  (0.0..1.0)^m  # m-dimensional real vectors with all components in [0., 1.]

```

**Distribution** From a programmatic standpoint, a distribution can be considered as a *stochastic* function that yields a random value in some domain. A distribution can accept zero or more arguments. A distribution should be *stochastically pure*, meaning that it always outputs the same value given the same arguments and the same state of the random number generator. Such purity makes it possible to reason about program structure and transform the model from one form to another.

**Factor** A factor is a pure real-valued function. Here, “pure” means that it always output the same value given the same inputs. Factors are the core building blocks of a probabilistic model. A complex distribution is typically formulated as a set of variables connected by factors. Even a simple distribution (e.g. normal distribution) consists of a factor that connects between generated variables and parameters (which may be considered as a variable with fixed value).

## 1.2 Getting Started: Gaussian Mixture Model

Here, we take the *Gaussian Mixture Model* as an example to illustrate how we can specify a probabilistic model using this Julia domain-specific language (DSL). A *Gaussian Mixture Model* is a generative model that combines several Gaussian components to approximate complex distributions (e.g. those with multiple modals.). A Gaussian mixture model is characterized by a prior distribution  $\pi$  and a set of Gaussian component parameters  $(\mu_1, \Sigma_1), \dots, (\mu_K, \Sigma_K)$ . The generative process is described as follows:

$$z_i \sim \pi$$

$$x_i | z_i \sim \mathcal{N}(\mu_{z_i}, \Sigma_{z_i})$$

### 1.2.1 Model Specification

The model specification is given by

```

1  @model GaussianMixtureModel begin
2      # constant declaration
3      @constant d::Int # vector dimension
4      @constant n::Int # number of observations
5      @hyperparam K::Int # number of components
6
7      # parameter declaration
8      @param pi :: (0.0..1.0)^K # prior proportions
9      for k in 1 : K
10         @param mu[k] :: Float64^d # component mean
11         @param sig[k] :: Float64^(d, d) # component covariance
12     end
13
14     # sample generation process
15     for i in 1 : n
16         z[i] ~ Categorical(pi)
17         x[i] ~ MultivariateNormal(mu[z[i]], sig[z[i]])
18     end
19 end

```

This model specification should be self-explanatory. However, it is still worth clarifying several aspects:

- The macro `@model` defines a model type named `GaussianMixtureModel`, and creates an environment (delimited by `begin` and `end`) for model formulation. All model types created by `@model` is a sub type of `AbstractModel`.
- The macro `@constant` declares `d`, `K`, and `n` as constants. The values of these constants need not be given in the specification. Instead, they are needed upon query. Particularly, to construct a model, one can write

```
mdl = GaussianMixtureModel()
```

One can *optionally* fix the value of constants through keyword arguments in model construction, as below

```
mdl = GaussianMixtureModel(d = 2, K = 5)
```

Note: fixing constants upon model construction is generally unnecessary. However, it might be useful to fix them under certain circumstances to to simplify queries or restrict its use. Once a constant is fixed, it need not be specified again in the query.

- The macro `@hyperparam` declares hyper parameters. Hyper parameters are similar to constant technically, except that they typically refer to model configurations that may be changed during cross validation.

- Variables can be defined using the syntax as `variable-name :: domain`. A for-loop can be used to declare multiple variables in the same domain. When the variable domain is clear from the context (e.g. the domain of  $z$  and  $x$  can be inferred from where they are drawn), the declaration can be omitted.
- The macro `@param` tags certain variables to be parameters. The information will be used in the learning algorithm to determine which variables are the parameters to estimate.
- The statement `variable-name ~ distribution` introduces a conditional distribution over variables, which will be translated into a factor during model compilation.

## 1.3 Generic Specification: Finite Mixture Model

The Gaussian mixture model can be considered as a special case in a generic family called *Finite mixture model*. Generally, the components of a finite mixture model can be arbitrary distributions. To capture the concept of *generic distribution family*, we introduce *generic specification* (or *parametric specification*), which can take type arguments.

The specification of the generic finite mixture model is given by

```

1  @model FiniteMixtureModel{G, ParamTypes} begin
2      # constant declaration
3      @hyperparam K::Int
4      @constant n::Int
5
6      # parameter declaration
7      @param pi :: (0.0..1.0)^K # prior proportions
8      for k = 1 : K
9          for j = 1 : length(ParamTypes)
10             @param theta[k][j] :: ParamTypes[j]
11         end
12     end
13
14     # sample generation process
15     for i in 1 : n
16         z[i] ~ Categorical(pi)
17         x[i] ~ G(theta[z[i]]...)
18     end
19 end

```

One may consider a generic specification above as a specification template. To obtain a Gaussian mixture model specification, we can use the `@modelalias` macro, as below:

```

1  @modelalias GaussianMixtureModel FiniteMixtureModel{G, ParamTypes} begin
2      @constant d::Int
3      @with G = MultivariateNormal
4      @with ParamTypes[1] = Float64^d # component mean
5      @with ParamTypes[2] = Float64^(d, d) # component covariance
6  end
7
8  mdl = GaussianMixtureModel()

```

The `@modelalias` macro allows introducing new constants and specializing the type parameters.

## 1.4 Queries

In machine learning, the most common queries that people would make include

- learning: estimate model parameters
- prediction: predict the value or marginal distribution over unknown variables, given a learned model and observed variables.
- evaluation: evaluate log-likelihood of observations with a given model
- sampling: draw a set of samples of certain variables

To simplify these common queries, we provide several functions.

### 1.4.1 Query

*Query* refers to the task of inferring the value or marginal distributions of unknown variables, given a set of known variables.

```
1 function query(rmdl::AbstractModel, knowns::Associative, qlist::Array, options)
2   set_variables!(rmdl, knowns)
3   q = compile_query(rmdl, qlist, options) # this returns a query function q
4   return q() # runs the query function and returns the results
5 end
6
7 function query(rmdl::AbstractModel, knowns::Associative, q)
8   infer(rmdl, knowns, q, default_options(rmdl))
9 end
```

`qlist` is a list of variables or functions over variables that you want to infer. The function `compile_query` actually performs model compilation, analyzing model structure, choosing appropriate inference algorithms, and generating a closure `q`, which, when executed, actually performs the inference.

This query function here is very flexible. One can use it for prediction and sampling, etc.

```
1 # let rmdl be a learned model
2
3 # predict the value of z given observation x
4 z_values = query(rmdl, {x=>columns(data)}, :z)
5
6 # infer the posterior marginal distributions over z given x
7 z_marginal = query(rmdl, {x=>columns(data)}, :(marginal(z)))
8
9 # you can simultaneously infer only selected variables in a flexible way
10 r = query(rmdl, {x=>columns(data)}, {(z[1]), (z[2]), (marginal(z[3]))})
11
12 # draw 100 samples of z
13 samples = query(rmdl, {x=>columns(data)}, :(samples(z, 100)))
```

Note that inputs to the function are symbols like `:z` or expressions like `:(marginal(z))`, which indicate *what we want to query*. It is incorrect to pass `z` or `marginal(z)` – the value of `z` or `marginal(z)` is unavailable before the inference.

### 1.4.2 Learning

*Learning* refers to the task of estimating model parameters given observed data. This can be considered as a special kind of query, which infers the values of model parameters, given observed data.

```
1 function learn_model(mdl::AbstractModel, data::Associative, options)
2   rmdl = copy(mdl)
3   set_variables!(rmdl, data)
```

```

4     q = compile_query(rmdl, parameters(rmdl), options)
5     set_variables!(rmdl, q())
6     return rmdl
7 end
8
9 function learn_model mdl::AbstractModel, data
10     learn_model(mdl, data, default_options(mdl))
11 end
12
13 # learn a GMM, a simple wrapper of learn_model
14 # suppose data is a d-by-n matrix
15 rmdl = learn_model(
16     GaussianMixtureModel(K = 3, d = size(data,1), n = size(data,2)),
17     {:x => columns(data)})

```

In the function `learn_model`, `parameters(rmdl)` returns a list of parameters as the query list. Then the statement `q = compile_query(rmdl, parameters(rmdl), options)` returns a query function `q`, such that `q()` executes the estimation procedure and returns the estimated model parameters. The following example shows how we can use this function to learn a Gaussian mixture model.

```

1 function learn_gmm(data::Matrix{Float64}, K::Int)
2     learn_model(
3         GaussianMixtureModel(K = K, d = size(data,1), n = size(data,2)),
4         {:x => columns(data)})
5 end
6
7 rmdl_K3 = learn_gmm(data, 3)
8 rmdl_K5 = learn_gmm(data, 5)

```

Here, `learn_gmm` is a light-weight wrapper of `learn_model`.

### 1.4.3 Evaluation

*Evaluation* refers to the task of evaluating log-pdf of samples with respect to a learned model.

```

# evaluate the logpdf of x with respect to a GMM
lp = query(rmdl, {:x=>columns(data)}, :(logpdf(x)))

```

### 1.4.4 Options

The compilation options that control how the query is compiled can be specified through the `options` argument in the `query` or `learn_model` function. The following is some examples

```

rmdl = learn_model(mdl, data, {:method=>"variational_em", :max_iter=>100, :tolerance=1.0e-6})

```

For sampling, we may use a different set of options

```

options = {
    :method=>"gibbs_sampling", # choose to use Gibbs sampling
    :burnin=>5000,             # the number of burn-in iterations
    :lag=>100                  # the interval between two samples to retain

    samples = query(rmdl, {x:=>columns(data)}, :(samples(z, 100)), options)

```

## 1.4.5 Query Functions with Arguments

It is often desirable in practice that a query function can be applied to different data sets without being re-compiled. For this purpose, we introduce a function `make_query_function`. The following example illustrates its use:

```
1 # suppose rmdl is a learned model
2
3 q = make_query_function(rmdl,
4   (:data,), # indicates that the function q would take one argument data
5   {:x=>:(columns(data))}, # indicates how the argument is set to the model as a known value
6   {:z},      # specifies what to query
7   options)  # compilation options
8
9 # q can be repeatedly use for different datasets (without being re-compiled)
10 z1 = q(x1)
11 z2 = q(x2)
```

Note that `q` is a closure that holds reference to the learned model, so you don't have to pass the model as an argument into `q`. The following code use this mechanism to generate a sampler:

```
1 q = make_query_function(rmdl,
2   (:data, :n), # q would take two arguments, the observed data and the number of samples
3   {:x=>:(columns(data))},
4   {:sample(z, n)},
5   options)
6
7 # draw 100 samples of z
8 zs1 = q(x, 100)
9
10 # draw another 200 samples of z
11 zs2 = q(x, 200)
```

---

## More Examples

---

This chapter shows several common examples to illustrate how the model specification DSL works in practice.

### 2.1 Latent Dirichlet Allocation

Latent Dirichlet Allocation (LDA) is a probabilistic model for topic modeling. The basic idea is to use a set of topics to describe documents. Key aspects of this model is summarized below:

1. Each document is considered as a *bag of words*, which means that only the frequencies of words matter, while the sequential order is ignored. In practice, each document is summarized by a histogram vector  $h$ .
2. The model comprises a set of topics. We denote the number of topics by  $K$ . Each topic is characterized by a distribution over vocabulary, denoted by  $\beta_k$ . It is a common practice to place a Dirichlet prior over these distributions.
3. Each document is associated with a topic proportion vector  $\theta$ , generated from a Dirichlet prior.
4. Each word in a document is generated independently. Specifically, a word is generated as follows

$$z_i \sim \theta$$

$$w_i | z_i \sim \beta_{z_i}$$

Here,  $z_i$  indicates the topic associated with the word  $w_i$ .

The model specification is then given by

```

1  @model LatentDirichletAllocation begin
2      # constants
3      @constant m::Int      # vocabulary size
4      @constant n::Int      # number of documents
5      @constant nw::Vector{Int} # numbers of words in documents
6
7      @hyperparam K::Int      # number of topics
8      @hyperparam a::Float64 # Dirichlet prior for beta
9
10     # parameters
11     @param alpha :: (0.0 .. Inf)^K
12     for k in 1 : K
13         @param beta[k] ~ Dirichlet(a)
14     end
15 
```

```

16   # documents
17   for i in 1 : n
18     theta[i] ~ Dirichlet(alpha)
19     let p = sum(beta[k] * theta[i][k] for k in 1 : K)
20     h[i] ~ Multinomial(nw[i], p)
21   end
22 end
23 end

```

Here, the keyword `let` introduces a local value `p` to denote the sum vector. It is important to note that `p` has a local scope (thus it is not visible outside the loop), and the value of `p` can be different for different `i`.

The following function learns an LDA model from word histograms of training documents.

```

1  function learn_lda(h::Matrix{Double}, K::Int, alpha::Float64)
2    learn_model(LatentDirichletAllocation(
3      m = size(h, 1), n = size(h, 2), nw = sum(h, 1), K = K, a = alpha,
4      {:h => columns(h)})
5  end
6
7  mdl = learn_lda(training_hists, K, alpha)

```

The following statement infers topic proportions of testing documents, given a learned LDA model.

```

1  # suppose mdl is a learned LDA model
2  theta = query(mdl, {:h => columns(testing_corpus)}, :theta)

```

## 2.2 Hidden Markov Model

Hidden Markov Model (HMM) is a popular model to describe dynamic processes. It assumes that the observed process is driven by a latent Markov chain, and the observation at each time step is independently generated conditioned on the latent states at the same time. A time-homogeneous Markov model is characterized by an initial distribution of states, denoted by  $\pi$ , a transition probability matrix, denoted by  $T$ , and an observation model that generates observations based on latent states. Each sample of a Hidden Markov model is a sequence  $x = (x_0, \dots, x_n)$ , where the observation  $x_t$  is associated with a latent state  $s_t$ . The joint distribution over both the observations and the states is given by

$$p(x, s) = \pi(x_0) \prod_{t=1}^n T(x_{t-1}, x_t) \prod_{t=0}^n p(x_t | s_t; \theta).$$

Here,  $\theta$  denotes the parameter of the observation model.

Generally, the component models associated with the observations can be any distributions. Therefore, HMM is actually a family of distributions that can be specified using a generic specification, as below:

```

1  @model HiddenMarkovModel{G, ParamTypes} begin
2    @constant n::Int # number of sequences
3    @constant len::Vector{Int} # the sequence lengths
4    @hyperparam K::Int # the size of latent state space
5
6    # parameters
7    @param pi :: (0.0 .. 1.0)^K # initial distribution
8    @param T :: (0.0 .. 1.0)^(K, K) # transition probability matrix
9
10   for k in 1 : K

```

```

11     for j in 1 : length(ParamTypes)
12         @param theta[k][j] :: ParamTypes[j]
13     end
14 end
15
16 # sequences
17 for i in 1 : n
18     z[i][1] ~ Categorical(pi)
19     for t = 2 : len[i]
20         z[i][t] ~ Categorical(T[z[t-1], :])
21     end
22
23     for t in 1 : len[i]
24         let s = z[i][t]
25             x[i][t] ~ G(theta[s]...)
26         end
27     end
28 end
29 end

```

To construct an HMM with K Gaussian components, one can write:

```

1 @modelalias HiddenMarkovGaussModel HiddenMarkovModel{G, Params} begin
2     @constant d::Int # vector space dimension
3     @with G = MultivariateNormal
4     @with Params[1] = Float64^d
5     @with Params[2] = Float64^(d, d)
6 end

```

The following query function learns a HMM (with Gaussian components):

```

1 function learn_hmm(seqs, K::Int)
2     # seqs is a collection of observed sequences
3     # K is the number of latent states
4
5     learn_model(HiddenMarkovGaussModel(
6         n = length(seqs), len = map(length, seqs), K = K),
7         {:x => seqs})
8 end
9
10 mdl = learn_hmm(seqs, K)

```

The following query draws samples of the latent state sequences, given a learned HMM model and a sequences of observations.

```

1 function hmm_sample(mdl::HiddenMarkovGaussModel, obs::Matrix{Float64}, ns::Int)
2     # obs is a sequence of observed features (each column for a time step)
3     # ns is the number of samples to d
4
5     query(mdl, {:x => obs}, :(sample(z, ns)))
6 end
7
8 # run the function to draw 100 sample state-sequences for x
9 x = rand(3, 100)
10 y = hmm_sample(mdl, x, 100)

```

## 2.3 Markov Random Fields

Unlike Bayesian networks, which can be factorized into a product of (conditional) distributions, Markov random fields are typically formulated in terms of potentials. Generally, a MRF formulation consists of two parts: identifying relevant cliques (small subsets of directly related variables) and assigning potential functions to them. In computer vision, Markov random fields are widely used in low level vision tasks, such as image recovery and segmentation. Deep Boltzmann machines, which become increasingly popular in recent years, are actually a special form of Markov random field. Here, we use a simple MRF model in the context of image denoising to demonstrate how one can use the model specification to describe an MRF.

From a probabilistic modeling standpoint, the task of image denoising can be considered as an inference problem based on an image model combined with an observation model. An image model captures the prior knowledge as to what an *clean* image may look like, while the observation model describes how the observed image is generated through a noisy imaging process. Here, we consider a simple setting: Gaussian MRF prior + white noise. A classical formulation of Gaussian MRF for image modeling is given below

$$p(x) = \frac{1}{Z} \exp(-E(x; \theta)).$$

Here, the distribution is formulated in the form of a *Gibbs distribution*, and  $E(x; \theta)$  is the energy function, which is controlled by a parameter  $\theta$ . The energy function  $E(x; \theta)$  can be devised in different ways. A typical design would encourage smoothness, that is, assign low energy value when the intensity values of neighboring pixels are close to each other. For example, a classical formulation uses the following energy function

$$E(x; \theta) = \theta \sum_{\{u,v\} \in \mathcal{C}} (x(u) - x(v))^2$$

Here,  $u$  and  $v$  are indices of pixels, and the clique set  $\mathcal{C}$  contains all edges between neighboring pixels. With the white noise assumption, the observed pixel values are given by

$$y(u) = x(u) + \varepsilon(u), \quad \text{with } \varepsilon(u) \sim \mathcal{N}(0, \sigma^2).$$

Below is the specification of the joint model:

```

1  @model SimpleMRF begin
2      @constant nimgs::Int # the number of images
3      @constant imgsizes::Vector{Int, Int}
4
5      # parameters
6      @param theta::Float64 # the Gaussian MRF parameter
7      @param sig::Float64 # the variance of white noise
8
9      for t in 1 : nimgs
10         let m = imgsizes[t][1], n = imgsizes[t][2]
11             x[t]::Float64^(m, n) # the true image
12             y[t]::Float64^(m, n) # the observed noisy image
13
14             let xt = x[t], yt = y[t]
15                 # the image prior (Gaussian MRF)
16                 for i in 2 : m-1, j in 2 : n-1
17                     @fac exp(-theta * (xt[i, j] - xt[i, j-1])^2)

```

```

18         @fac exp(-theta * (xt[i,j] - xt[i,j+1])^2)
19         @fac exp(-theta * (xt[i,j] - xt[i-1,j])^2)
20         @fac exp(-theta * (xt[i,j] - xt[i+1,j])^2)
21     end
22
23     # the observation model
24     for i in 1 : m, j in 1 : n
25         yt[i,j] ~ Normal(xt[i,j], sig)
26     end
27 end
28 end
29 end
30 end

```

The following statement learns the model from a set of uncorrupted images

```

# suppose imgs is an array of images
mdl = learn_model(SimpleMRF(nimgs=length(imgs), imgsizes=map(size, imgs)), {x=>imgs})

```

In this specification, four potentials are used to connect a pixel to its left, right, upper, and lower neighbors. This approach would become quite cumbersome as the neighborhood grows. Many state-of-the-art denoising algorithms use much larger neighborhood (e.g.  $5 \times 5$ ,  $9 \times 9$ , etc) to capture high order texture structure. A representative example is the *Field of Experts*, where the MRF prior is defined using a set of filters as follows:

$$p(x) = \frac{1}{Z} \exp \left( \sum_{k=1}^K \sum_{c \in \mathcal{C}} \rho(J_k^T x_c, \alpha_k) \right), \quad \text{with } \rho(v, \alpha) := -\alpha \log(1 + v^2).$$

Here,  $\mathcal{C}$  is the set of all patches of certain size (say  $5$  times  $5$ ), and  $x_c$  is the pixel values over a small patch  $c$ . Here,  $K$  filters  $J_1, \dots, J_K$  are used, and  $J_k^T x_c$  is the filter response at patch  $c$ .  $\rho$  is a robust potential function that maps the filter responses to potential values, controlled by a parameter  $\alpha$ . The specification below describes this more sophisticated model, where local functions and local variables are used to simplify the specification.

```

1  @model FieldOfExperts begin
2      @constant K::Int          # the number of filters
3      @constant w::Int         # patch size (w = 5 for 5 x 5 patches)
4      @constant ew::Int = (w - 1) / 2 # half patch dimension
5      @constant nimgs::Int     # the number of images
6      @constant imgsizes::Vector{Int, Int}
7      @constant sig :: Float64 # variance of white noise
8
9      # parameters
10     for k = 1 : K
11         @param J[k] :: Float64^(w, w) # filter kernel
12         @param alpha[k] :: Float64    # filter coefficient
13     end
14
15     # the robust potential function
16     rho(v, a) = -a * log(1 + v * v)
17
18     for t in 1 : nimgs
19         let m = imgsizes[t][1], n = imgsizes[t][2]
20             x[t]::Float64^(m, n) # the true image
21             y[t]::Float64^(m, n) # the observed noisy image
22
23             let xt = x[t], yt = y[t]
24                 # the image prior

```

```

25         for k in 1 : K, i in 1+ew : m-ew, j in 1+ew : n-ew
26             let c = vec(xt(i-ew:i+ew, j-ew:j+ew))
27                 @fac exp(rho(dot(J[k], c), alpha[k]))
28             end
29         end
30
31         # the observation model
32         for i in 1 : m, j in 1 : n
33             yt[i,j] ~ Normal(xt[i,j], sig)
34         end
35     end
36 end
37 end
38 end

```

Below is a query function that learns a field-of-experts model.

```

1  function learn_foe(imgs, w::Int, K::Int)
2      # imgs: an array of images
3      # w: the patch dimension
4      # K: the number of filters
5
6      mdl = FieldOfExperts(
7          K = K, w = w, nimgs = length(imgs),
8          imgsize = map(size, imgs))
9
10     learn_model(mdl, {x=>imgs})
11 end

```

Given a learned model, the following query function performs image denoising.

```

1  function foe_denoise(mdl::FieldOfExperts, sig::Float64, noisy_im::Matrix{Float64})
2      # sig: the noise variance (which is typically given in denoising tasks)
3      # noisy_im: the observed noisy image
4
5      query(mdl, {sig=>sig, :y=>[noisy_im]}, :x)
6  end
7
8  denoised_im = foe_denoise(mdl, 0.1, noisy_im)

```

## 2.4 Conditional Random Fields

Structured prediction, which exploits the statistical dependencies between multiple entities within an instance, has become an important area in machine learning and related fields. Conditional random field is a popular model in this area. Here, I consider a simple application of CRF in computer vision. A visual scene usually comprises multiple objects, and there exist statistical dependencies between the scene category and the objects therein. For example, a bed is more likely in the bedroom than in a forest. A conditional random field that takes advantage of such relations can be formulated as follows

$$p(s, o|x, y) = \frac{1}{Z(\alpha, \beta, \theta)} \exp \left( \psi_s(s, x; \alpha) + \sum_{i=1}^n \psi_o(o_i, y_i; \beta) + \sum_{i=1}^n \varphi(s, o_i; \theta) \right)$$

This formulation contains three potentials:

- $\psi_s(s, x; \alpha) := \alpha_s^T x$  connects the scene class  $s$  to the observed scene feature  $x$ ,
- $\psi_o(o_i, y_i; \beta) := \beta_o^T y_i$  connects the object label  $o_i$  to the corresponding object feature  $y_i$ ,
- $\varphi(s, o_i; \theta) := \theta(s, o_i)$  captures the statistical dependencies between scene classes and object classes.

In addition,  $Z$  is the normalization constant, whose value depends on the parameters  $\alpha, \beta$ , and  $\theta$ . Below is the model specification:

```

1  @model SceneObjectCRF begin
2      @constant M::Int      # the number of scene classes
3      @constant N::Int      # the number of object classes
4      @constant p::Int      # the scene feature dimension
5      @constant q::Int      # the object feature dimension
6      @constant nscenes    # the number of scenes
7      @constant nobjs::Vector{Int} # numbers of objects in each scene
8
9      for k in 1 : M
10         @param alpha[k] :: Float64^p
11     end
12     for k in 1 : N
13         @param beta[k] :: Float64^q
14     end
15     @param theta :: Float64^(p, q)
16
17     for i in 1 : nscenes
18         let n = nobjs[i]
19             s[i] :: 1 .. M      # the scene class label
20             o[i] :: (1 .. N)^n # the object class labels
21             x[i] :: Float64^p   # the scene feature vector
22
23             for j in 1 : n
24                 y[i][j] :: Float64^q # the object features
25             end
26
27             @fac dot(alpha[s[i]], x)
28             for j in 1 : n
29                 let k = o[i][j]
30                     @expfac dot(beta[k], y[i][j])
31                     @expfac theta[s[i], k]
32                 end
33             end
34         end
35     end
36 end

```

Note here that `@expfac f(x)` is equivalent to `@fac exp(f(x))`. The introduction of `@expfac` is to simplify the syntax in cases where factors are specified in log-scale.

## 2.5 Deep Boltzmann Machines

A *Boltzmann machine (BM)* is a generative probabilistic model that describes data through hidden layers. In particular, a *deep belief network* and a *deep Boltzmann machine*, which becomes increasingly popular in machine learning and its application domains, can be constructed by stacking multiple layers of BMs. In a generic Boltzmann machine, the joint distributions over both hidden units  $\mathbf{h}$  and visible units  $\mathbf{v}$  are given by

$$p(\mathbf{v}, \mathbf{h}; \theta) = \frac{1}{Z(\theta)} \exp\left(\frac{1}{2}\mathbf{v}^T \mathbf{L} \mathbf{v} + \frac{1}{2}\mathbf{h}^T \mathbf{J} \mathbf{h} + \mathbf{v}^T \mathbf{W} \mathbf{h}\right)$$

When  $\mathbf{L}$  and  $\mathbf{J}$  are zero matrices, this reduces to a *restricted Boltzmann machine*. By stacking multiple layers of BMs, we obtain a *deep Boltzmann machine* as follows

$$p(\mathbf{v}, \mathbf{h}; \theta) = \frac{1}{Z} \exp\left(\frac{1}{2}\mathbf{v}^T \mathbf{L} \mathbf{v} + \mathbf{v}^T \mathbf{W}_0 \mathbf{h}_1 + \sum_{l=1}^L \mathbf{h}_l^T \mathbf{J}_l \mathbf{h}_l + \sum_{l=1}^{L-1} \mathbf{h}_l^T \mathbf{W}_l \mathbf{h}_{l+1}\right)$$

This probabilistic network, despite its complex internal structure, can be readily specified using the DSL as below

```

1  @model DeepBoltzmannMachine begin
2      @hyperparam L::Int # the number of latent layers
3      @hyperparam nnodes::Vector{Int} # the number of nodes in each layer
4      @constant d::Int # the dimension of observed sample
5      @constant n::Int # the number of observed samples
6
7      # declare coefficient matrices
8      @param L :: Float64^(d, d)
9      @param W0 :: Float64^(d, nnodes[1])
10
11     for k = 1 : L
12         @param J[k] :: Float64^(nnodes[k], nnodes[k])
13     end
14
15     for k = 1 : L-1
16         @param W[k] :: Float64^(nnodes[k], nnodes[k+1])
17     end
18
19     # declare of variables
20     for i = 1 : n
21         obs[i] :: Float64^d
22         for k = 1 : L
23             latent[i][k] :: Float64^(nnodes[k])
24         end
25     end
26
27     # samples
28     for i = 1 : n
29         let h = samples[i], v = obs[i]
30             # intra-layer connections
31             @expfac v' * L * v
32
33             for k = 1 : L
34                 @expfac h[k]' * J[k] * h[k]
35             end
36
37             # inter-layer connections
38             @expfac v' * W0 * h[1]
39
40             for k = 1 : L-1
41                 @expfac h[k]' * W[k] * h[k+1]
42             end
43         end

```

```
44     end
45 end
```

To learn this model from a set of samples, one can write

```
1  function learn_deepbm(x::Matrix{Float64}, nnodes::Vector{Int})
2      # each column of x is a sample
3      # nnodes specifies the number of nodes at each layer
4
5      learn_model(DeepBoltzmannMachine(L = length(nnodes), nnodes=nnodes,
6          d=size(x,1), n=size(x,2)), {:obs=>columns(x)})
7  end
8
9  mdl = learn_deepbm(x, [100, 50, 20])
```



---

## Nonparametric Models

---

Bayesian nonparametric models, such as DP mixture models, provide a flexible approach in which model structure (e.g. the number of components) can be adapted to data. The capability and effectiveness of such models have been proven in many applications.

This flexibility, however, leads to new questions to probabilistic programming – how to express models whose size/structure can vary. The Julia macro system makes it possible to address this problem in an elegant way due to its lazy evaluation nature.

### 3.1 Dirichlet Process Mixture Model

Dirichlet process mixture model (DPMM) is one of the most widely used in Bayesian nonparametrics. The formulation of DPMM is given by

$$\begin{aligned} D &\sim DP(\alpha B) \\ \theta_i &\sim D, \quad x_i \sim G(\theta_i), \quad \forall i = 1, 2, \dots, n \end{aligned}$$

Here,  $\alpha$  is the concentration parameter,  $B$  is the base measure, and  $G$  denotes the component models that generate the observed samples. The model specification for a DPMM is given as below. It has been shown that  $D$  is almost surely discrete, and can be expressed in the following form:

$$D = \sum_{k=1}^{\infty} \pi_k \delta_{\phi_k}$$

Hence, there exists positive probability that some of the components will be repeatedly sampled, and thus the number of distinct components is usually smaller than the number of samples. In practice, it is useful to construct a pool of components  $\phi_1, \phi_2, \dots$ , and introduce an indicator  $z_i$  for each sample  $x_i$ , such that  $\theta_i = \phi_{z_i}$ . To make this explicit, we can reformulate the model as below

$$\begin{aligned} \pi &\sim \text{StickBreak}(\alpha) \\ \phi_k &\sim B, \quad \forall k = 1, 2, \dots \\ z_i &\sim \pi, \quad x_i \sim G(\phi_{z_i}), \quad \forall i = 1, 2, \dots, n \end{aligned}$$

Here,  $\pi$ , a sample from the stick breaking process, is an infinite sequence that sums to unity (i.e.  $\sum_{i=1}^{\infty} \pi_i = 1$ ). The values of  $\pi_i$  are defined as

$$v_k \sim \text{Beta}(1, \alpha), \quad \forall k = 1, 2, \dots$$

$$\pi_1 = v_1, \quad \pi_k = v_k \prod_{l=1}^{k-1} (1 - v_l), \quad \forall k = 1, 2, \dots$$

The model specification for this formulation is given below

```

@model DPMixtureModel{B, G} begin
  @constant n::Int # the number of observed samples
  @hyperparam alpha::Float64 # the concentration parameter

  pi ~ StickBreak(alpha)
  for k = 1 : Inf
    phi[k] ~ B
  end

  # samples
  for i = 1 : n
    z[i] ~ pi
    x[i] ~ G(phi[z[i]])
  end
end
end

```

Remarks:

- The parametric setting makes it possible to use arbitrary base distribution  $B$  and component  $G$  here, with the same generic formulation.
- In theory,  $\pi$  is an infinite sequence, and therefore it is not feasible to completely instantiate a sample path of  $\pi$  in computer. This variable may be marginalized out in the inference, and thus directly querying  $\pi$  is not allowed.
- $\phi$  has infinitely many components. However, only a finite number of them are needed in inference. The compiler should generate a lazy data structure that only memorizes the subset of components needed during the inference. In particular, `phi[k]` is constructed and memorized when there is an `i` such that `k = z[i]`. Some efforts (e.g. the *LazySequences.jl*) have demonstrated that lazy data structure can be implemented efficiently in Julia.

## 3.2 Hierarchical Dirichlet Processes

HDP is an extension of the DP mixture models, which allows groups of data to be modeled by different DPs that share components. The formulation of HDP is given below

$$D_0 \sim DP(\alpha B)$$

$$D_k \sim DP(\gamma_k D_0)$$

$$\theta_{ki} \sim D_k, \quad x_{ki} \sim G(\theta_{ki}), \quad \forall k = 1, \dots, m, \quad i = 1, \dots, n_k$$

Using  $D_0$  as a base measure for the DP associated with each group, all groups share components in  $D_0$  while allowing potentially infinite number of components. This formulation can be re-written (equivalently) using Pitman-Yor process, as follows

$$\pi_0 \sim \text{StickBreak}(\alpha)$$

$$\psi_j \sim B, \quad \forall j = 1, 2, \dots$$

Then for each  $k = 1, \dots, m$ ,

$$\begin{aligned}\pi_k &\sim \text{StickBreak}(\gamma_k) \\ u_{kt} &\sim \pi_0, \quad \phi_{kt} = \psi_{u_{kt}}, \quad t = 1, 2, \dots \\ z_{ki} &\sim \pi_k, \quad x_{ki} \sim G(\phi_{kz_{ki}}), \quad i = 1, 2, \dots, n_k\end{aligned}$$

Here is a brief description of this procedure:

1. To generate  $D_0$ , we first draw an *infinite* multinomial distribution  $\pi_0$  from a Pitman-Yor process with concentration parameter  $\alpha$ , and draw each component  $\psi_j$  from  $B$ . Then  $D_0 = \sum_{j=1}^{\infty} \pi_j \psi_j$ .
2. Then for each group (say the  $k$ -th one), we draw  $\pi_k$  from a stick breaking process and draw each component from  $D_0$ . Note that drawing a component  $\phi_{kt}$  from  $D_0$  is equivalent to choosing one of the atoms in  $D_0$ , which can be done in two steps: draw  $u_{kt}$  from  $\pi_0$  and then set  $\phi_{kt} = \psi_{u_{kt}}$ . In other words, the  $t$ -th component in the  $k$ -th group is identical to the  $u_{kt}$ -th component in  $D_0$ .
3. Finally, to generate the  $i$ -th sample in the  $k$ -th group, denoted by  $x_{ki}$ , we first draw  $z_{ki}$  from  $\pi_k$  and use the corresponding component  $\phi_{kz_{ki}}$  to generate the sample.

This formulation can be expressed using the DSL as below:

```

1  @model HierarchicalDP{B, G} begin
2      @constant m::Int          # the number of groups
3      @constant ns::Vector{Int} # the number of samples in each group
4      @hyperparam alpha::Float64 # the base concentration
5      @hyperparam gamma::Float64 # the group specific concentration
6
7      # for D0
8      pi0 ~ StickBreak(alpha)
9      for j = 1 : Inf
10         psi[j] ~ B
11     end
12
13     # each group
14     for k = 1 : m
15         pi[k] ~ StickBreak(gamma)
16
17         # Dk
18         for t = 1 : Inf
19             u[k][t] ~ pi0
20             phi[k][t] = psi[u[k][t]]
21         end
22
23         # samples
24         for i = 1 : ns[k]
25             z[k][t] ~ pi[k]
26             x[k][i] ~ G(phi[z[k][t]])
27         end
28     end
29 end

```

### 3.3 Gaussian Processes

*Gaussian process (GP)* is another important stochastic process that is widely used in Bayesian modeling. Formally, a Gaussian process is defined to be a function-valued distribution  $X_t : t \in T$ , where  $T$  can be arbitrary domain, such

that any finite subset of values in  $X_t$  is normally distributed. A Gaussian process is characterized by a mean function  $\mu : T \rightarrow R$  and a positive definite covariance function  $\kappa : T \times T \rightarrow R$ . The covariance function is typically given in a parametric form. The following is one that is widely used

$$\kappa(s, t; \theta) = \theta_0 \delta_{s,t} + \theta_1 \exp(-\theta_2(s - t)^2)$$

In many application, the GP is considered to be hidden, and observations are a noisy transformation of the samples generated from the GP, as

$$g \sim GP(\mu, \kappa) \\ x_i \sim B, \quad y_i \sim F(g(x_i)), \quad i = 1, \dots, n$$

The following model specification describes this model.

```
1 @model TransformedGaussProcess{B, F} begin
2   @constant n::Int          # the number of observed samples
3
4   # define mean and covariance function
5   @param theta::Float^3
6   mu(x) = 0.
7   kappa(x, y) = theta[1] * delta(x, y) + theta[2] * exp(- theta[3] * abs2(x - y))
8
9   # GP
10  g ~ GaussianProcess(mu, kappa)
11
12  # samples
13  for i = 1 : n
14    x[i] ~ B
15    y[i] ~ F(g[x[i]])
16  end
17 end
```

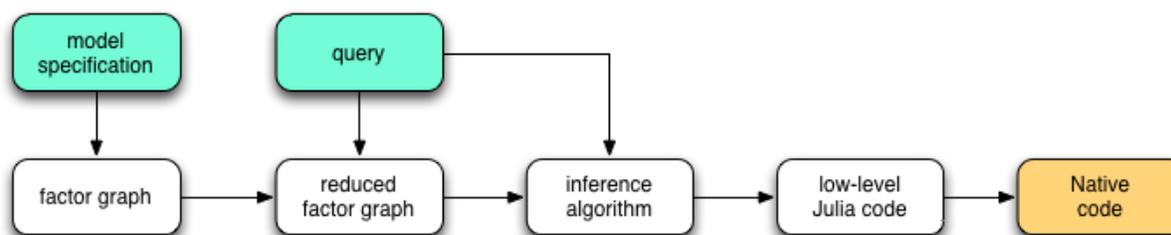
The `GaussianProcess` distribution here is a high-order stochastic function, which takes into two function arguments and generates another function. This is readily implementable in Julia, where functions are first-class citizens like in many functional programming languages.

---

## The Inference Framework

---

The diagram below outlines the overall architecture of the inference framework.



The entire process of generating inference codes from model and query specification consists of four stages:

1. The `@model` macro will convert a model specification into a model class, of which the internal representation is a factor graph.
2. The `@query` macro will reduce the factor graph based on the query. The reduction may involve following steps:
  - (a) Simplify factors by absorbing known variables. For example, a second-order factor (i.e. a factor with two arguments)  $f(x, y)$  can be reduced to a first-order factor if the value of one variable (say  $y$ ) is given.
  - (b) Eliminate irrelevant variables and factors: variables and factors that do not influence the conditional distribution of the queried variables can be safely removed. For example, consider a joint distribution  $p(x, y, z) = p_1(x|z)p_2(y|z)$ . When the value of  $y$  is given, the variable  $y$  is conditionally independent from  $x$ . Therefore, the variable  $y$  can be ignored in the inference for  $x$ .

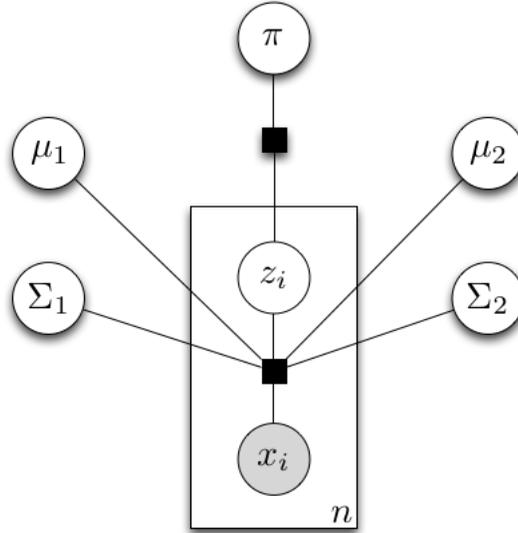
The macro `@query` also generates a high-level description of the inference algorithm.

3. An *inference compiler* will compile the inference algorithm into low-level Julia codes, taking into account the computational architecture of the target platform (e.g. CPU cores, GPU, cluster, cloud, etc).
4. Finally, the Julia compiler will emit LLVM instructions, which will then be compiled into native codes by the LLVM compiler.

Here, we focus on the first two steps, that is, compilation of model and query specifications into inference algorithms.

### 4.1 Gaussian Mixture Model Revisited

To illustrate the procedure of model compilation, let's revisit the Gaussian mixture model (GMM). Given a model specification, the `@model` macro creates a *factor graph*, which is a hyper-graph with factors connecting between variables. The following diagram is the factor graph that represents a GMM with two components.



In this graph, each sample is associated with two factors, a mixture factor that connects between observed samples, components, and component indicator  $z_i$ , and a factor that specifies the prior of each component indicator. These two factors were directly specified in the model specification. In the *model learning query*, the data  $\mathbf{x}$  are given. Therefore the order of each mixture factor is reduced from  $2K + 2$  to  $2K + 1$ .

The most widely used learning algorithm for this purpose is *Expectation Maximization*, which comprises three updating steps, as outlined below.

- Update the posterior probabilities of  $z_i$  conditioned on prior  $\pi$ , component parameters  $(\mu_k, \Sigma_k)$ , and the corresponding observed sample  $x_i$ , as

$$\begin{aligned} q_i(k) &:= p(z_i = k | \pi, \theta, x_i) \propto \pi(k) \mathcal{N}(x_i; \mu_k, \Sigma_k) \\ &= \exp(\log \pi(k) + \log \mathcal{N}(x_i; \mu_k, \Sigma_k)) \end{aligned}$$

This can be interpreted as an integration of a message from the prior factor (that is,  $\log \pi(k)$ ) and a message from the corresponding mixture factor (that is,  $\log \mathcal{N}(x_i; \mu_k, \Sigma_k)$ ). Here,  $\log \mathcal{N}(x_i; \mu_k, \Sigma_k)$  denotes the pdf at  $x_i$  with respect to a Gaussian component.

- Update the maximum likelihood estimation of  $\pi$ , as

$$\pi(k) = \frac{1}{n} \sum_{i=1}^n p(z_i = k) = \frac{1}{n} \sum_{i=1}^n \exp(\log q_i(k))$$

This can be interpreted as a computational process that combines messages from each of  $z_i$  (that is,  $\log q_i(k)$ ).

- Update the maximum likelihood estimation of component parameters  $\mu_k$  and  $\Sigma_k$ , as

$$\begin{aligned} \mu_k &= \left( \sum_{i=1}^n q_i(k) x_i \right) \left( \sum_{i=1}^n q_i(k) \right)^{-1} \\ \Sigma_k &= \left( \sum_{i=1}^n q_i(k) (x_i - \mu_k)(x_i - \mu_k)^T \right) \left( \sum_{i=1}^n q_i(k) \right)^{-1} \end{aligned}$$

Again, this is also a specific way that combines messages from each of the mixture factors associated with them.

The analysis above shows that the E-M algorithm can be viewed as a message-passing scheme that iteratively update the states associated with each variable by exchanging messages between factors and variables. Actually, the message-passing scheme is a very general paradigm. Many widely used inference algorithms, including *mean-field based variational inference*, *belief propagation*, *expectation propagation*, *Gibbs sampling*, and *Hamiltonian Monte Carlo*, can be implemented as certain forms of message passing schemes.

## 4.2 Inference via Message Passing

Generally, a message passing procedure consists of several stages:

1. Initialize the states associated with each variable. Depending on the chosen algorithm, the states associated with a variable can be in different forms. For example, for a discrete variable, its associated state can be a value (in Gibbs sampling or maximum-likelihood estimation), variational distribution (in variational inference), or a marginal distribution (in belief propagation or expectation propagation).
2. Iteratively invoke the following steps until convergence or other termination criteria are met.
  - (a) update the message from a factor to some of its incident variables (based on updated status of other incident variables).
  - (b) update variable states based on incoming messages.
3. Compute queried quantities based on variable states. For example, if a Gibbs sampling algorithm is used, the expectation of a variable can be approximated by the sample mean.

It is possible to generate such an inference procedure according to the structure of the factor graph. However, there remains several challenges to be addressed:

1. Some intermediate quantities may be used in the computation of several different messages. For example,  $\sum_{i=1}^n q_i(k)$  appears in multiple updating formulas for GMM. It is desirable to identify these quantities and avoid unnecessary re-computation of the same value.
2. Each message depends on the states of several variables, while the states of a variable may depend on several messages. A message/variable state only needs to be updated when its depending values have changed. To identify whether a variable/message needs to be updated, a natural idea is to build a dependency graph, where each node corresponds to either a variable state or a message. By time-stamping each node, it is not difficult to see whether a node can be updated by looking at the time-stamps of each neighboring nodes.
3. Updating steps can be scheduled in numerous ways. Poor scheduling may result in slow convergence. Therefore, deriving a reasonable schedule is also important to achieve high efficiency.