

QCB 408 / 508 – Notes on Week 06

Juechun Tang

2020-03-22

Summary

In week 6, we continued our discussion on Bayesian Estimation and classification, different ways of doing priors, numerical methods and likelihood; and the following concepts are discussed:

- Bayesian Estimation
 - Posterior Distribution
 - Posterior Expectation
 - Posterior Intervals
 - Maximum a posteriori Probability
 - Loss Function
 - Bayes Risk
 - Bayes Estimation
- Bayes Classification
 - Posterior Probability
 - loss Function
 - Bayes Risk
 - Bayes Rule
- Priors
 - Conjugate Priors
 - Jeffreys Prior
 - Improper Prior
- Empirical Bayes
- Numerical methods
 - Latent variable models, EM models
 - Markov chain Monte Carlo

Estimation

Bayesian estimation is an analog of point estimation, such as posterior expectation, posterior median, MAP, etc.

Posterior distribution

Assume we have $(X_1, X_2, \dots, X_n) | \theta \stackrel{\text{iid}}{\sim} F_\theta$ with prior distribution $\theta \sim F_\tau$. The posterior distribution of $\theta | \mathbf{X}$ is obtained through Bayes theorem:

$$f(\theta | \mathbf{X}) = \frac{f(\mathbf{X} | \theta) f(\theta)}{f(\mathbf{X})} = \frac{f(\mathbf{X} | \theta) f(\theta)}{\int f(\mathbf{X} | \theta^*) f(\theta^*) d\theta^*} \propto L(\theta; \mathbf{x}) f(\theta)$$

Note \mathbf{X} is bold and capital, representing r.v. Now θ is a random variable, so we write in term of $f(\mathbf{X} | \theta)$. Here we use a condensed notation (not writing τ). $f(\mathbf{X})$ is the marginal pdf or pmf of the data. The denominator is what usually difficult to calculate. This posterior pdf $\propto L(\theta; \mathbf{x}) f(\theta)$ because 1. $f(\mathbf{X} | \theta)$ is the likelihood function and 2. the marginal on the denominator is not a function of θ .

Posterior Expectation

Once we have the posterior pdf or pmf, we can do any probability calculation we want. The the posterior expected value is:

$$E[\theta | \mathbf{x}] = \int \theta f(\theta | \mathbf{x}) d\theta \tag{1}$$

$$= \frac{\int \theta L(\theta; \mathbf{x}) f(\theta) d\theta}{\int L(\theta; \mathbf{x}) f(\theta) d\theta} \tag{2}$$

Equation 2 is got from substituting the Posterior distribution $f(\theta | \mathbf{x})$ in and $f(\mathbf{x})$ cancels out eventually. This is equivalent to the motivating example from last week.

Posterior Interval

A Bayesian analog of the confidence interval: an interval within which the parameter of interest falls within. This is straight probability calculation. Note θ is posterior distribution conditioned on \mathbf{x}

$$1 - \alpha = \Pr(C_\ell \leq \theta \leq C_u | \mathbf{x})$$

Maximum A Posteriori Probability (MAP)

This is maximum likelihood analog. We are maximizing the likelihood times the prior (weighting) over θ .

$$\begin{aligned} \hat{\theta}_{\text{MAP}} &= \operatorname{argmax}_\theta \Pr(\theta | \mathbf{x}) f(\theta) \\ &= \operatorname{argmax}_\theta L(\theta; \mathbf{x}) f(\theta) \end{aligned}$$

Note for many machine learning settings, if there are infinite amount of data $f(\theta)$ peaks to a true value.

Loss functions (Error in estimate)

Let $\mathcal{L}(\theta, \tilde{\theta})$ be the loss function for a given estimator $\tilde{\theta}$. We can have for example squared error loss or absolute error loss:

$$\mathcal{L}(\theta, \tilde{\theta}) = (\theta - \tilde{\theta})^2 \text{ or } \mathcal{L}(\theta, \tilde{\theta}) = |\theta - \tilde{\theta}|.$$

Fields like social sciences and economics commonly use loss functions.

$$\begin{aligned} \mathbb{E} \left[(\theta - \tilde{\theta})^2 \right] &= (\mathbb{E} [\tilde{\theta}] - \theta)^2 + \text{Var} (\tilde{\theta}) \\ &= \text{bias}^2 + \text{variance} \end{aligned}$$

Let's work through how we get the result:

$$\begin{aligned} \mathbb{E} \left[(\theta - \tilde{\theta})^2 \right] &= \mathbb{E} [(\theta^2 - 2\theta\tilde{\theta} + \tilde{\theta}^2)] \\ &= \theta^2 - 2\theta \mathbb{E}[\tilde{\theta}] + \mathbb{E}[\tilde{\theta}^2] \\ &= \mathbb{E}[\tilde{\theta}]^2 + \theta^2 - 2\theta \mathbb{E}[\tilde{\theta}] + \mathbb{E}[\tilde{\theta}^2] - \mathbb{E}[\tilde{\theta}]^2 \\ &= (\mathbb{E} [\tilde{\theta}] - \theta)^2 + \text{Var} (\tilde{\theta}) \end{aligned}$$

Bayes Risk

The **Bayes risk**, $R(\theta, \tilde{\theta})$, is the expected loss with respect to the posterior:

$$\mathbb{E} [\mathcal{L}(\theta, \tilde{\theta}) | \mathbf{x}] = \int \mathcal{L}(\theta, \tilde{\theta}) f(\theta | \mathbf{x}) d\theta$$

The intergral is taken over the posterior distribution $\theta | \mathbf{x}$. What we have here is that we first conditioned on observed data \mathbf{x} , and then do the probability calculation. Note both θ and $\tilde{\theta}$ are conditioned on data.

Bayes Estimators

Once we have Bayes risk from loss function, we can calculate Bayes Estimators to minimize the Bayes risk. Note everything is in finite sample size, n is not going to ∞ . The challenge is what the prior distribution is, which we will discuss more in details.

The posterior expectation $\mathbb{E}[\theta | \mathbf{x}]$ minimizes the Bayes risk of $\mathcal{L}(\theta, \tilde{\theta}) = (\theta - \tilde{\theta})^2$ (squared loss).

The posterior median of $f(\theta | \mathbf{x})$, calculated by $F_{\theta | \mathbf{x}}^{-1}(1/2)$, minimizes the Bayes risk of $\mathcal{L}(\theta, \tilde{\theta}) = |\theta - \tilde{\theta}|$ (absolute loss).

Bayes Classification

Bayes classification is an analog of hypothesis testing. The **Neyman-Pearson Lemma** and **Generalized LRT** we covered before have a direct anaalog in bayesian inference. Typically we deal with composite hypothesis where ther are more than one value in each hypothesis.

We will set up a composite vs. composite in the following section.

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 $\theta \in \Theta_1$.

The null hypothesis is $\theta \in \Theta_0$ and the alternative hypotehsis is $\theta \in \Theta_1$.

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

A criticism for the Bayesian is that they often think parameters as a interval instead of a single value. From the Null hypothesis significance testing (NHST), it's not reasonable to have a simple null hypothesis. One exception that's totally scientific is genetics mapping where we consider if the recombination is 1/2 or less than 1/2.

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

Loss Function

We can only make 2 types of error for classification: 0 when it's 1 or 1 when it's 0.

Let \tilde{H} be estimate for H , $\mathcal{L}(\tilde{H}, H)$ 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 $c_I, c_{II} > 0$, where c_I indicates Type I error and c_{II} indicates Type II error.

Bayes Risk

Now that we get loss function from posterior and prior, we will calculate Bayes Risk as we did for estimation, and then Bayes Estimate which minimize the Bayes Risk.

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

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

Note that the equation is interpreted as the sum of Type I penalty times Type I error rate times the prior plus the Type II penalty times Type II error rate times the prior.

The challenge is to determine what c_I , c_{II} and prior are. If there are trillions of data points, we can get values for them.

Bayes Rule (Bayes Estimator)

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.

Priors

The challenge to get the Posterior distribution is how we get the prior distribution. Sometimes the prior is chosen for mathematical convenience, other types are estimated using data.

Conjugate Priors

A **conjugate prior** is a prior distribution for a data generating distribution so that the posterior distribution is of the same type as the prior.

Conjugate priors are useful for obtaining straightforward calculations of the posterior. If we can find a conjugate prior, we can avoid calculating marginal on the denominator.

There is a systematic method for calculating conjugate priors for exponential family distributions.

Example: Beta-Bernoulli

Suppose $\mathbf{X}|p \stackrel{\text{iid}}{\sim} \text{Bernoulli}(p)$ and suppose that the prior $p \sim \text{Beta}(\alpha, \beta)$. The posterior is:

$$\begin{aligned} f(p|\mathbf{x}) &\propto L(p;\mathbf{x})f(p) \\ &= p^{\sum x_i} (1-p)^{\sum (1-x_i)} p^{\alpha-1} (1-p)^{\beta-1} \\ &= p^{\alpha-1+\sum x_i} (1-p)^{\beta-1+\sum (1-x_i)} \\ &\propto \text{Beta}(\alpha + \sum x_i, \beta + \sum (1-x_i)) \end{aligned}$$

Therefore,

$$\mathbb{E}[p|\mathbf{x}] = \frac{\alpha + \sum x_i}{\alpha + \beta + n}.$$

If there's a true p , note that as $n \rightarrow \infty$, p concentrates to true p (Strong Law of large numbers).

Example: Normal-Normal

Suppose $\mathbf{X}|\mu \stackrel{\text{iid}}{\sim} \text{Normal}(\mu, \sigma^2)$, where σ^2 is known, and suppose that $\mu \sim \text{Normal}(a, b^2)$.

Then it can be shown that $\mu|\mathbf{x} \sim \text{Normal}(\mathbb{E}[\mu|\mathbf{x}], \text{Var}(\mu|\mathbf{x}))$

$$\mathbb{E}[\mu|\mathbf{x}] = \frac{b^2}{\frac{\sigma^2}{n} + b^2} \bar{x} + \frac{\frac{\sigma^2}{n}}{\frac{\sigma^2}{n} + b^2} a$$

$$\text{Var}(\mu|\mathbf{x}) = \frac{b^2 \frac{\sigma^2}{n}}{\frac{\sigma^2}{n} + b^2}$$

Let's work through:

$$\begin{aligned}
f(\mu|\mathbf{x}) &\propto L(\mu; \mathbf{x})f(\mu) \\
&= e^{-\frac{\sum (x_i - \mu)^2}{2\sigma^2}} e^{-\frac{(\mu - a)^2}{2b^2}} \\
&= e^{-\frac{\sum x_i^2 - n\mu^2 + 2\mu n\bar{x}}{2\sigma^2} - \frac{\mu^2 + a^2 - 2\mu a}{2b^2}} \\
&\propto e^{-\frac{n\mu^2 + 2\mu n\bar{x}}{2\sigma^2} - \frac{\mu^2 - 2\mu a}{2b^2}} \\
&= e^{-\frac{\mu^2}{2}(\frac{n}{\sigma^2} + \frac{1}{b^2}) + \mu(\frac{n\bar{x}}{\sigma^2} + \frac{a}{b^2})}
\end{aligned}$$

Let's try to get a quadratic function in terms of μ :

Let $\lambda^2 = (\frac{n}{\sigma^2} + \frac{1}{b^2})^{-1}$, $\frac{\lambda'}{\lambda^2} = \frac{n\bar{x}}{\sigma^2} + \frac{a}{b^2}$, We have:

$$\begin{aligned}
f(\mu|\mathbf{x}) &= e^{-\frac{\mu^2}{2\lambda^2} + \frac{\lambda'}{\lambda^2}\mu} \\
&\propto e^{-\frac{\mu^2}{2\lambda^2} + \frac{\lambda'}{\lambda^2}\mu + \frac{\lambda'^2}{2\lambda^2}} \\
&= e^{-\frac{(\mu - \lambda')^2}{2\lambda^2}} \\
&\propto \text{Normal}(\lambda', \lambda^2)
\end{aligned}$$

where

$$\begin{aligned}
\lambda' &= (\frac{n\bar{x}}{\sigma^2} + \frac{a}{b^2}) (\frac{b^2 \frac{\sigma^2}{n}}{\frac{\sigma^2}{n} + b^2}) = \frac{b^2}{\frac{\sigma^2}{n} + b^2} \bar{x} + \frac{\frac{\sigma^2}{n}}{\frac{\sigma^2}{n} + b^2} a \\
\lambda^2 &= \frac{\sigma^2 b^2}{nb^2 + \sigma^2} = \frac{b^2 \frac{\sigma^2}{n}}{\frac{\sigma^2}{n} + b^2}
\end{aligned}$$

Note that the posterior mean is a weighted average of the sample mean and the prior a . It's weighted by the weighted averaged of the two variances. Again, as n goes to inf, if I'm a frequentist, then we will get the true value (variance goes to 0). Note the conjugate prior only tells which distribution to pick, but we still need to specify what the parameter values are.

Jeffreys Prior

Another way to pick prior for mathematical convinience is through **Jeffreys Prior**. After reparameterization, we get the same answer when doing inference.

If we do inference based on prior $\theta \sim F_\tau$ to obtain $f(\theta|\mathbf{x}) \propto L(\theta; \mathbf{x})f(\theta)$, it follows that this inference may *not* be invariant to transformations of θ , such as $\eta = g(\theta)$.

If we utilize a **Jeffreys prior**, which means it is such that

$$f(\theta) \propto \sqrt{I(\theta)}$$

then the prior will be invariant to transformations of θ . We would want to show that $f(\theta) \propto \sqrt{I(\theta)}$ implies $f(\eta) \propto \sqrt{I(\eta)}$.

Examples: Jeffreys Priors

Normal(μ, σ^2), σ^2 known: $f(\mu) \propto 1$. Note $f(\mu) \propto 1$ a problem because this is not a proper probability distribution

Normal(μ, σ^2), μ known: $f(\sigma) \propto \frac{1}{\sigma}$

Poisson(λ): $f(\lambda) \propto \frac{1}{\sqrt{\lambda}}$

Bernoulli(p): $f(p) \propto \frac{1}{\sqrt{p(1-p)}}$

Improper Prior

An **improper prior** is a prior such that $\int f(\theta)d\theta = \infty$. Therefore, it is not a proper probability distribution, but we can view it as a weighting function and the resulting posterior $f(\theta|\mathbf{x}) \propto L(\theta;\mathbf{x})f(\theta)$ yields a probability distribution.

Take for example the case where $\mathbf{X}|\mu \stackrel{\text{iid}}{\sim} \text{Normal}(\mu, \sigma^2)$, where σ^2 is known, and suppose that $f(\mu) \propto 1$. Then $\int f(\theta)d\theta = \infty$, but

$$f(\theta|\mathbf{x}) \propto L(\theta;\mathbf{x})f(\theta) \sim \text{Normal}(\bar{x}, \sigma^2/n)$$

which is a proper probability distribution.

Empirical Bayes

Empirical Bayes uses the data to estimate the **prior parameter**. A lot of ML methods are **Empirical Bayes** in terms of still treating parameter as a random variable, but everything is evaluated from data. It is commonly used to estimate multiple parameters.

This is especially useful for high-dimensional data when many parameters are simultaneously drawn from a prior with multiple observations drawn per parameter realization.

There are both parametric and nonparametric empirical bayes. Nonparametric means doing inference in a way that we are making minimal assumptions of the distribution. And we focus on parametric bayes here.

Approach

The usual approach is to integrate out the parameter to obtain

$$f(\mathbf{x};\tau) = \int f(\mathbf{x}|\theta)f(\theta;\tau)d\theta.$$

Now we have pdf or pmf my data in terms of prior parameters. An estimation method (such as MLE) is then applied to estimate τ . Then inference proceeds as usual under the assumption that $\theta \sim f(\theta;\hat{\tau})$.

Problem is that we tend to overfit the data. In ML, what people do in practice is to keep some test data, perform cross-validation to avoid overfitting.

MLE tends to fail when doing high-dimensional data. Some empirical bayes takes overfitting into account and gets much better results.

Example: Normal

Suppose that $X_i|\mu_i \sim \text{Normal}(\mu_i, 1)$ for $i = 1, 2, \dots, n$ where these rv's are independent. Also suppose that $\mu_i \stackrel{\text{iid}}{\sim} \text{Normal}(a, b^2)$.

Note here we have a different μ for every observation. Let this be a toy model for gene expression where each X_i is the sample mean from gene i , and we model this in aggregates. Note it's sufficient for μ_i per sufficient statistics. Now we integrate out μ_i

$$f(x_i; a, b) = \int f(x_i | \mu_i) f(\mu_i; a, b) d\mu_i \sim \text{Normal}(a, 1 + b^2).$$

My MLE are:

$$\implies \hat{a} = \bar{x}, \quad 1 + \hat{b}^2 = \frac{\sum_{k=1}^n (x_k - \bar{x})^2}{n}$$

If the $1 + \hat{b}^2 < 1$, this term gets truncated (Restricted).

From **conjugate prior**, we have:

$$E[\mu_i | x_i] = \frac{1}{1 + b^2} a + \frac{b^2}{1 + b^2} x_i \implies$$

$$\begin{aligned} \hat{E}[\mu_i | x_i] &= \frac{1}{1 + \hat{b}^2} \hat{a} + \frac{\hat{b}^2}{1 + \hat{b}^2} x_i \\ &= \frac{n}{\sum_{k=1}^n (x_k - \bar{x})^2} \bar{x} + \left(1 - \frac{n}{\sum_{k=1}^n (x_k - \bar{x})^2}\right) x_i \end{aligned}$$

This is called shrinkage estimator. We just take sample mean for gene 1 as μ_1 and so on. For empirical Bayes, we take weighted average with the overall pulled sample mean. Sum of squared error loss and empirical Bayes are uniformly better than the MLE even if the truncated version.

Numerical Methods

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

Some numerical approach for likelihood based inference:

- Expectation-maximization (EM) algorithm
- Variational inference: an extension of the EM algorithm
- Markov chain Monte Carlo (MCMC)
 - Metropolis sampling
 - Metropolis-Hastings sampling
 - Gibbs sampling

In class, we mainly discussed EM and MCMC.

Latent Variable Models

One way think about EM algorithm is through latent variable models.

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.

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

Let $X_1, X_2, \dots, X_n \stackrel{\text{iid}}{\sim} F_{\theta}$ be a sample of n independent observations from a mixture of k normal distributions with pdf

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

where π_k represents the probability of which component from the mixture is drawn from and $\sum \pi_k = 1$. In the simple case of 2 Gaussian mixture, $\pi_1 + \pi_2 = 1$. Note this is a messy function which is not normal.

The goal is to estimate the unknown parameters $\theta = (\pi_1, \dots, \pi_K, \mu_1, \dots, \mu_K, \sigma_1^2, \dots, \sigma_K^2)$. In this case, the MLEs of the unknown parameters cannot be found analytically. Mixture models allow us to estimate the parameters.

Let $\mathbf{Z}_1, \mathbf{Z}_2, \dots, \mathbf{Z}_n \stackrel{\text{iid}}{\sim} \text{Multinomial}_K(1, \boldsymbol{\pi})$ be the latent variables that determine the component k from which the observation is drawn from, where $\boldsymbol{\pi} = (\pi_1, \dots, \pi_K)$. Note that $Z_{ik} \in \{0, 1\}$ and $\sum_{k=1}^K Z_{ik} = 1$. Note each \mathbf{Z}_i is a vector of length K .

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

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

from which we can form an analytic solution.

Note that

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

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

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

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 pmf

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

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 pmf is

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

EM Algorithm

We start from an intractable pdf or pmf. After we introduce a latent variable, the likelihood becomes tractable and then we can derive MLE or MAP. If we have clinical studies and have missing data, this framework works well.

Conceptually, what we basically do is that suppose we have some data points that may come from a mixture of known distributions, we randomly assign parameters values and do soft-clustering, that is, calculate the probability of which component each data point come from. Then the parameters are optimized and iterate until convergence.

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.

To do EM, we first need to know the distribution, probably from EDA to justify the validity of the probabilistic model selected.

The Algorithm

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

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

where $\log f(\mathbf{x}, \mathbf{Z}; \boldsymbol{\theta})$ is the complete data log likelihood. It's a function of Z which we replace with $\mathbb{E}_{\mathbf{Z}|\mathbf{X}=\mathbf{x}}$ given current estimate of $\boldsymbol{\theta}^{(t)}$. As we integrate out Z , we can then do MLE or MAP.

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

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

where $f(\boldsymbol{\theta})$ is the prior distribution on $\boldsymbol{\theta}$ as we shown before in the bayesian framework.

EM Examples

Normal Mixture Model

Returning to the Normal mixture model, we first calculate the log likelihood of the complete data:

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

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

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 caculating

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

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

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

We got the result from Bayes rule the the numerator is the prior times the normal pdf and the denominator is the mariginal $f(x)$.

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 paramaters 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 $\boldsymbol{\theta}^{(t)}$. Note μ_k, σ_k^2, π_k are all free parameters.

M-Step

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

To get $\pi_k^{(t+1)}$, we use a Lagrange multiplier. Let $\sum_{i=1}^n \hat{z}_{ik}^{(t)} = a_k$, we have

$$\max \sum_{k=1}^K a_k \log \pi_k$$

s.t. $\pi_k \geq 0, \sum \pi_k = 1$

$$\max_{\pi_k} \min_{\lambda} \sum_{k=1}^K a_k \log \pi_k + \lambda (\sum (\pi_k - 1)) \leq \max_{\pi_k} \sum_{k=1}^K a_k \log \pi_k + \lambda^* (\sum (\pi_k - 1)) = f(\pi)$$

where λ^* is a random λ

$$\begin{aligned} \frac{\partial f(\pi)}{\partial \pi_k} &= \frac{a_k}{\pi_k} + \lambda^* = 0 \quad \forall k \\ \sum \pi_k &= \sum -\frac{a_k}{\lambda^*} = 1 \end{aligned}$$

we get $\lambda^* = -\sum a_k$, $\pi_k = \frac{a_k}{\sum a_k} \rightarrow \pi_k^{(t+1)} = \frac{\sum_{i=1}^n \hat{z}_{ik}^{(t)}}{n}$.

Similarly, we can get;

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

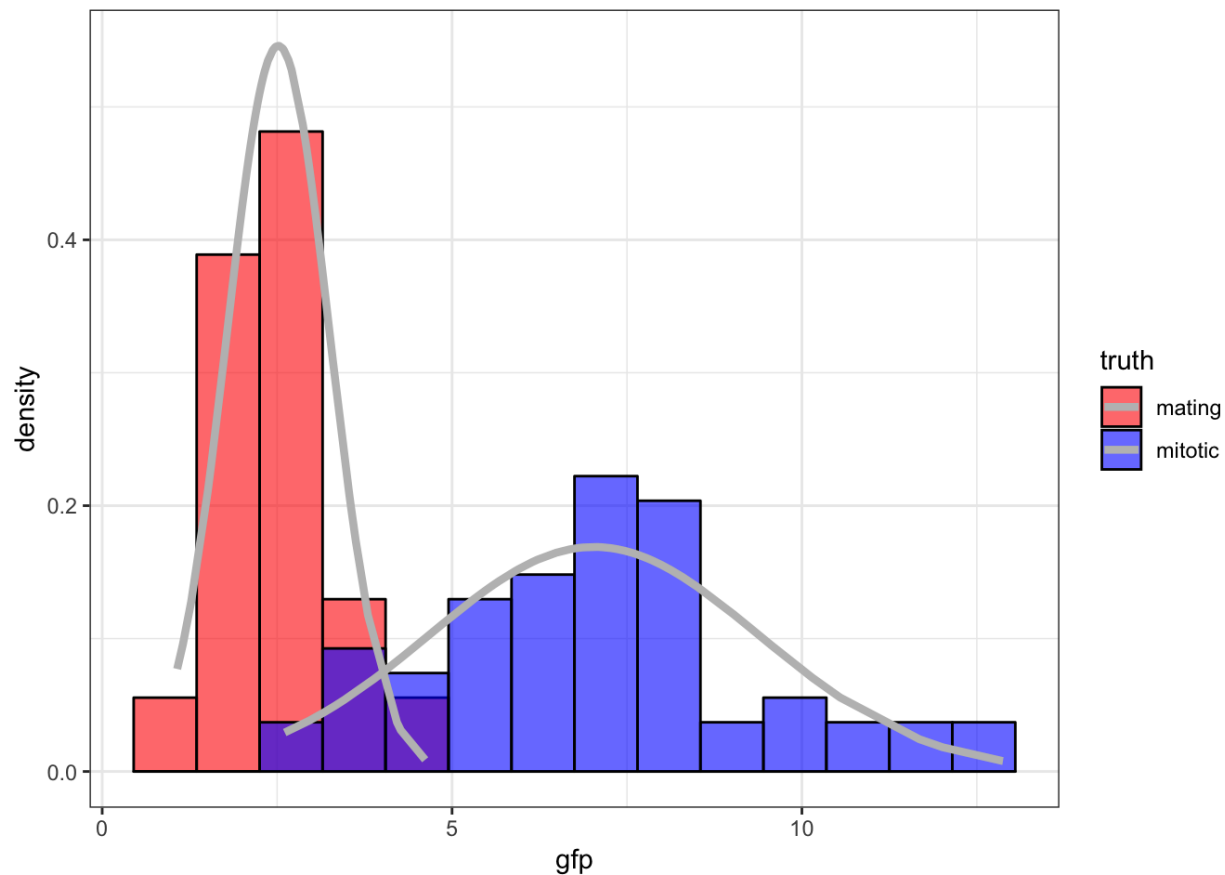
Caveat

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. We will run an EM and then compare with the the true distribution.



Initialize Values

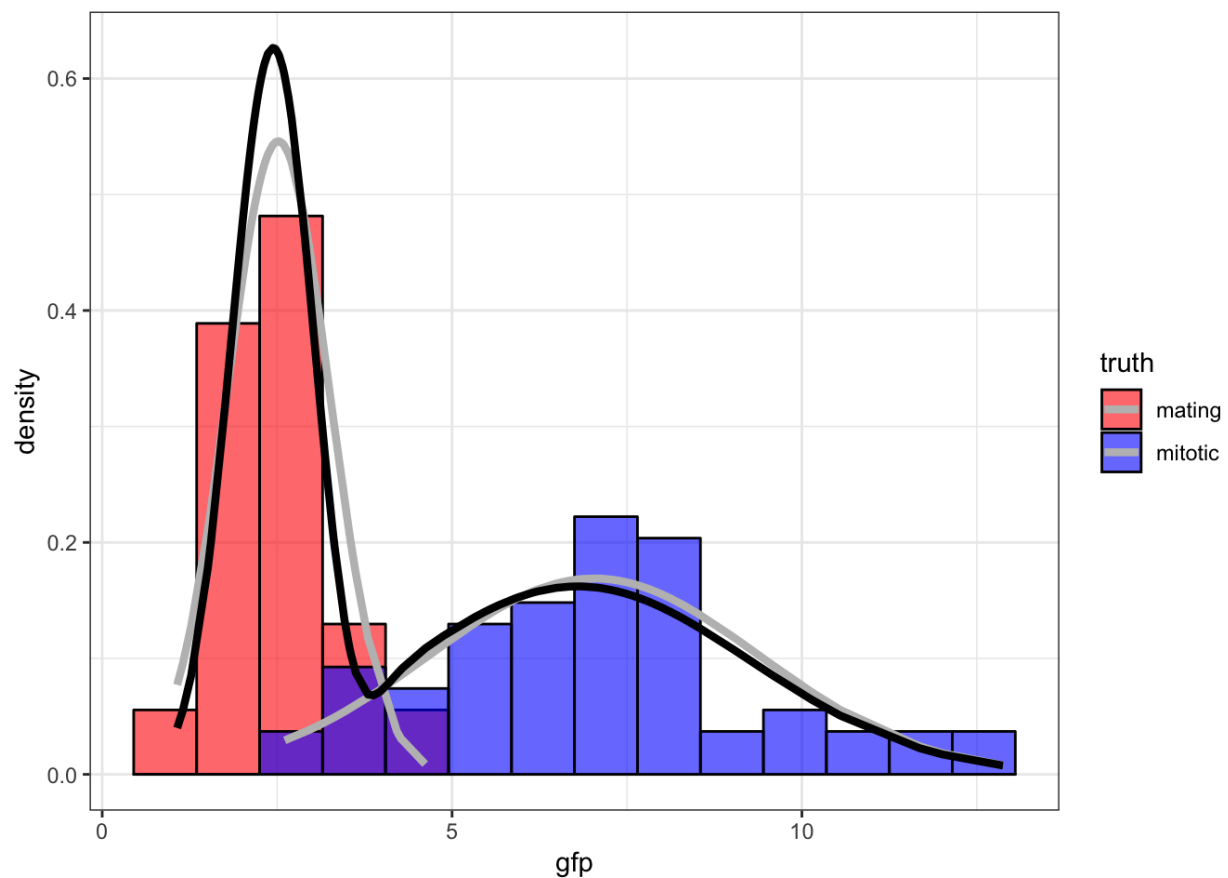
```
> set.seed(508)
> B <- 100 #number of iteration, although need to check if converge at 100
> p <- rep(0,B) #pi vector
> x <- rep(0,B) #added just to compile, don't have the data locally
> mu1 <- rep(0,B)
> mu2 <- rep(0,B)
> s1 <- rep(0,B)
> s2 <- rep(0,B)
> p[1] <- runif(1, min=0.1, max=0.9)
> mu.start <- sample(x, size=2, replace=FALSE)
> #x is a vector of data, draw 2 data potin randomly and set as initial mu
> mu1[1] <- min(mu.start)
> mu2[1] <- max(mu.start)
> s1[1] <- var(sort(x)[1:60]) #bottom half from component 1
> s2[1] <- var(sort(x)[61:120])
> z <- rep(0,120) #indication variable whether it's from population 2 or 1
```

A smarter way of setting of mean can be done as the same as variance

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]))))
+   # z is indicator of whether it is from component 2.
+   # It's 1 if from 2 0 if from 1
+   #denominator is the marginal distribution of x
+   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,]  NA NA  NA NA NA
[99,]  NA NA  NA NA NA
[100,] NA NA  NA NA NA
> #see if convergence
```

Fitted Mixture Distribution



Grey line is when we know the truth (complete data likelihood estimates); Black line is when we don't know the truth (observed likelihood estimates). EM performed really well in this case.

Other Applications of EM

- Dealing with missing data:
- Multiple imputation of missing data: what missing data looks like under repeated sampling
- Truncated observations
- Bayesian hyperparameter estimation: when have a Bayesian model, treat unobserved variables as latent variable (Empirical Bayes)
- Hidden Markov models: model data with some underlying state e.g. a Hidden Markov state: gene in/out of open reading frame (ORF).

EM Increases Likelihood

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

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

which is a non-decreasing iteration

Also, by the properties of KL divergence (broader way to do numerical method) stated above (more details in Ch 43 in the reading), we have

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

Putting these together we have

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

Markov Chain Monte Carlo

A very popular technique in Bayesian inference, but it's very slow and doesn't scale well to high-dimensional datasets; The tuning steps are very tricky and need a lot of experience.

Motivation

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

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

compared with

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

Say we have a plot of the

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

There will be areas that are very high and very low. High areas are with high posterior probability. In order to take the product and turn into a probability distribution, we need the denominator, which is often hard to calculate, that's where MCMC comes in.

MCMC tries to spend more time in high probability areas. If we have two θ values and I take the ratio of the two pdfs, now the marginal disappears. Markov chain Monte Carlo is a method for simulating data approximately from $f(\theta|\mathbf{x})$ with knowledge of only $L(\theta;\mathbf{x})f(\theta)$.

But once we get into high-dimensions, say the gaussian mixture model with multiple parameters, it would be hard to calculate by MCMC.

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 (a sequence of R.V. that the current variable only depends on the previous one) of θ values so that, in some asymptotic sense, these are equivalent to iid draws from $f(\theta|\mathbf{x})$.

The draws are done competitively so that the next draw of a realization of θ depends on the current value.

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

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

Metropolis-Hastings Algorithm

1. Initialize $\theta^{(0)}$
2. Generate $\theta^* \sim q(\theta|\theta^{(b)})$ for some pdf or pmf $q(\cdot|\cdot)$ which we have to choose the distribution
3. With probability

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

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

This is taking the ratio of the posterior pdf of new θ^* vs. old $\theta^{(b)}$, and weigh that by $\frac{q(\theta^{(*)}|\theta^{(b)})}{q(\theta^{(b)}|\theta^{(*)})}$. The q function here is to introduce some stochastic behavior so it's not deterministic and stuck at local maxima.

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

Metropolis Algorithm

It's a simplification of Metropolis-Hastings Algorithm. The Metropolis algorithm restricts $q(\cdot, \cdot)$ to be symmetric so that $q(\theta^{(b)}|\theta^*) = q(\theta^*|\theta^{(b)})$ and

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

Utilizing MCMC Output

Two common uses of the output from MCMC are as follows: 1. The posterior expectation $E[f(\boldsymbol{\theta})|\mathbf{x}]$ 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 $\boldsymbol{\theta}$ with some μ and σ^2 .
- 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.
 - There’s a burn in time to reach stable state.
- The generated $\boldsymbol{\theta}^{(b)}$ are typically “thinned” (only sampled every so often) to reduce Markov dependence.
 - If we want to approximate iid draws from the posterior, since $\boldsymbol{\theta}^{(b_2)}|\boldsymbol{\theta}^{(b_1)}, \boldsymbol{\theta}^{(b_3)}|\boldsymbol{\theta}^{(b_2)}, \dots$, sometimes we need to do thinning, that is taking, say every 100th, $\boldsymbol{\theta}$.

Session Information

```
> sessionInfo()
R version 3.6.0 (2019-04-26)
Platform: x86_64-apple-darwin15.6.0 (64-bit)
Running under: macOS 10.15.1

Matrix products: default
BLAS: /Library/Frameworks/R.framework/Versions/3.6/Resources/lib/libRblas.0.dylib
LAPACK: /Library/Frameworks/R.framework/Versions/3.6/Resources/lib/libRlapack.dylib

locale:
[1] en_US.UTF-8/en_US.UTF-8/en_US.UTF-8/C/en_US.UTF-8/en_US.UTF-8

attached base packages:
[1] stats      graphics  grDevices  utils      datasets  methods   base

other attached packages:
[1] gapminder_0.3.0 forcats_0.4.0 stringr_1.4.0 dplyr_0.8.4
[5] purrr_0.3.3     readr_1.3.1  tidyr_1.0.2  tibble_2.1.3
[9] ggplot2_3.2.1   tidyverse_1.3.0 knitr_1.28
```

loaded via a `namespace` (and not attached):

```
[1] tidyselect_1.0.0 xfun_0.12      haven_2.2.0      lattice_0.20-38
[5] colorspace_1.4-1 vctrs_0.2.2      generics_0.0.2   htmltools_0.4.0
[9] yaml_2.2.1       rlang_0.4.4      pillar_1.4.3     withr_2.1.2
[13] glue_1.3.1       DBI_1.1.0        dbplyr_1.4.2     modelr_0.1.5
[17] readxl_1.3.1     lifecycle_0.1.0  munsell_0.5.0    gtable_0.3.0
[21] cellranger_1.1.0 rvest_0.3.5      evaluate_0.14    fansi_0.4.0
[25] broom_0.5.4      Rcpp_1.0.3       scales_1.1.0     backports_1.1.5
[29] jsonlite_1.6.1   fs_1.3.1         hms_0.5.3        digest_0.6.23
[33] stringi_1.4.3    grid_3.6.0       cli_2.0.0        tools_3.6.0
[37] magrittr_1.5     lazyeval_0.2.2   crayon_1.3.4     pkgconfig_2.0.3
[41] xml2_1.2.2       reprex_0.3.0     lubridate_1.7.4  assertthat_0.2.1
[45] rmarkdown_2.1    httr_1.4.1       rstudioapi_0.11  R6_2.4.1
[49] nlme_3.1-143     compiler_3.6.0
```