

# Clustering SNVs for Tumor Heterogeneity

Math-CS Sc.B Thesis Proposal

David Liu

January 9, 2017

## Problem statement

Cancer results from an evolutionary process where somatic mutations occur and accumulate in a population of cells. A mutation causes genetic variation at a genomic site called a single nucleotide variant (SNV) for the cell which obtained the mutation and all of its progeny. Thus, combinations of SNVs correspond to different subpopulations of cells (clones) which can be placed on a phylogenetic tree. This mixture of clones is a phenomenon known as intratumor heterogeneity. Each clone is commonly modeled by a binary feature vector, where each feature is an SNV.

So suppose we take biopsy samples from a tumor separated spatially or temporally. Each sample will have different proportions of clones whose relationships are unknown. Our goal is to characterize the evolutionary history of the clones (invariant across samples) and their mixing proportions (variant across samples).

The data typically used for tree inference is called the variant allele frequency (VAF), observed as follows. Suppose we bulk sequence each sample separately, so that the reads are from a sample's mixture of clones. By comparing to a control sample, if a read contains the mutated allele at a SNV, it is called a variant read; otherwise it is called a reference read. The VAF is defined for each SNV in each sample, as, in each sample, the number of variant reads at an SNV divided by the number of total reads at that SNV.

There exist algorithms to infer a tree and its population mixing proportions given perfectly accurate VAFs. However in the real world, data is noisy and this isn't so simple. Two VAFs that look different may actually be from the same clone; likewise, two VAFs that look the same may be from different clones. Thus, it is common to perform clustering to assign mutations to clusters which correspond to clones. Tree inference is then performed on these clusters.

In this thesis, I propose and implement a method to cluster mutations using the Dirichlet Process and variational inference based on a binomial mixture model, which is suited for the clone mixing problem. It could also have applications in other clustering problems.

---

## Model

Let  $m = 1, \dots, M$  index the samples and let  $n = 1, \dots, N$  index the SNVs. For each SNV  $n$  in sample  $m$ , we observe two quantities: the total number  $d_{mn}$  of reads and the number  $v_{mn}$  of variant reads. Let  $\mathbf{x}$  be the vector which encapsulates the data. Let  $k = 1, \dots, K$  index the clusters and let  $z_n$  denote the cluster assignment of SNV  $n$ , where  $z_n$  is a 1-of- $K$  indicator vector. In addition, let  $\mathbf{z}$  be the vector of all  $z_n$ .

Suppose that each clone emits variant reads according to a binomial distribution. Thus, for cluster  $k$  and some reads  $d_{mn}$ , we have variant reads distributed with  $\text{Bin}(d_{mn}, \phi_k z_{nk})$ . (**Note: May change this into a negative binomial model.**) We can consider the joint probability of a site by the product of each sample, since we assume they are independent. That is,

$$p(\mathbf{x}_{\mathbf{n}}|\phi_{\mathbf{n}}) = \prod_{m=1}^M \text{Bin}(d_{mn}, \phi_{mn}) \quad (1)$$

Further suppose that the cluster memberships  $z_n$  and weights  $\pi_k$  are generated by a Dirichlet process prior. As we can see by inspecting the graphical model, the likelihood of  $\mathbf{x}$  depends on the latent variables in a straightforward way:

$$\begin{aligned} p(\mathbf{x}_{\mathbf{n}}|\mathbf{z}, \phi) &= \prod_{k=1}^K p(\mathbf{x}_{\mathbf{n}}|\phi_{\mathbf{n}})^{z_{nk}} \\ &= \prod_{k=1}^K \prod_{m=1}^M \text{Bin}(d_{mn}, \phi_{mn})^{z_{nk}} \end{aligned} \quad (2)$$

Now consider the joint likelihood of the observed data and latent variables, which follows from (2):

$$p(\mathbf{v}_{\mathbf{n}}, \mathbf{z}|\pi, \phi_{\mathbf{n}}) = \prod_{k=1}^K \prod_{m=1}^M (\pi_k \text{Bin}(d_{mn}, \phi_{mn}))^{z_{nk}} \quad (3)$$

Figure 1 on the next page shows the graphical model.

Notation:

$n = 1, \dots, N$ : SNVs  
 $m = 1, \dots, M$ : samples  
 $k = 1, \dots, K$ : clusters

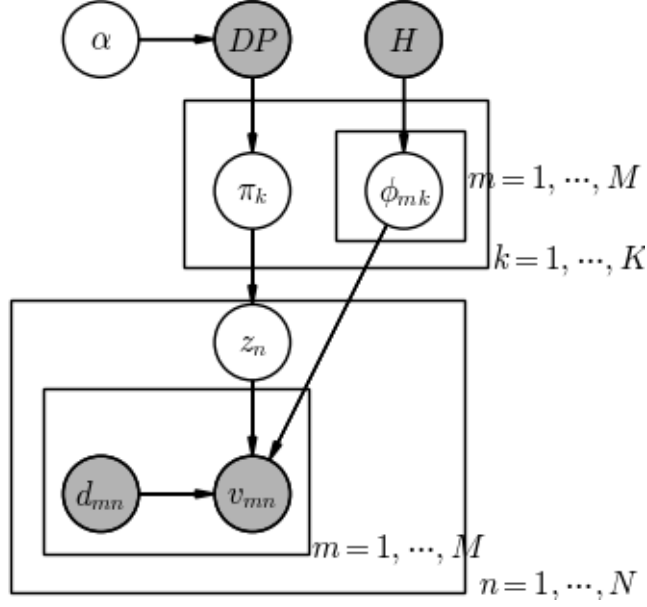


Figure 1: Graphical model for the VAFs.

$DP$  = Dirichlet process RV (stick-breaking construction)  
 $H \sim U(0, 1)$  = Base distribution for parameters  $\phi_{nk}$   
 $\pi_k$  = Weights for categorical distribution  
 $z_n \in \{1, \dots, K\} \sim \text{Categorical}(\pi_1, \dots, \pi_K)$  = Cluster membership for SNV  $n$   
 $\phi_{mk}$  = Cluster frequency  
 $v_{mn} \sim \text{Binom}(d_{m,n}, \phi_{m,c_n})$  = Observed variant reads  
 $d_{mn}$  = Observed total reads

The VAF of SNV  $n$  in sample  $m$  is equal to  $f_{m,n} = \frac{v_{mn}}{d_{mn}}$ . We can encode all sample VAFs in a  $M \times N$  matrix  $F = [f_{mn}]$  which we call the VAF matrix, which is the input for the clustering problem.

---

Suppose we have some set of clusters  $\mathcal{C} = \{C_1, \dots, C_K\}$ . Let  $C_k = \{n : c_n = k\}$ , that is  $C_k$  is all  $n$  that are in cluster  $k$ . Further suppose that we have some matrix of inferred VAFs,  $\Phi$ . Pool reads as follows:

$$\tilde{d}_{m,k} = \sum_{n \in C_k} d_{m,n} \quad (4)$$

$$\tilde{v}_{m,k} = \sum_{n \in C_k} v_{m,n} \quad (5)$$

Then the likelihood of some VAF matrix  $F$  under some clustering  $\mathcal{C}$  is

$$P(F|\mathcal{C}, \Phi) = \prod_{m=1}^M \prod_{k=1}^K \binom{\tilde{d}_{m,k}}{\tilde{v}_{m,k}} (\Phi_{m,k})^{\tilde{v}_{m,k}} (1 - \Phi_{m,k})^{\tilde{d}_{m,k} - \tilde{v}_{m,k}} \quad (6)$$

where the outer product is over all samples and the inner product is over the clusters in a sample.

The posterior involves a Dirichlet Process prior, which is analytically intractable. We must use some sort of computational technique to perform inference on this posterior.

## Variational Inference

Variational inference is an alternative to MCMC-based inference methods. At a high level, variational inference factors a posterior using the mean-field approximation, which approximates the posterior in a higher-dimensional space using simpler independent functions. Then a simple coordinate ascent can be performed in order to infer the model parameters.

Now we wish to use variational inference to infer a  $\mathcal{C}, \Phi$  for our posterior. This requires the posterior to be in the exponential family (source).

**Lemma 1.** *The product of two exponential family functions is also exponential.*

*Proof.* Exponential family functions can be parameterized by

$$f_X(x | \theta) = h(x) \exp(\theta^T \cdot T(x) - A(\theta)).$$

So suppose we have  $f_1, f_2$  which are correspondingly parameterized. Note that  $T(x)$  is the same because the sufficient statistic is invariant across a distribution. Then the product is

$$\begin{aligned} f_1 \times f_2 &= h_1(x) \exp(\theta_1^T \cdot T(x) - A(\theta_1)) \times h_2(x) \exp(\theta_2^T \cdot T(x) - A(\theta_2)) \\ &= \tilde{h}(x) \exp((\theta_1 + \theta_2)^T \cdot T(x) - \tilde{A}(\theta_1, \theta + 2)) \end{aligned}$$

which is also an exponential family distribution.  $\square$

**Claim 1.** *The posterior is also of the exponential family.*

*Proof.* The DP prior is also exponential (source). Then the claim is trivial by lemma 1.  $\square$

---

## Details on Variational Inference

### The ELBO

Let  $\mathbf{z}$  denote the latent variables, and  $\mathbf{x}$  denote the data. We seek to approximate the posterior  $p(\mathbf{z}|\mathbf{x})$  from a family of distributions  $\mathcal{D}$  by solving the following optimization problem:

$$q^*(z) = \arg \min_{q(\mathbf{z}) \in \mathcal{D}} \text{KL}((q(\mathbf{z})||p(\mathbf{z}|\mathbf{x}))). \quad (7)$$

where KL is the KL-divergence, which measures the “distance” between two distributions.

However, (4) requires us to compute the log evidence (which is intractable over the space of all  $\mathbf{z}$ ) since

$$\text{KL}(q(\mathbf{z})||p(\mathbf{z}|\mathbf{x})) = \text{E}[\log q(\mathbf{z})] - \text{E}[\log p(\mathbf{z}, \mathbf{x})] + \log p(\mathbf{x}). \quad (8)$$

Instead, we optimize an objective function which is not dependent on  $\log p(\mathbf{x})$ . We call this the evidence lower bound (ELBO), which is equal to the negative KL-divergence plus the log evidence.

$$\text{ELBO}(q) = \text{Elog } p(\mathbf{z}, \mathbf{x}) - \text{Elog } q(\mathbf{z}). \quad (9)$$

and thus we see that  $\log p(\mathbf{x})$  is a constant with respect to  $q$ . The ELBO gets its name from the fact that it is a lower bound for the log evidence.

### The mean-field variational family

And what family of distributions do we use for  $\mathcal{D}$ ? The standard technique is to use a simple one from physics, the mean-field variational family. In this family, the latent variables  $\mathbf{z}$  are mutually independent so that the joint distribution factorizes:

$$q(\mathbf{z}) = \prod_{j=1}^m q_j(z_j). \quad (10)$$

where  $q_j$  is a bounded variation dependent only on  $z_j$ . The structure of the model will dictate the optimal form of  $q_j$ .

### Mean-field assumptions break dependency to give us coordinate ascent

The optimization is solved using a coordinate ascent algorithm, where the independence of the latent variables gives us orthogonality. Let  $z_{-j}$  denote the set of latent variables  $z_l$  such that  $l \neq j$ . Consider the complete conditional of  $z_j$ , which is a function of the other latent variables and the data,  $p(z_j | \mathbf{z}_{-j}, \mathbf{x})$ . Since the expectation in the ELBO is with respect to  $q(\mathbf{z})$ , which we have assumed factorizes, then we can dissect out the dependence with respect to  $\mathbf{z}_j$  by using (6) and (7):

$$\begin{aligned} \text{ELBO}(q) &= \int \prod q_i(\mathbf{z}_i) \left( \log p(\mathbf{z}, \mathbf{x}) - \sum_i \log q_i(\mathbf{z}_i) \right) d\mathbf{z} \\ &\propto \int q_j(z_j) \text{E}_{-j} [\log p(\mathbf{x}, \mathbf{z})] d\mathbf{z}_j - \int q_j(z_j) \log q_j(\mathbf{z}_j) d\mathbf{z}_j \end{aligned} \quad (11)$$

Now suppose that we fix  $z_{-j}$  and maximize the ELBO. Then the ELBO is maximized when  $\log q_j(\mathbf{z}_j) \propto \mathbb{E}_{-j} [\log p(\mathbf{x}, \mathbf{z})]$ , by the positivity of the KL-divergence. Thus the optimal  $q^*(\mathbf{z}_j)$  occurs when

$$q_j^*(\mathbf{z}_j) \propto \exp(\mathbb{E}_{-j} [\log p(\mathbf{x}, \mathbf{z})]) \quad (12)$$

(9) underlies the coordinate-ascent variational inference algorithm. By iterating through each variational factor, fixing the others, and performing coordinate ascent (similar to Gibbs sampling), then we eventually reach a local optimum of the ELBO.

CAVI

---

**Algorithm 1:** CAVI

---

**Input:** A model  $p(\mathbf{x}, \mathbf{z})$ , a data set  $\mathbf{x}$

**Output:** A variational density  $q(\mathbf{z}) = \prod_{j=1}^m q_j(z_j)$

**Initialize:** Variational factors  $q_j(z_j)$

**while** the ELBO has not converged **do**

**for**  $j \in \{1, \dots, m\}$  **do**

        Set  $q_j(z_j) \propto \exp\{\mathbb{E}_{-j}[\log p(z_j | \mathbf{z}_{-j}, \mathbf{x})]\}$

**end**

    Compute  $\text{ELBO}(q) = \mathbb{E}[\log p(\mathbf{z}, \mathbf{x})] + \mathbb{E}[\log q(\mathbf{z})]$

**end**

**return**  $q(\mathbf{z})$

---

## Variational Inference on Multi-sample Binomial Model

So suppose that our posterior factors as follows:

$$q(\mathbf{z}, \Phi | F) = \underbrace{\prod_{k=1}^{\infty} q(\vec{\phi}_k)}_{\text{Observation parameters}} \times \underbrace{\prod_{k=1}^{\infty} q(\vec{w}_k)}_{\text{Cluster proportions}} \times \underbrace{\prod_{n=1}^N q(z_n)}_{\text{Assignment of data to clusters}} \quad (13)$$

$$= \underbrace{\prod_{k=1}^{\infty} \prod_{m=1}^M q(\phi_{km})}_{\text{Observation parameters}} \times \underbrace{\prod_{k=1}^{\infty} \prod_{m=1}^M q(w_{km})}_{\text{Cluster proportions}} \times \underbrace{\prod_{n=1}^N q(z_n)}_{\text{Assignment of data to clusters}} \quad (14)$$

Where we make the variational approximations based on sufficient statistic parameters (everything with a hat):

$$q(\phi_{km}) = P(\phi_{km} | \hat{\tau}_{km}, \hat{\nu}_{km})$$

$$q(w_{km}) = \text{Beta}(w_{km} | \hat{\eta}_{km0}, \hat{\eta}_{km1})$$

$$q(z_n) = \text{Cat}_{\infty}(z_n | \hat{r}_{n1}, \dots, \hat{r}_{nk})$$

$q(\phi_{km}) = P(\phi_{km} | \hat{\tau}_{km}, \hat{\nu}_{km})$  is something I still have to figure out.

The actual update equations are also something I still have to figure out.

---

## **Implmentation**

This will either be implemented in Python or R.

## **Evaluating results**

The clusterings will be evaluated on simulated and real tumor data. It may also be used in conjunction with other algorithms which use such clusterings as input, to see if this improves their performance.