

Clustering SNVs for Tumor Heterogeneity

Math-CS Sc.B Thesis Proposal

David Liu

September 20, 2016

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 (groups) and let $n = 1, \dots, N$ index the SNVs (features). For each SNV n in sample m , we observe two quantities: the total number $d_{m,n}$ of reads and the number $v_{m,n}$ of variant reads. Let $k = 1, \dots, K$ index the clusters and let c_n denote the cluster assignment of SNV n .

Suppose that each clone emits variant reads according to a binomial distribution with according to $\text{Binomial}(d_{m,n}, \phi_{m,c_n})$. (**Note: May change this into a negative binomial model.**)

Further suppose that the cluster memberships are generated by a Dirichlet process, which is the prior for the clustering model. The cluster membership c_i of an SNV i is shared across samples and is distributed according to $\text{Categorical}(w_1, \dots, w_K)$. Figure 1 on the next page shows the graphical model.

Notation:

$n = 1, \dots, N$: SNVs

$m = 1, \dots, M$: samples (groups)

$k = 1, \dots, K$: clusters

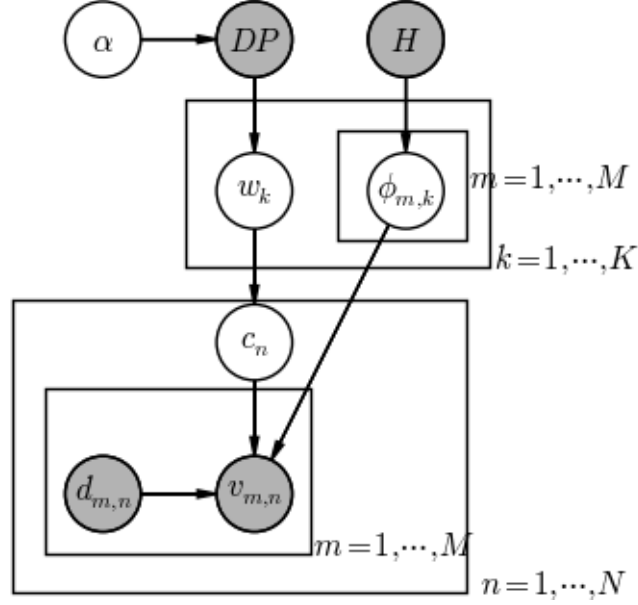


Figure 1: Graphical model for the VAFs.

DP = Dirichlet process RV (stick-breaking construction)

$H \sim U(0, 1)$ = Base distribution

w_k = Weights for categorical distribution

$c_n \in \{1, \dots, K\} \sim \text{Categorical}(w_1, \dots, w_K)$ = Cluster membership for SNV n

$\phi_{m,k}$ = Cluster frequency

$v_{m,n} \sim \text{Binom}(d_{m,n}, \phi_{m,c_n})$ = Observed variant reads

$d_{m,n}$ = Observed total reads

The VAF of SNV n in sample m is equal to $f_{m,n} = \frac{v_{m,n}}{d_{m,n}}$. 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, Φ . Pool reads as follows:

$$\tilde{d}_{m,k} = \sum_{n \in C_k} d_{m,n} \quad (1)$$

$$\tilde{v}_{m,k} = \sum_{n \in C_k} v_{m,n} \quad (2)$$

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

$$P(F|\mathcal{C}, \Phi) = \prod_{m=1}^M \prod_{k=1}^K \binom{\tilde{d}_{m,k}}{\tilde{v}_{m,k}} (\Phi_{m,k})^{\tilde{v}_{m,k}} (1 - \Phi_{m,k})^{\tilde{d}_{m,k} - \tilde{v}_{m,k}} \quad (3)$$

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. 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}, Φ 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

Variational Inference on Multi-sample Binomial Model

So suppose that our posterior factors as follows:

$$q(\mathcal{C}, \Phi | F) = \underbrace{\prod_{k=1}^{\infty} q(\vec{\phi}_k)}_{\text{Observation parameters}} \times \underbrace{\prod_{k=1}^{\infty} q(\vec{w}_k)}_{\text{Cluster proportions}} \times \underbrace{\prod_{n=1}^N q(z_n)}_{\text{Assignment of data to clusters}} \quad (4)$$

$$= \underbrace{\prod_{k=1}^{\infty} \prod_{m=1}^M q(\phi_{km})}_{\text{Observation parameters}} \times \underbrace{\prod_{k=1}^{\infty} \prod_{m=1}^M q(w_{km})}_{\text{Cluster proportions}} \times \underbrace{\prod_{n=1}^N q(z_n)}_{\text{Assignment of data to clusters}} \quad (5)$$

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

$$\begin{aligned} q(\phi_{km}) &= P(\phi_{km} | \hat{\tau}_{km}, \hat{\nu}_{km}) \\ q(w_{km}) &= \text{Beta}(w_{km} | \hat{\eta}_{km0}, \hat{\eta}_{km1}) \\ q(z_n) &= \text{Cat}_{\infty}(z_n | \hat{r}_{n1}, \dots, \hat{r}_{nk}) \end{aligned}$$

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