

Likelihood Ratio Tests

General Set-up

Most hypothesis testing procedures can be formulated so that a test statistic, $S(\mathbf{x})$ is applied to the data $\mathbf{x} = (x_1, x_2, \dots, x_n)^T$ so that:

1. $S(\mathbf{x}) \geq 0$
2. The larger $S(\mathbf{x})$ is, the more significant the test is (i.e., the more evidence against the null in favor of the alternative)
3. The p-value is $p(\mathbf{x}) = \Pr(S(\mathbf{X}^*) \geq S(\mathbf{x}))$ where $S(\mathbf{X}^*)$ is distributed according to the null distribution

Significance Regions

A level α test is a significance rule (i.e., a rule for calling a test statistically significant) that results in a false positive rate (i.e., Type I error rate) of α . Under our set-up, significance regions take the form:

$$\Gamma_\alpha = \{\mathbf{x} : S(\mathbf{x}) \geq c_{1-\alpha}\},$$

where $c_{1-\alpha}$ is the $(1 - \alpha)$ percentile of $S(\mathbf{X}^*)$ so that $\Pr(S(\mathbf{X}^*) \geq c_{1-\alpha}) = \alpha$. We restrict $0 < \alpha < 1$.

Note that if $\alpha' \leq \alpha$ then $\Gamma_{\alpha'} \subseteq \Gamma_\alpha$.

P-values

A p-value can be defined in terms of significance regions:

$$p(\mathbf{x}) = \min\{\alpha : \mathbf{x} \in \Gamma_\alpha\}$$

Example: Wald Test

Consider the hypothesis test, $H_0 : \theta = \theta_0$ vs $H_1 : \theta \neq \theta_0$. Let $\hat{\theta}_n(\mathbf{x})$ be the MLE of θ . We have

$$S(\mathbf{x}) = \frac{|\hat{\theta}_n(\mathbf{x}) - \theta_0|}{\text{sê}(\hat{\theta}_n(\mathbf{x}))},$$

$$\Gamma_\alpha = \{\mathbf{x} : S(\mathbf{x}) \geq |z_{\alpha/2}|\},$$

where $z_{\alpha/2}$ is the $\alpha/2$ percentile of the Normal(0, 1) distribution.

Neyman-Pearson Lemma

Suppose we are testing $H_0 : \theta = \theta_0$ vs $H_1 : \theta = \theta_1$ where in practice θ_0 and θ_1 are known, fixed quantities. The most powerful test has statistic and significance regions:

$$S(\mathbf{x}) = \frac{f(\mathbf{x}; \theta_1)}{f(\mathbf{x}; \theta_0)} = \frac{L(\theta_1; \mathbf{x})}{L(\theta_0; \mathbf{x})}$$

$$\Gamma_\alpha = \{\mathbf{x} : S(\mathbf{x}) \geq c_{1-\alpha}\},$$

where $c_{1-\alpha}$ is the $(1 - \alpha)$ percentile of $S(\mathbf{X}^*)$ so that $\Pr_{\theta_0}(S(\mathbf{X}^*) \geq c_{1-\alpha}) = \alpha$.

Simple vs. Composite Hypotheses

A simple hypothesis is defined in terms of a single value, e.g.,

- $H_0 : \mu = 0$
- $H_0 : p = p_0$ where p_0 is a placeholder for a known, fixed number in practice
- $H_1 : \lambda = 5$

A composite hypothesis is defined by multiple values, e.g.,

- $H_0 : \mu \leq 0$ vs $H_1 : \mu > 0$
- $H_0 : p_1 = p_2$, where p_1 and p_2 are two unknown parameters corresponding to two populations
- $H_1 : \mu \neq 0$

General Hypothesis Tests

Let $X_1, X_2, \dots, X_n \stackrel{\text{iid}}{\sim} F_\theta$ where $\theta \in \Theta$. Let $\Theta_0, \Theta_1 \subseteq \Theta$ so that $\Theta_0 \cap \Theta_1 = \emptyset$ and $\Theta_0 \cup \Theta_1 = \Theta$. The hypothesis test is:

$$H_0 : \theta \in \Theta_0 \text{ vs } H_1 : \theta \in \Theta_1$$

If Θ_0 or Θ_1 contain more than one value then the corresponding hypothesis is composite.

Composite H_0

The significance regions indexed by their level α are determined so that:

$$\Gamma_\alpha = \{\mathbf{x} : S(\mathbf{x}) \geq c_{1-\alpha}\},$$

where $c_{1-\alpha}$ is such that

$$\max_{\theta \in \Theta_0} \Pr(S(\mathbf{X}^*) \geq c_{1-\alpha}) = \alpha.$$

In this case,

$$\begin{aligned} p(\mathbf{x}) &= \min \{\alpha : \mathbf{x} \in \Gamma_\alpha\} \\ &= \max_{\theta \in \Theta_0} \Pr(S(\mathbf{X}^*) \geq S(\mathbf{x})) \end{aligned}$$

Generalized LRT

Let $X_1, X_2, \dots, X_n \stackrel{\text{iid}}{\sim} F_\theta$ where $\theta \in \Theta$ and we are testing $H_0 : \theta \in \Theta_0$ vs $H_1 : \theta \in \Theta_1$.

The generalized LRT utilizes test statistic and significance regions:

$$\lambda(\mathbf{x}) = \frac{\max_{\theta \in \Theta} L(\theta; \mathbf{x})}{\max_{\theta \in \Theta_0} L(\theta; \mathbf{x})} = \frac{L(\hat{\theta}; \mathbf{x})}{L(\hat{\theta}_0; \mathbf{x})}$$

$$\Gamma_\alpha = \{\mathbf{x} : \lambda(\mathbf{x}) \geq c_{1-\alpha}\}$$

Null Distribution of Gen. LRT

The null distribution of $\lambda(\mathbf{x})$ under “certain regularity assumptions” can be shown to be such that, as $n \rightarrow \infty$,

$$2 \log \lambda(\mathbf{x}) \xrightarrow{D} \chi_v^2$$

where $v = \dim(\Theta) - \dim(\Theta_0)$.

The significance regions can be more easily written as $\Gamma_\alpha = \{\mathbf{x} : 2 \log \lambda(\mathbf{x}) \geq c_{1-\alpha}\}$ where $c_{1-\alpha}$ is the $1 - \alpha$ percentile of the χ_v^2 distribution.

Example: Poisson

Let $X_1, X_2, \dots, X_n \stackrel{\text{iid}}{\sim} \text{Poisson}(\theta)$ where $\theta > 0$ and we are testing $H_0 : \theta = \theta_0$ vs $H_1 : \theta \neq \theta_0$. The unconstrained MLE is $\hat{\theta} = \bar{x}$. The generalized LRT statistic

$$2 \log \lambda(\mathbf{x}) = 2 \log \frac{e^{-n\hat{\theta}} \hat{\theta}^{\sum x_i}}{e^{-n\theta_0} \theta_0^{\sum x_i}} = 2n \left[(\theta_0 - \hat{\theta}) - \hat{\theta} \log(\theta_0/\hat{\theta}) \right]$$

which has an asymptotic χ_1^2 null distribution.

Example: Normal

Let $X_1, X_2, \dots, X_n \stackrel{\text{iid}}{\sim} \text{Normal}(\mu, \sigma^2)$ and we are testing $H_0 : \mu = \mu_0$ vs $H_1 : \mu \neq \mu_0$. The generalized LRT can be applied for multidimensional parameter spaces Θ as well. The statistic, which has asymptotic null distribution χ_1^2 , is

$$2 \log \lambda(\mathbf{x}) = 2 \log \left(\frac{\hat{\sigma}_0^2}{\hat{\sigma}^2} \right)^{n/2}$$

where

$$\hat{\sigma}_0^2 = \frac{\sum_{i=1}^n (x_i - \mu_0)^2}{n}, \quad \hat{\sigma}^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n}.$$

Bayesian Classification

Assumptions

Let $(X_1, X_2, \dots, X_n) | \theta \stackrel{\text{iid}}{\sim} F_\theta$ where $\theta \in \Theta$ and $\theta \sim F_\tau$. Let $\Theta_0, \Theta_1 \subseteq \Theta$ so that $\Theta_0 \cap \Theta_1 = \emptyset$ and $\Theta_0 \cup \Theta_1 = \Theta$.

Given observed data \mathbf{x} , we wish to classify whether $\theta \in \Theta_0$ or $\theta \in \Theta_1$.

This is the Bayesian analog of hypothesis testing.

Prior Probability on H

Let H be a rv such that $H = 0$ when $\theta \in \Theta_0$ and $H = 1$ when $\theta \in \Theta_1$.

From the prior distribution on θ , we can calculate

$$\Pr(H = 0) = \int_{\theta \in \Theta_0} f(\theta) d\theta$$

and $\Pr(H = 1) = 1 - \Pr(H = 0)$.

Posterior Probability

Using Bayes theorem, we can also calculate

$$\begin{aligned} \Pr(H = 0 | \mathbf{x}) &= \frac{f(\mathbf{x} | H = 0) \Pr(H = 0)}{f(\mathbf{x})} \\ &= \frac{\int_{\theta \in \Theta_0} f(\mathbf{x} | \theta) f(\theta) d\theta}{\int_{\theta \in \Theta} f(\mathbf{x} | \theta) f(\theta) d\theta} \end{aligned}$$

where note that $\Pr(H = 1 | \mathbf{x}) = 1 - \Pr(H = 0 | \mathbf{x})$.

Loss Function

Let $\mathcal{L}(\tilde{H}, H)$ be such that

$$\begin{aligned} \mathcal{L}(\tilde{H} = 1, H = 0) &= c_I \\ \mathcal{L}(\tilde{H} = 0, H = 1) &= c_{II} \end{aligned}$$

for some $c_I, c_{II} > 0$.

Bayes Risk

The Bayes risk, $R(\tilde{H}, H)$, is

$$\begin{aligned} \mathbb{E}[\mathcal{L}(\tilde{H}, H) | \mathbf{x}] &= c_I \Pr(\tilde{H} = 1, H = 0) + c_{II} \Pr(\tilde{H} = 0, H = 1) \\ &= c_I \Pr(\tilde{H} = 1 | H = 0) \Pr(H = 0) \\ &\quad + c_{II} \Pr(\tilde{H} = 0 | H = 1) \Pr(H = 1) \end{aligned}$$

Notice how this balances what frequentists call Type I error and Type II error.

Bayes Rule

The estimate \tilde{H} that minimizes $R(\tilde{H}, H)$ is

$$\tilde{H} = 1 \text{ when } \Pr(H = 1|\mathbf{x}) \geq \frac{c_I}{c_I + c_{II}}$$

and $\tilde{H} = 0$ otherwise.

Numerical Methods for Likelihood Functions

Why Numerical Methods for Likelihood

Challenges

Frequentist model:

$$X_1, X_2, \dots, X_n \stackrel{\text{iid}}{\sim} F_{\theta}$$

Bayesian model:

$$X_1, X_2, \dots, X_n | \theta \stackrel{\text{iid}}{\sim} F_{\theta} \text{ and } \theta \sim F_{\tau}$$

Sometimes it's not possible to find formulas for $\hat{\theta}_{\text{MLE}}$, $\hat{\theta}_{\text{MAP}}$, $E[\theta|\mathbf{x}]$, or $f(\theta|\mathbf{x})$. We have to use numerical methods instead.

Approaches

We will discuss the following numerical approaches to likelihood based inference:

- Expectation-maximization (EM) algorithm
- Variational inference
- Markov chain Monte Carlo (MCMC)
 - Metropolis sampling
 - Metropolis-Hastings sampling
 - Gibbs sampling

Latent Variable Models

Definition

Latent variables (or hidden variables) are random variables that are present in the model, but unobserved.

We will denote latent variables by Z , and we will assume

$$(X_1, Z_1), (X_2, Z_2), \dots, (X_n, Z_n) \stackrel{\text{iid}}{\sim} F_{\theta}.$$

A realized value of Z is z , $\mathbf{Z} = (Z_1, Z_2, \dots, Z_n)^T$, etc.

The EM algorithm and variational inference involve latent variables.

Bayesian models are a special case of latent variable models: the unobserved random parameters are latent variables.

Empirical Bayes Revisited

In the earlier EB example, we supposed that $X_i | \mu_i \sim \text{Normal}(\mu_i, 1)$ for $i = 1, 2, \dots, n$ where these rv's are independent, and also that $\mu_i \stackrel{\text{iid}}{\sim} \text{Normal}(a, b^2)$.

The unobserved parameters $\mu_1, \mu_2, \dots, \mu_n$ are latent variables. In this case, $\theta = (a, b^2)$.

Normal Mixture Model

Suppose $X_1, X_2, \dots, X_n \stackrel{\text{iid}}{\sim} F_{\boldsymbol{\theta}}$ where $\boldsymbol{\theta} = (\pi_1, \dots, \pi_K, \mu_1, \dots, \mu_K, \sigma_1^2, \dots, \sigma_K^2)$ with pdf

$$f(\mathbf{x}; \boldsymbol{\theta}) = \prod_{i=1}^n \sum_{k=1}^K \pi_k \frac{1}{\sqrt{2\pi\sigma_k^2}} \exp\left\{-\frac{(x_i - \mu_k)^2}{2\sigma_k^2}\right\}.$$

The MLEs of the unknown parameters cannot be found analytically. This is a mixture common model to work with in applications, so we need to be able to estimate the parameters.

There is a latent variable model that produces the same marginal distribution and likelihood function. Let $\mathbf{Z}_1, \mathbf{Z}_2, \dots, \mathbf{Z}_n \stackrel{\text{iid}}{\sim} \text{Multinomial}_K(1, \boldsymbol{\pi})$ where $\boldsymbol{\pi} = (\pi_1, \dots, \pi_K)$. Note that $Z_{ik} \in \{0, 1\}$ and $\sum_{k=1}^K Z_{ik} = 1$. Let $[X_i | Z_{ik} = 1] \sim \text{Normal}(\mu_k, \sigma_k^2)$, where $\{X_i | \mathbf{Z}_i\}_{i=1}^n$ are jointly independent.

The joint pdf is

$$f(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}) = \prod_{i=1}^n \prod_{k=1}^K \left[\pi_k \frac{1}{\sqrt{2\pi\sigma_k^2}} \exp\left\{-\frac{(x_i - \mu_k)^2}{2\sigma_k^2}\right\} \right]^{z_{ik}}.$$

Note that

$$f(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}) = \prod_{i=1}^n f(x_i, \mathbf{z}_i; \boldsymbol{\theta}).$$

It can be verified that $f(\mathbf{x}; \boldsymbol{\theta})$ is the marginal distribution of this latent variable model:

$$f(x_i; \boldsymbol{\theta}) = \sum_{\mathbf{z}_i} f(x_i, \mathbf{z}_i; \boldsymbol{\theta}) = \sum_{k=1}^K \pi_k \frac{1}{\sqrt{2\pi\sigma_k^2}} \exp\left\{-\frac{(x_i - \mu_k)^2}{2\sigma_k^2}\right\}.$$

Bernoulli Mixture Model

Suppose $X_1, X_2, \dots, X_n \stackrel{\text{iid}}{\sim} F_{\boldsymbol{\theta}}$ where $\boldsymbol{\theta} = (\pi_1, \dots, \pi_K, p_1, \dots, p_K)$ with pdf

$$f(\mathbf{x}; \boldsymbol{\theta}) = \prod_{i=1}^n \sum_{k=1}^K \pi_k p_k^{x_i} (1 - p_k)^{1-x_i}.$$

As in the Normal mixture model, the MLEs of the unknown parameters cannot be found analytically.

As before, there is a latent variable model that produces the same marginal distribution and likelihood function. Let $\mathbf{Z}_1, \mathbf{Z}_2, \dots, \mathbf{Z}_n \stackrel{\text{iid}}{\sim} \text{Multinomial}_K(1, \boldsymbol{\pi})$ where $\boldsymbol{\pi} = (\pi_1, \dots, \pi_K)$. Note that $Z_{ik} \in \{0, 1\}$ and $\sum_{k=1}^K Z_{ik} = 1$. Let $[X_i | Z_{ik} = 1] \sim \text{Bernoulli}(p_k)$, where $\{X_i | \mathbf{Z}_i\}_{i=1}^n$ are jointly independent.

The joint pdf is

$$f(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}) = \prod_{i=1}^n \prod_{k=1}^K [p_k^{x_i} (1 - p_k)^{1-x_i}]^{z_{ik}}.$$

EM Algorithm

Rationale

For any likelihood function, $L(\boldsymbol{\theta}; \mathbf{x}) = f(\mathbf{x}; \boldsymbol{\theta})$, there is an abundance of optimization methods that can be used to find the MLE or MAP. However:

- Optimization methods can be messy to implement
- There may be probabilistic structure that we can use to simplify the optimization process and also provide theoretical guarantees on its convergence
- Optimization isn't necessarily the only goal, but one may also be interested in point estimates of the latent variable values

Requirement

The expectation-maximization (EM) algorithm allows us to calculate MLEs and MAPs when certain geometric properties are satisfied in the probabilistic model.

In order for the EM algorithm to be a practical approach, then we should have a latent variable model $f(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})$ that is used to do inference on $f(\mathbf{x}; \boldsymbol{\theta})$ or $f(\boldsymbol{\theta}|\mathbf{x})$.

Note: Sometimes (\mathbf{x}, \mathbf{z}) is called the **complete data** and \mathbf{x} is called the **observed data** when we are using the EM as a method for dealing with missing data.

The Algorithm

1. Choose initial value $\boldsymbol{\theta}^{(0)}$
2. Calculate $f(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}^{(t)})$
3. Calculate

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t)}) = E_{\mathbf{Z}|\mathbf{X}=\mathbf{x}} \left[\log f(\mathbf{x}, \mathbf{Z}; \boldsymbol{\theta}); \boldsymbol{\theta}^{(t)} \right]$$

4. Set

$$\boldsymbol{\theta}^{(t+1)} = \operatorname{argmax}_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t)})$$

5. Iterate until convergence and set $\hat{\boldsymbol{\theta}} = \boldsymbol{\theta}^{(\infty)}$

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t)})$$

Continuous \mathbf{Z} :

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t)}) = \int \log f(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}) f(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)}) d\mathbf{z}$$

Discrete \mathbf{Z} :

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t)}) = \sum_{\mathbf{z}} \log f(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}) f(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{(t)})$$

EM for MAP

If we wish to calculate the MAP we replace $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t)})$ with

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t)}) = E_{\mathbf{Z}|\mathbf{X}=\mathbf{x}} \left[\log f(\mathbf{x}, \mathbf{Z}; \boldsymbol{\theta}); \boldsymbol{\theta}^{(t)} \right] + \log f(\boldsymbol{\theta})$$

where $f(\boldsymbol{\theta})$ is the prior distribution on $\boldsymbol{\theta}$.

EM Examples

Normal Mixture Model

Returning to the Normal mixture model introduced earlier, we first calculate

$$\log f(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}) = \sum_{i=1}^n \sum_{k=1}^K z_{ik} \log \pi_k + z_{ik} \log \phi(x_i; \mu_k, \sigma_k^2)$$

where

$$\phi(x_i; \mu_k, \sigma_k^2) = \frac{1}{\sqrt{2\pi\sigma_k^2}} \exp \left\{ -\frac{(x_i - \mu_k)^2}{2\sigma_k^2} \right\}.$$

In calculating

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t)}) = \mathbb{E}_{\mathbf{Z}|\mathbf{X}=\mathbf{x}} \left[\log f(\mathbf{x}, \mathbf{Z}; \boldsymbol{\theta}); \boldsymbol{\theta}^{(t)} \right]$$

we only need to know $\mathbb{E}_{\mathbf{Z}|\mathbf{X}=\mathbf{x}}[Z_{ik}|\mathbf{x}; \boldsymbol{\theta}]$, which turns out to be

$$\mathbb{E}_{\mathbf{Z}|\mathbf{X}=\mathbf{x}}[Z_{ik}|\mathbf{x}; \boldsymbol{\theta}] = \frac{\pi_k \phi(x_i; \mu_k, \sigma_k^2)}{\sum_{j=1}^K \pi_j \phi(x_i; \mu_j, \sigma_j^2)}.$$

Note that we take

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t)}) = \mathbb{E}_{\mathbf{Z}|\mathbf{X}=\mathbf{x}} \left[\log f(\mathbf{x}, \mathbf{Z}; \boldsymbol{\theta}); \boldsymbol{\theta}^{(t)} \right]$$

so the parameter in $\log f(\mathbf{x}, \mathbf{Z}; \boldsymbol{\theta})$ is a free $\boldsymbol{\theta}$, but the parameters used to take the conditional expectation of \mathbf{Z} are fixed at $\boldsymbol{\theta}^{(t)}$. Let's define

$$\hat{z}_{ik}^{(t)} = \mathbb{E} \left[z_{ik} | \mathbf{x}; \boldsymbol{\theta}^{(t)} \right] = \frac{\pi_k^{(t)} \phi(x_i; \mu_k^{(t)}, \sigma_k^{2,(t)})}{\sum_{j=1}^K \pi_j^{(t)} \phi(x_i; \mu_j^{(t)}, \sigma_j^{2,(t)})}.$$

E-Step

We calculate

$$\begin{aligned} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t)}) &= \mathbb{E}_{\mathbf{Z}|\mathbf{X}=\mathbf{x}} \left[\log f(\mathbf{x}, \mathbf{Z}; \boldsymbol{\theta}); \boldsymbol{\theta}^{(t)} \right] \\ &= \sum_{i=1}^n \sum_{k=1}^K \hat{z}_{ik}^{(t)} \log \pi_k + \hat{z}_{ik}^{(t)} \log \phi(x_i; \mu_k, \sigma_k^2) \end{aligned}$$

At this point the parameters making up $\hat{z}_{ik}^{(t)}$ are fixed at $\boldsymbol{\theta}^{(t)}$.

M-Step

We now calculate $\theta^{(t+1)} = \operatorname{argmax}_{\theta} Q(\theta, \theta^{(t)})$, which yields:

$$\begin{aligned}\pi_k^{(t+1)} &= \frac{\sum_{i=1}^n \hat{z}_{ik}^{(t)}}{n} \\ \mu_k^{(t+1)} &= \frac{\sum_{i=1}^n \hat{z}_{ik}^{(t)} x_i}{\sum_{i=1}^n \hat{z}_{ik}^{(t)}} \\ \sigma_k^{2,(t+1)} &= \frac{\sum_{i=1}^n \hat{z}_{ik}^{(t)} (x_i - \mu_k^{(t+1)})^2}{\sum_{i=1}^n \hat{z}_{ik}^{(t)}}\end{aligned}$$

Note: You need to use a Lagrange multiplier to obtain $\{\pi_k^{(t+1)}\}_{k=1}^K$.

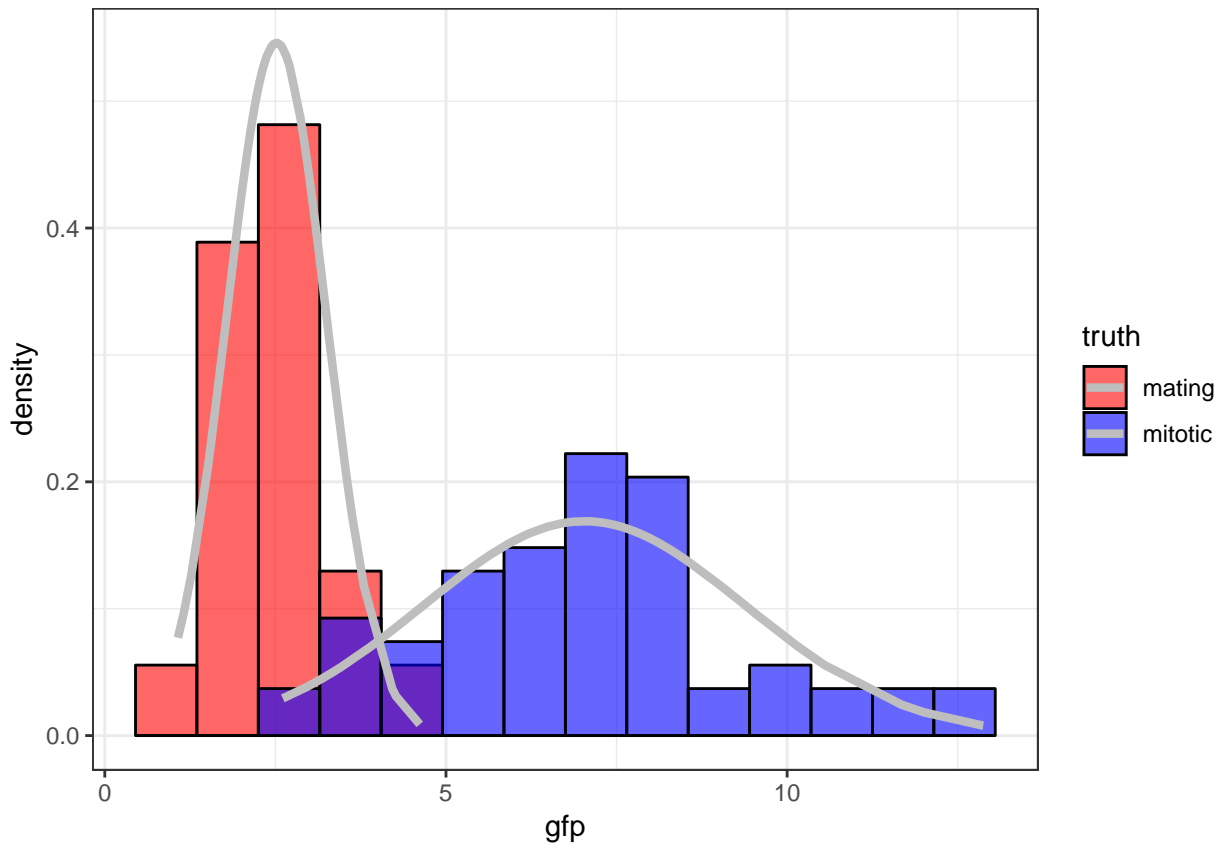
Caveat

If we assign one and only one data point to mixture component k , meaning $\mu_k^{(t)} = x_i$ and $\hat{z}_{ik}^{(t)} = 1$ for some k and i , then as $\sigma_k^{2,(t)} \rightarrow 0$, the likelihood goes to ∞ .

Therefore, when implementing the EM algorithm for this particular Normal mixture model, we have to be careful to bound all $\sigma_k^{2,(t)}$ away from zero and avoid this scenario.

Yeast Gene Expression

Measured ratios of the nuclear to cytoplasmic fluorescence for a protein-GFP construct that is hypothesized as being nuclear in mitotic cells and largely cytoplasmic in mating cells.



Initialize Values

```
> set.seed(508)
> B <- 100
> p <- rep(0,B)
> mu1 <- rep(0,B)
> mu2 <- rep(0,B)
> s1 <- rep(0,B)
> s2 <- rep(0,B)
> p[1] <- runif(1, min=0.1, max=0.9)
> mu.start <- sample(x, size=2, replace=FALSE)
> mu1[1] <- min(mu.start)
> mu2[1] <- max(mu.start)
> s1[1] <- var(sort(x)[1:60])
> s2[1] <- var(sort(x)[61:120])
> z <- rep(0,120)
```

Run EM Algorithm

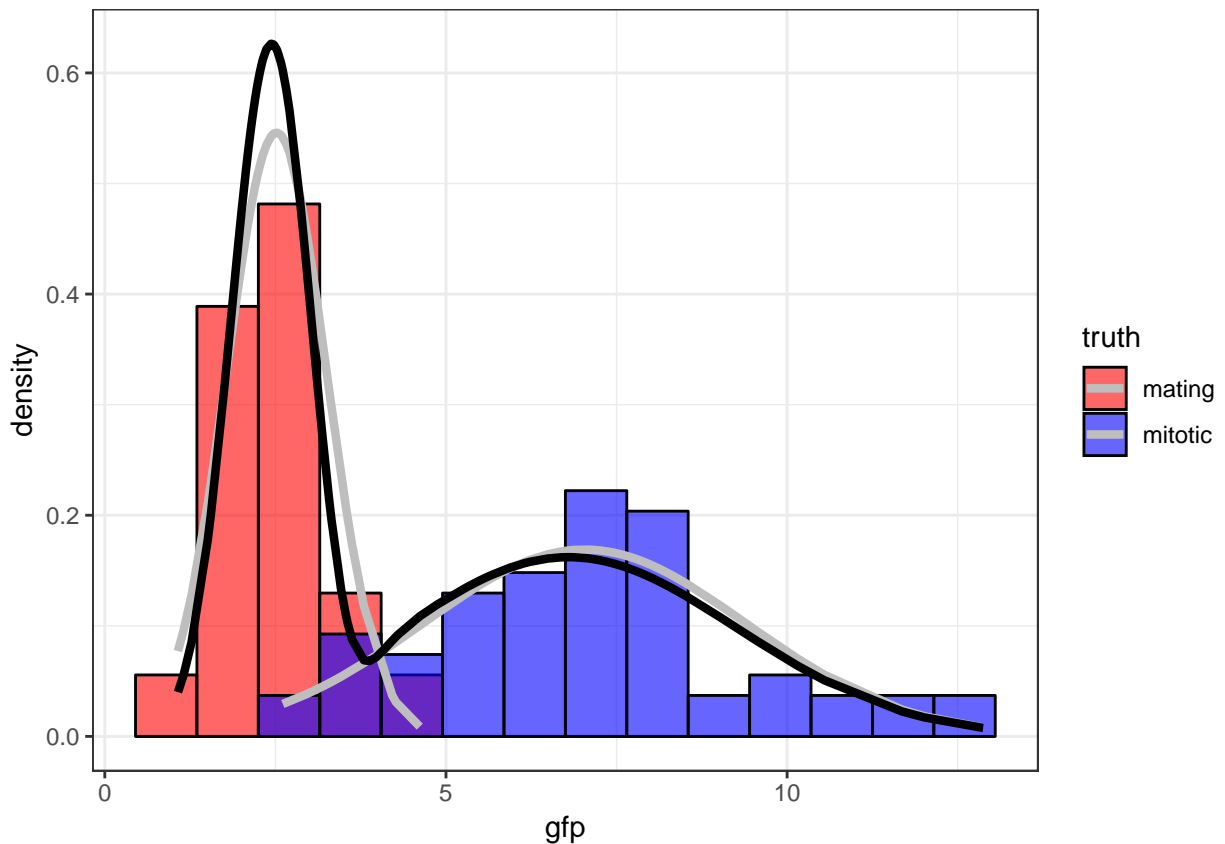
```
> for(i in 2:B) {
+   z <- (p[i-1]*dnorm(x, mean=mu2[i-1], sd=sqrt(s2[i-1])))/
+     (p[i-1]*dnorm(x, mean=mu2[i-1], sd=sqrt(s2[i-1])) +
+      (1-p[i-1])*dnorm(x, mean=mu1[i-1], sd=sqrt(s1[i-1])))
+   mu1[i] <- sum((1-z)*x)/sum(1-z)
+   mu2[i] <- sum(z*x)/sum(z)
```

```

+ s1[i] <- sum((1-z)*(x-mu1[i])^2)/sum(1-z)
+ s2[i] <- sum(z*(x-mu2[i])^2)/sum(z)
+ p[i] <- sum(z)/length(z)
+ }
>
> tail(cbind(mu1, s1, mu2, s2, p), n=3)
      mu1      s1      mu2      s2      p
[98,] 2.455325 0.3637967 6.7952 6.058291 0.5340015
[99,] 2.455325 0.3637967 6.7952 6.058291 0.5340015
[100,] 2.455325 0.3637967 6.7952 6.058291 0.5340015

```

Fitted Mixture Distribution



Bernoulli Mixture Model

As an exercise, derive the EM algorithm of the Bernoulli mixture model introduced earlier.

Hint: Replace $\phi(x_i; \mu_k, \sigma_k^2)$ with the appropriate Bernoulli pmf.

Other Applications of EM

- Dealing with missing data
- Multiple imputation of missing data
- Truncated observations
- Bayesian hyperparameter estimation
- Hidden Markov models

Theory of EM

Decomposition

Let $q(\mathbf{z})$ be a probability distribution on the latent variables, \mathbf{z} . Consider the following decomposition:

$$\log f(\mathbf{x}; \boldsymbol{\theta}) = \mathcal{L}(q(\mathbf{z}), \boldsymbol{\theta}) + \text{KL}(q(\mathbf{z}) \| f(\mathbf{z} | \mathbf{x}; \boldsymbol{\theta}))$$

where

$$\begin{aligned}\mathcal{L}(q(\mathbf{z}), \boldsymbol{\theta}) &= \int q(\mathbf{z}) \log \left(\frac{f(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})}{q(\mathbf{z})} \right) d\mathbf{z} \\ \text{KL}(q(\mathbf{z}) \| f(\mathbf{z} | \mathbf{x}; \boldsymbol{\theta})) &= - \int q(\mathbf{z}) \log \left(\frac{f(\mathbf{z} | \mathbf{x}; \boldsymbol{\theta})}{q(\mathbf{z})} \right) d\mathbf{z}\end{aligned}$$

Kullback-Leibler Divergence

The KL divergence provides an asymmetric measure of the difference between two probability distributions.

The KL divergence is such that $\text{KL}(q \| f) \geq 0$ where $\text{KL}(q \| f) = 0$ if and only if $q = f$. This property is known as **Gibbs inequality**.

Lower Bound

Note that $\mathcal{L}(q(\mathbf{z}), \boldsymbol{\theta})$ provides a lower bound on the likelihood function:

$$\log f(\mathbf{x}; \boldsymbol{\theta}) \geq \mathcal{L}(q(\mathbf{z}), \boldsymbol{\theta})$$

If we set $q(\mathbf{z}) = f(\mathbf{z} | \mathbf{x}; \boldsymbol{\theta}^{(t)})$, then for a fixed $\boldsymbol{\theta}^{(t)}$ and as a function of $\boldsymbol{\theta}$,

$$\begin{aligned}\mathcal{L}(q(\mathbf{z}), \boldsymbol{\theta}) &\propto \int f(\mathbf{z} | \mathbf{x}; \boldsymbol{\theta}^{(t)}) \log f(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}) d\mathbf{z} \\ &= Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t)})\end{aligned}$$

EM Increases Likelihood

Since $\boldsymbol{\theta}^{(t+1)} = \text{argmax}_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t)})$, it follows that

$$Q(\boldsymbol{\theta}^{(t+1)}, \boldsymbol{\theta}^{(t)}) \geq Q(\boldsymbol{\theta}^{(t)}, \boldsymbol{\theta}^{(t)}).$$

Also, by the properties of KL divergence stated above, we have

$$\text{KL}(f(\mathbf{z} | \mathbf{x}; \boldsymbol{\theta}^{(t+1)}) \| f(\mathbf{z} | \mathbf{x}; \boldsymbol{\theta}^{(t)})) \geq \text{KL}(f(\mathbf{z} | \mathbf{x}; \boldsymbol{\theta}^{(t)}) \| f(\mathbf{z} | \mathbf{x}; \boldsymbol{\theta}^{(t)})).$$

Putting these together we have

$$\log f(\mathbf{x}; \boldsymbol{\theta}^{(t+1)}) \geq \log f(\mathbf{x}; \boldsymbol{\theta}^{(t)}).$$

Variational Inference

Rationale

Performing the EM algorithm required us to be able to compute $f(\mathbf{z}|\mathbf{x};\boldsymbol{\theta})$ and also optimize $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t)})$. Sometimes this is not possible. Variational inference takes advantage of the decomposition

$$\log f(\mathbf{x};\boldsymbol{\theta}) = \mathcal{L}(q(\mathbf{z}), \boldsymbol{\theta}) + \text{KL}(q(\mathbf{z})\|f(\mathbf{z}|\mathbf{x};\boldsymbol{\theta}))$$

and instead considers other forms of $q(\mathbf{z})$ to identify a more tractable optimization.

Optimization Goal

Since

$$\log f(\mathbf{x};\boldsymbol{\theta}) = \mathcal{L}(q(\mathbf{z}), \boldsymbol{\theta}) + \text{KL}(q(\mathbf{z})\|f(\mathbf{z}|\mathbf{x};\boldsymbol{\theta}))$$

it follows that the closer $q(\mathbf{z})$ is to $f(\mathbf{z}|\mathbf{x};\boldsymbol{\theta})$, the term $\mathcal{L}(q(\mathbf{z}), \boldsymbol{\theta})$ grows larger while $\text{KL}(q(\mathbf{z})\|f(\mathbf{z}|\mathbf{x};\boldsymbol{\theta}))$ becomes smaller. The goal is typically to identify a restricted form of $q(\mathbf{z})$ that maximizes $\mathcal{L}(q(\mathbf{z}), \boldsymbol{\theta})$, which serves as an approximation to the posterior distribution $f(\mathbf{z}|\mathbf{x};\boldsymbol{\theta})$.

Mean Field Approximation

A mean field approximation implies we restrict $q(\mathbf{z})$ to be

$$q(\mathbf{z}) = \prod_{k=1}^K q_k(\mathbf{z}_k)$$

for some partition $\mathbf{z} = (\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_K)$. This partition is very context specific and is usually driven by the original model and what is tractable.

Optimal $q_k(\mathbf{z}_k)$

Under the above restriction, it can be shown that the $\{q_k(\mathbf{z}_k)\}$ that maximize $\mathcal{L}(q(\mathbf{z}), \boldsymbol{\theta})$ have the form:

$$q_k(\mathbf{z}_k) \propto \exp \left\{ \int \log f(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}) \prod_{j \neq k} q_j(\mathbf{z}_j) d\mathbf{z}_j \right\}.$$

These pdf's or pmf's can be calculated iteratively by cycling over $k = 1, 2, \dots, K$ after initializing them appropriately. Note that convergence is guaranteed.

Remarks

- If $\boldsymbol{\theta}$ is also random, then it can be included in \mathbf{z} .
- The estimated $\hat{f}(\mathbf{z}|\mathbf{x})$ is typically concentrated around the high density region of the true $f(\mathbf{z}|\mathbf{x})$, so it is useful for calculations such as the MAP, but it is not guaranteed to be a good overall estimate of $f(\mathbf{z}|\mathbf{x})$.
- Variational inference is typically faster than MCMC (covered next).
- Given this is an optimization procedure, care can be taken to speed up convergence and avoid unintended local maxima.

Markov Chain Monte Carlo

Motivation

When performing Bayesian inference, it is often (but not always) possible to calculate

$$f(\boldsymbol{\theta}|\mathbf{x}) \propto L(\boldsymbol{\theta}; \mathbf{x})f(\boldsymbol{\theta})$$

but it is typically much more difficult to calculate

$$f(\boldsymbol{\theta}|\mathbf{x}) = \frac{L(\boldsymbol{\theta}; \mathbf{x})f(\boldsymbol{\theta})}{f(\mathbf{x})}.$$

Markov chain Monte Carlo is a method for simulating data approximately from $f(\boldsymbol{\theta}|\mathbf{x})$ with knowledge of only $L(\boldsymbol{\theta}; \mathbf{x})f(\boldsymbol{\theta})$.

Note

MCMC can be used to approximately simulate data from any distribution that is only proportionally characterized, but it is probably most well known for doing so in the context of Bayesian inference.

We will explain MCMC in the context of Bayesian inference.

Big Picture

We draw a Markov chain of $\boldsymbol{\theta}$ values so that, in some asymptotic sense, these are equivalent to iid draws from $f(\boldsymbol{\theta}|\mathbf{x})$.

The draws are done competitively so that the next draw of a realization of $\boldsymbol{\theta}$ depends on the current value.

The Markov chain is set up so that it only depends on $L(\boldsymbol{\theta}; \mathbf{x})f(\boldsymbol{\theta})$.

A lot of practical decisions need to be made by the user, so utilize MCMC carefully.

Metropolis-Hastings Algorithm

1. Initialize $\boldsymbol{\theta}^{(0)}$
2. Generate $\boldsymbol{\theta}^* \sim q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(b)})$ for some pdf or pmf $q(\cdot|\cdot)$
3. With probability

$$A(\boldsymbol{\theta}^*, \boldsymbol{\theta}^{(b)}) = \min \left(1, \frac{L(\boldsymbol{\theta}^*; \mathbf{x})f(\boldsymbol{\theta}^*)q(\boldsymbol{\theta}^{(b)}|\boldsymbol{\theta}^*)}{L(\boldsymbol{\theta}^{(b)}; \mathbf{x})f(\boldsymbol{\theta}^{(b)})q(\boldsymbol{\theta}^*|\boldsymbol{\theta}^{(b)})} \right)$$

set $\boldsymbol{\theta}^{(b+1)} = \boldsymbol{\theta}^*$. Otherwise, set $\boldsymbol{\theta}^{(b+1)} = \boldsymbol{\theta}^{(b)}$

4. Continue for $b = 1, 2, \dots, B$ iterations and *carefully* select which $\boldsymbol{\theta}^{(b)}$ are utilized to approximate iid observations from $f(\boldsymbol{\theta}|\mathbf{x})$

Metropolis Algorithm

The Metropolis algorithm restricts $q(\cdot, \cdot)$ to be symmetric so that $q(\boldsymbol{\theta}^{(b)}|\boldsymbol{\theta}^*) = q(\boldsymbol{\theta}^*|\boldsymbol{\theta}^{(b)})$ and

$$A(\boldsymbol{\theta}^*, \boldsymbol{\theta}^{(b)}) = \min \left(1, \frac{L(\boldsymbol{\theta}^*; \mathbf{x})f(\boldsymbol{\theta}^*)}{L(\boldsymbol{\theta}^{(b)}; \mathbf{x})f(\boldsymbol{\theta}^{(b)})} \right).$$

Utilizing MCMC Output

Two common uses of the output from MCMC are as follows:

1. $E[f(\boldsymbol{\theta})|\mathbf{x}]$ is approximated by

$$\hat{E}[f(\boldsymbol{\theta})|\mathbf{x}] = \frac{1}{B} \sum_{b=1}^B f(\boldsymbol{\theta}^{(b)}).$$

2. Some subsequence $\boldsymbol{\theta}^{(b_1)}, \boldsymbol{\theta}^{(b_2)}, \dots, \boldsymbol{\theta}^{(b_m)}$ from $\{\boldsymbol{\theta}^{(b)}\}_{b=1}^B$ is utilized as an empirical approximation to iid draws from $f(\boldsymbol{\theta}|\mathbf{x})$.

Remarks

- The random draw $\boldsymbol{\theta}^* \sim q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(b)})$ perturbs the current value $\boldsymbol{\theta}^{(b)}$ to the next value $\boldsymbol{\theta}^{(b+1)}$. It is often a Normal distribution for continuous $\boldsymbol{\theta}$.
- Choosing the variance of $q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(b)})$ is important as it requires enough variance for the theory to be applicable within a reasonable number of computations, but it cannot be so large that new values of $\boldsymbol{\theta}^{(b+1)}$ are rarely generated.
- $A(\boldsymbol{\theta}^*, \boldsymbol{\theta}^{(b)})$ is called the acceptance probability.
- The algorithm must be run for a certain number of iterations (“burn in”) before observed $\boldsymbol{\theta}^{(b)}$ can be utilized.
- The generated $\boldsymbol{\theta}^{(b)}$ are typically “thinned” (only sampled every so often) to reduce Markov dependence.

Full Conditionals

Suppose that $\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_K)$. Define the subset vector as $\boldsymbol{\theta}_{a:b} = (\theta_a, \theta_{a+1}, \dots, \theta_{b-1}, \theta_b)$ for any $1 \leq a \leq b \leq K$.

The full conditional of θ_k is

$$\Pr(\theta_k | \boldsymbol{\theta}_{1:k-1}, \boldsymbol{\theta}_{k+1:K}, \mathbf{x})$$

Gibbs Sampling

Gibbs sampling a special type of Metropolis-Hasting MCMC. The algorithm samples one coordinate of $\boldsymbol{\theta}$ at a time.

1. Initialize $\boldsymbol{\theta}^{(0)}$.
2. Sample:
 $\theta_1^{(b+1)} \sim \Pr(\theta_1 | \boldsymbol{\theta}_{2:K}^{(b)}, \mathbf{x})$
 $\theta_2^{(b+1)} \sim \Pr(\theta_2 | \theta_1^{(b+1)}, \boldsymbol{\theta}_{3:K}^{(b)}, \mathbf{x})$
 $\theta_3^{(b+1)} \sim \Pr(\theta_3 | \boldsymbol{\theta}_{1:2}^{(b+1)}, \boldsymbol{\theta}_{3:K}^{(b)}, \mathbf{x})$
 \vdots
 $\theta_K^{(b+1)} \sim \Pr(\theta_K | \boldsymbol{\theta}_{1:K-1}^{(b+1)}, \mathbf{x})$

3. Continue for $b = 1, 2, \dots, B$ iterations.

Gibbs and MH

As an exercise, show that Gibbs sampling is a special case of the Metropolis-Hastings algorithm where $A(\boldsymbol{\theta}^*, \boldsymbol{\theta}^{(b)}) = 1$.

Latent Variables

Note that MCMC is often used to calculate a posterior distribution on latent variables.

This makes sense because unobserved random parameters are a special type of latent variable.

Theory

The goal of MCMC is to construct a Markov chain that converges to a stationary distribution that is equivalent to the target probability distribution.

Under reasonably general assumptions, one can show that the Metropolis-Hastings algorithm produces a Markov chain that is *homogeneous* and achieves *detailed balance*, which implies the Markov chain is *ergodic* so that $\boldsymbol{\theta}^{(B)}$ converges in distribution to $f(\boldsymbol{\theta}|\mathbf{x})$ as $B \rightarrow \infty$ and that

$$\hat{E}[f(\boldsymbol{\theta})|\mathbf{x}] = \frac{1}{B} \sum_{b=1}^B f(\boldsymbol{\theta}^{(b)}) \xrightarrow{B \rightarrow \infty} E[f(\boldsymbol{\theta})|\mathbf{x}].$$

Software

Stan is probably the currently most popular software for doing Bayesian computation, including MCMC and variational inference.

There are also popular R packages, such as `MCMCpack`.

MCMC Example

Single Nucleotide Polymorphisms

PSD Admixture Model

PSD model proposed in Pritchard, Stephens, Donnelly (2000) *Genetics*.

Gibbs Sampling Approach

The Bayesian Gibbs sampling approach to inferring the PSD model touches on many important ideas, such as conjugate priors and mixture models.

We will focus on a version of this model for diploid SNPs.

The Data

\mathbf{X} , a $L \times N$ matrix consisting of the genotypes, coded as 0, 1, 2. Each row is a SNP, each column is an individual.

In order for this model to work, the data needs to be broken down into “phased” genotypes. For the 0 and 2 cases, it’s obvious how to do this, and for the 1 case, it’ll suffice for this model to randomly assign the alleles to chromosomes. We will explore phasing more on HW4.

Thus, we wind up with two $\{0, 1\}$ *binary* matrices \mathbf{X}_A and \mathbf{X}_B , both $L \times N$. We will refer to allele A and allele B . Note $\mathbf{X} = \mathbf{X}_A + \mathbf{X}_B$.

Model Components

- K , the number of populations that we model the genotypes as admixtures of. This is chosen before inference.

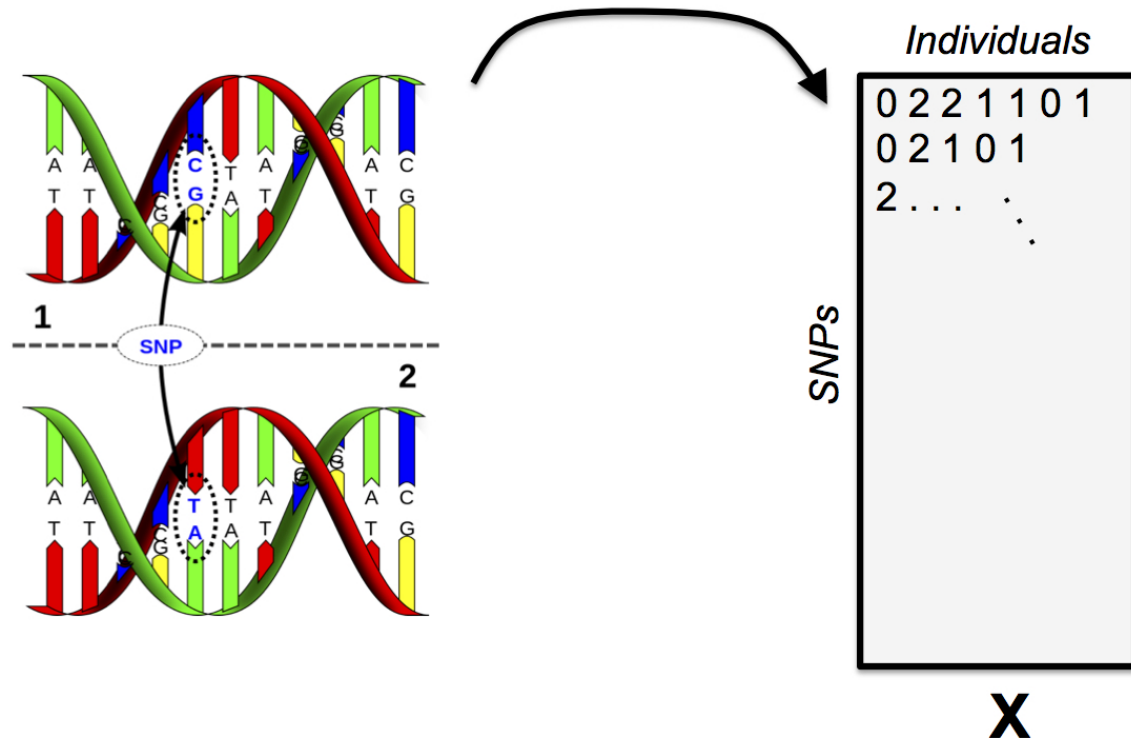


Figure 1: SNPs

- \mathbf{Q} , a $N \times K$ matrix, the admixture proportions, values are in the interval $[0, 1]$ and rows are constrained to sum to 1.
- \mathbf{P} , a $L \times K$ matrix, the allele frequencies for each population, values are in the interval $[0, 1]$.
- \mathbf{Z}_A and \mathbf{Z}_B , two $L \times N$ matrices that tell us which population the respective allele is from. Elements consist of the integers between 1 and K . This is a hidden variable.

The Model

- Each allele (elements of \mathbf{X}_A and \mathbf{X}_B) is a Bernoulli random variable, with success probability determined by which population that allele is assigned to (i.e., depends on \mathbf{Z}_A , \mathbf{Z}_B , and \mathbf{P}).
- We put a uniform Beta prior, i.e., $\text{Beta}(1, 1)$, on each element of \mathbf{P} .
- We put a uniform Dirichlet prior, i.e., $\text{Dirichlet}(1, \dots, 1)$, on each *row* of \mathbf{Q} .
- \mathbf{Z}_A and \mathbf{Z}_B are K -class Multinomial draws where the probability of drawing each class is determined by each row of \mathbf{Q} .

Conditional Independence

The key observation is to understand which parts of the model are dependent on each other in the data generating process.

- The data \mathbf{X}_A and \mathbf{X}_B depends directly on \mathbf{Z}_A , \mathbf{Z}_B , and \mathbf{P} (not \mathbf{Q} !).
- The latent variable \mathbf{Z}_A and \mathbf{Z}_B depend only on \mathbf{Q} and they're conditionally independent given \mathbf{Q} .
- \mathbf{Q} and \mathbf{P} depend only on their priors.

$$\Pr(\mathbf{X}_A, \mathbf{X}_B, \mathbf{Z}_A, \mathbf{Z}_B, \mathbf{P}, \mathbf{Q}) = \Pr(\mathbf{X}_A, \mathbf{X}_B | \mathbf{Z}_A, \mathbf{Z}_B, \mathbf{P}) \Pr(\mathbf{Z}_A | \mathbf{Q}) \Pr(\mathbf{Z}_B | \mathbf{Q}) \Pr(\mathbf{P}) \Pr(\mathbf{Q})$$

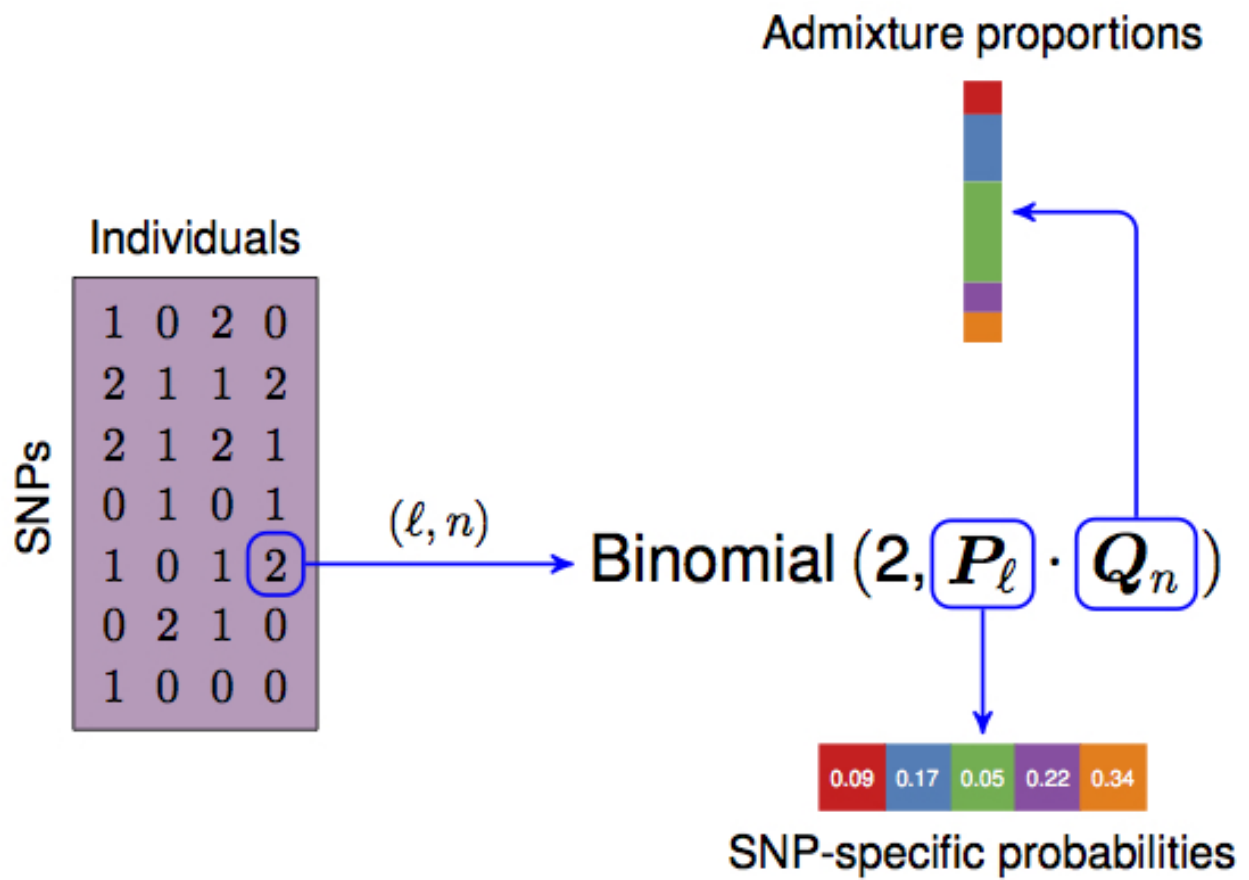


Figure 2: PSD

The Posterior

We desire to compute the posterior distribution $\Pr(\mathbf{P}, \mathbf{Q}, \mathbf{Z}_A, \mathbf{Z}_B | \mathbf{X}_A, \mathbf{X}_B)$. Gibbs sampling tells us if we can construct conditional distributions for each random variable in our model, then iteratively sampling and updating our model parameters will result in a stationary distribution that is the same as the posterior distribution.

Gibbs sampling is an extremely powerful approach for this model because we can utilize conjugate priors as well as the independence of various parameters in the model to compute these conditional distributions.

Full Conditional for \mathbf{Q}

Note that \mathbf{Z}_A and \mathbf{Z}_B are the only parts of this model that directly depend on \mathbf{Q} .

$$\begin{aligned} & \Pr(Q_n | \mathbf{Q}_{-n}, \mathbf{Z}_A, \mathbf{Z}_B, \mathbf{P}, \mathbf{X}_A, \mathbf{X}_B) \\ &= \Pr(Q_n | \mathbf{Z}_A, \mathbf{Z}_B) \\ &\propto \Pr(Z_{An}, Z_{Bn} | Q_n) \Pr(Q_n) \\ &= \Pr(Z_{An} | Q_n) \Pr(Z_{Bn} | Q_n) \Pr(Q_n) \\ &\propto \left(\prod_{\ell=1}^L \prod_{k=1}^K Q_{nk}^{\mathbb{1}(Z_{An\ell}=k) + \mathbb{1}(Z_{Bn\ell}=k)} \right) \\ &= \prod_{k=1}^K Q_{nk}^{S_{nk}} \end{aligned}$$

where S_{nk} is simply the count of the number of alleles for individual n that got assigned to population k .

Thus, $Q_n | \mathbf{Z}_A, \mathbf{Z}_B \sim \text{Dirichlet}(S_{j1} + 1, \dots, S_{jk} + 1)$.

We could have guessed that this distribution is Dirichlet given that \mathbf{Z}_A and \mathbf{Z}_B are multinomial! Let's use conjugacy to help us in the future.

Full Conditional for \mathbf{P}

$$\begin{aligned} & \Pr(P_\ell | \mathbf{P}_{-\ell}, \mathbf{Z}_A, \mathbf{Z}_B, \mathbf{Q}, \mathbf{X}_A, \mathbf{X}_B) \\ &\propto \Pr(\mathbf{X}_{A\ell}, \mathbf{X}_{B\ell} | P_\ell, \mathbf{Z}_{A\ell}, \mathbf{Z}_{B\ell}) \Pr(P_\ell) \end{aligned}$$

We know $\Pr(\mathbf{X}_{A\ell}, \mathbf{X}_{B\ell} | P_\ell, \mathbf{Z}_{A\ell}, \mathbf{Z}_{B\ell})$ will be Bernoulli and $\Pr(P_\ell)$ will be beta, so the full conditional will be beta as well. In fact, the prior is uniform so it vanishes from the RHS.

Thus, all we have to worry about is the Bernoulli portion $\Pr(\mathbf{X}_{A\ell}, \mathbf{X}_{B\ell} | P_\ell, \mathbf{Z}_{A\ell}, \mathbf{Z}_{B\ell})$. Here, we observe that if the \mathbf{Z}_A and \mathbf{Z}_B are "known", then we know which value of P_ℓ to plug into our Bernoulli for \mathbf{X}_A and \mathbf{X}_B . Following the Week 6 lectures, we find that the full conditional for \mathbf{P} is:

$$P_{\ell k} | \mathbf{Z}_A, \mathbf{Z}_B, \mathbf{X}_A, \mathbf{X}_B \sim \text{Beta}(1 + T_{\ell k 0}, 1 + T_{\ell k 1})$$

where $T_{\ell k 0}$ is the total number of 0 alleles at SNP ℓ for population k , and $T_{\ell k 1}$ is the analogous quantity for the 1 allele.

Full Conditional Z_A & Z_B

We'll save some math by first noting that alleles A and B are independent of each other, so we can write this for only Z_A without losing any information. Also, all elements of Z_A are independent of each other. Further, note that each element of Z_A is a single multinomial draw, so we are working with a discrete random variable.

$$\begin{aligned} & \Pr(Z_{A\ell n} = k | \mathbf{X}_A, \mathbf{Q}, \mathbf{P}) \\ &= \Pr(Z_{A\ell n} = k | X_{A\ell n}, Q_n, P_\ell) \\ &\propto \Pr(X_{A\ell n} | Z_{A\ell n} = k, Q_n, P_\ell) \Pr(Z_{A\ell n} = k | Q_n, P_\ell) \end{aligned}$$

We can look at the two factors. First:

$$\Pr(Z_{A\ell n} = k | Q_n, P_\ell) = \Pr(Z_{A\ell n} = k | Q_n) = Q_{nk}$$

Then:

$$\Pr(X_{A\ell n} | Z_{A\ell n} = k, Q_n, P_\ell) = P_{\ell k}$$

Thus, we arrive at the formula:

$$\Pr(Z_{A\ell n} = k | \mathbf{X}_A, \mathbf{Q}, \mathbf{P}) \propto P_{\ell k} Q_{nk}$$

Gibbs Sampling Updates

It's neat that we wind up just iteratively counting the various discrete random variables along different dimensions.

$$\begin{aligned} Q_n | \mathbf{Z}_A, \mathbf{Z}_B &\sim \text{Dirichlet}(S_{j_1} + 1, \dots, S_{j_k} + 1) \\ P_{\ell k} | \mathbf{Z}_A, \mathbf{Z}_B, \mathbf{X}_A, \mathbf{X}_B &\sim \text{Beta}(1 + T_{\ell k 0}, 1 + T_{\ell k 1}) \\ Z_{A\ell n} | \mathbf{X}_A, \mathbf{Q}, \mathbf{P} &\sim \text{Multinomial}\left(\frac{P_\ell * Q_n}{P_\ell \cdot Q_n}\right) \end{aligned}$$

where $*$ means element-wise vector multiplication.

Implementation

The Markov chain property means that we can't use vectorization forward in time, so R is not the best way to implement this algorithm.

That being said, we can vectorize the pieces that we can and demonstrate what happens.

Matrix-wise `rdirichlet` Function

Drawing from a Dirichlet is easy and vectorizable because it consists of normalizing independent gamma draws.

```
> rdirichlet <- function(alpha) {
+   m <- nrow(alpha)
+   n <- ncol(alpha)
+   x <- matrix(rgamma(m * n, alpha), ncol = n)
+   x/rowSums(x)
+ }
```

Inspect Data

```
> dim(Xa)
[1] 400 24
> X[1:3,1:3]
      NA18516 NA19138 NA19137
rs2051075    0      1      2
rs765546     2      2      0
rs10019399   2      2      2
> Xa[1:3,1:3]
      NA18516 NA19138 NA19137
rs2051075    0      0      1
rs765546     1      1      0
rs10019399   1      1      1
```

Model Parameters

```
> L <- nrow(Xa)
> N <- ncol(Xa)
>
> K <- 3
>
> Za <- matrix(sample(1:K, L*N, replace=TRUE), L, N)
> Zb <- matrix(sample(1:K, L*N, replace=TRUE), L, N)
> P <- matrix(0, L, K)
> Q <- matrix(0, N, K)
```

Update P

```
> update_P <- function() {
+   Na_0 <- Za * (Xa==0)
+   Na_1 <- Za * (Xa==1)
+   Nb_0 <- Zb * (Xb==0)
+   Nb_1 <- Zb * (Xb==1)
+   for(k in 1:K) {
+     N0 <- rowSums(Na_0==k)+rowSums(Nb_0==k)
+     N1 <- rowSums(Na_1==k)+rowSums(Nb_1==k)
+     P[,k] <- rdirichlet(1+cbind(N1, N0))[,1]
+   }
+   P
+ }
```

Update Q

```
> update_Q <- function() {
+   M_POP0 <- apply(Za, 2, function(x) {tabulate(x, nbins=K)})
+   M_POP1 <- apply(Zb, 2, function(x) {tabulate(x, nbins=K)})
+
+   rdirichlet(t(1+M_POP0+M_POP1))
+ }
```

Update (Each) Z

```
> update_Z <- function(X) {
+   Z <- matrix(0, nrow(X), ncol(X))
+   for(n in 1:N) {
+     PZ0 <- t(t((1-P)) * Q[n,])
+     PZ1 <- t(t(P) * Q[n,])
+     PZ <- X[,n]*PZ1 + (1-X[,n])*PZ0
+     Z[,n] <- apply(PZ, 1, function(p){sample(1:K, 1, prob=p)})
+   }
+   Z
+ }
```

Model Log-likelihood Function

```
> model_ll <- function() {
+   AFa <- t(sapply(1:L, function(i){P[i,][Za[i,]]}))
+   AFb <- t(sapply(1:L, function(i){P[i,][Zb[i,]]}))
+   # hint, hint, HW3
+   sum(dbinom(Xa, 1, AFa, log=TRUE)) +
+   sum(dbinom(Xb, 1, AFb, log=TRUE))
+ }
```

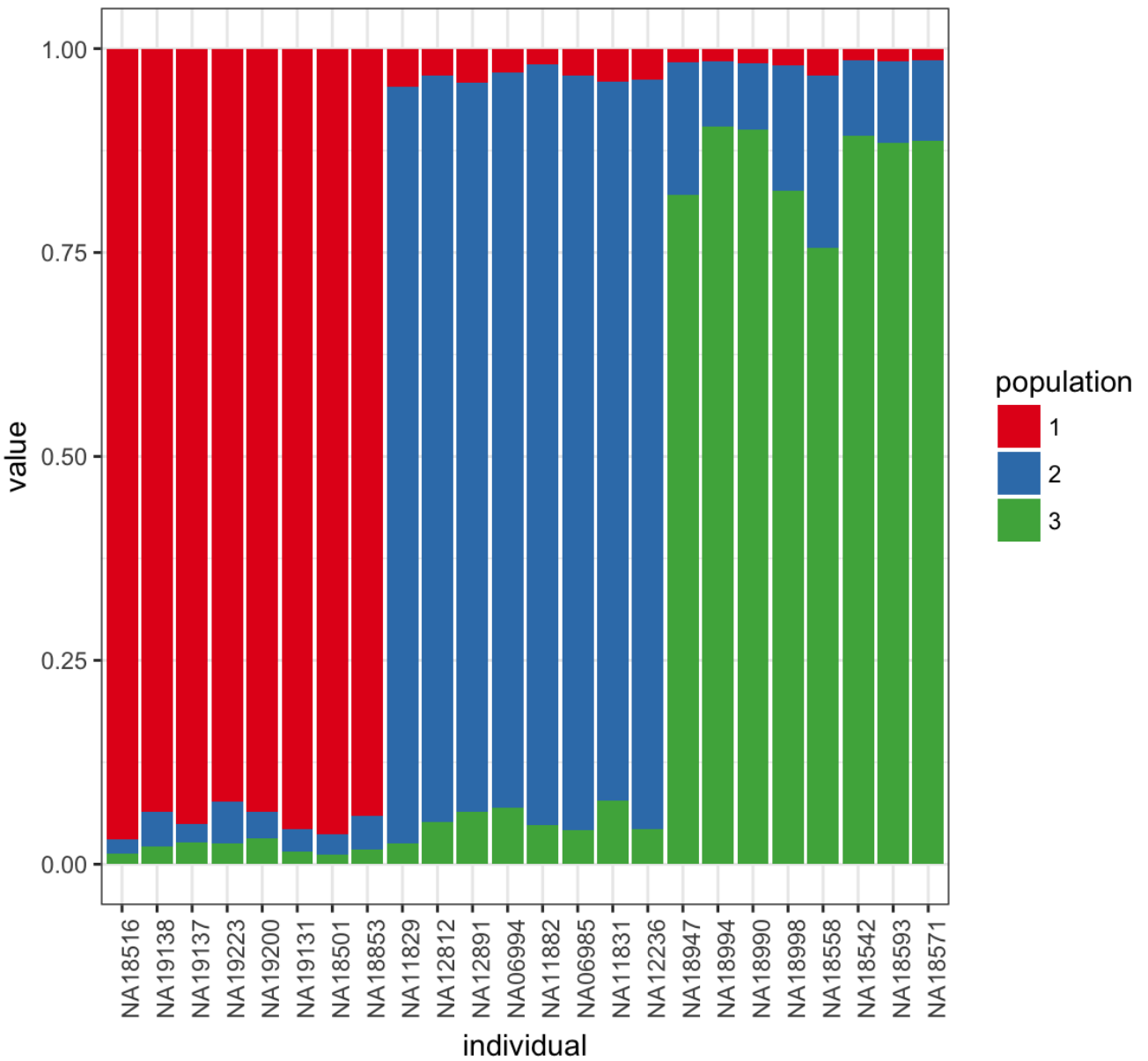
MCMC Configuration

```
> MAX_IT <- 20000
> BURNIN <- 5000
> THIN <- 20
>
> QSUM <- matrix(0, N, K)
>
> START <- 200
> TAIL <- 500
> LL_start <- rep(0, START)
> LL_end <- rep(0, TAIL)
```

Run Sampler

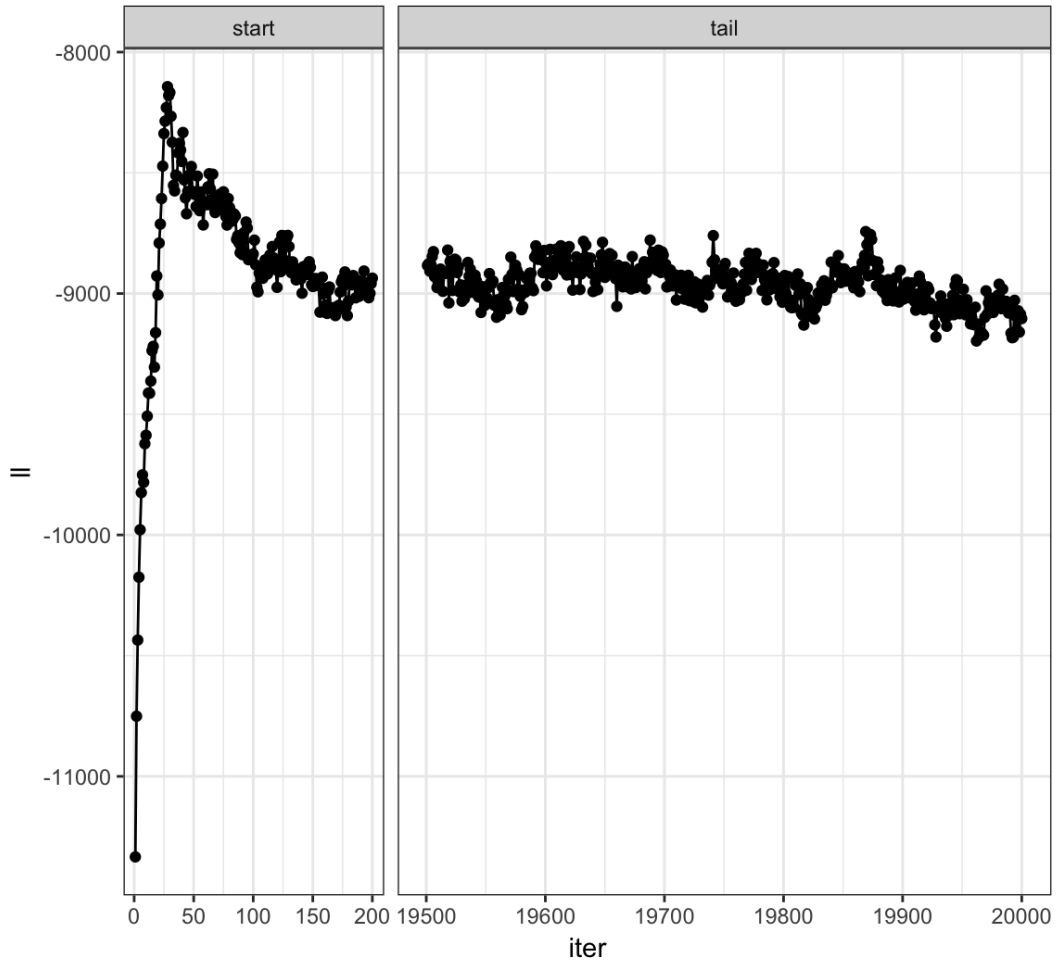
```
> set.seed(1234)
>
> for(it in 1:MAX_IT) {
+   P <- update_P()
+   Q <- update_Q()
+   Za <- update_Z(Xa)
+   Zb <- update_Z(Xb)
+
+   if(it > BURNIN && it %% THIN == 0) {QSUM <- QSUM+Q}
+   if(it <= START) {LL_start[it] <- model_ll()}
+   if(it > MAX_IT-TAIL) {LL_end[it-(MAX_IT-TAIL)] <- model_ll()}
+ }
>
> Q_MEAN <- QSUM/((MAX_IT-BURNIN)/THIN)
```

Posterior Mean of Q



Plot Log-likelihood Steps

Note both the needed burn-in and thinning.



What Happens for K=4?

```

> K <- 4
> Za <- matrix(sample(1:K, L*N, replace=TRUE), L, N)
> Zb <- matrix(sample(1:K, L*N, replace=TRUE), L, N)
> P <- matrix(0, L, K)
> Q <- matrix(0, N, K)
> QSUM <- matrix(0, N, K)

```

Run Sampler Again

```

> for(it in 1:MAX_IT) {
+   P <- update_P()
+   Q <- update_Q()
+   Za <- update_Z(Xa)
+   Zb <- update_Z(Xb)
+
+   if(it > BURNIN && it %% THIN == 0) {
+     QSUM <- QSUM+Q
+   }
+ }

```

```
>  
> Q_MEAN <- QSUM/((MAX_IT-BURNIN)/THIN)
```

Posterior Mean of Q

