

## Prob 1. Derive the formulae on p.31.

Gaussian Mixture Model EM

$\theta = \{\mu_1, \dots, \mu_m, \sigma_1, \dots, \sigma_m, \alpha_1, \dots, \alpha_m\}$ , and so from the first section of this note, our likelihood is:

$$L(\theta | X_1, \dots, X_n) = \prod_{i=1}^n \sum_{k=1}^m \alpha_k N(x_i; \mu_k, \sigma_k^2)$$

So our log-likelihood is:

$$\ell(\theta) = \sum_{i=1}^n \log \left( \sum_{k=1}^m \alpha_k N(x_i; \mu_k, \sigma_k^2) \right) \quad (1)$$

Take initial guesses for the parameters  $\hat{\Theta}$

*E* step: compute the posterior distribution of  $z_i$  given the observations

$$\Lambda_{i,k} = P_{\hat{\theta}}(z_i = k | x_i) = \frac{\hat{\alpha}_k P_{\hat{\theta}_k}(x_i)}{P_{\hat{\theta}}(x_i)} = \frac{\hat{\alpha}_k P_{\hat{\theta}_k}(x_i)}{\sum_{j=1}^m \hat{\alpha}_j P_{\hat{\theta}_j}(x_i)} \quad (2)$$

$$Q(\Theta | \hat{\theta}) = \sum_{i=1}^n \sum_{k=1}^m \Lambda_{i,k} \log(\alpha_k P_{\theta_k}(x_i))$$

*M* step: update  $\alpha_k, \mu_k, \sigma_k$  for each Gaussian  $k = 1, \dots, m$

$$\hat{\alpha}_k = \frac{1}{n} \sum_{i=1}^n \Lambda_{i,k}, \hat{\mu}_k = \frac{\sum_{i=1}^n x_i \Lambda_{i,k}}{\sum_{i=1}^n \Lambda_{i,k}}, \hat{\sigma}_k^2 = \frac{\sum_{i=1}^n \Lambda_{i,k} \|x_i - \hat{\mu}_k\|^2}{\sum_{i=1}^n \Lambda_{i,k}}$$

### Proof

Let the partial derivative of (1) about  $\mu_k$  equal to

$$\sum_{i=1}^n \left\{ \frac{1}{\sum_{k=1}^m \alpha_k N(x_i; \mu_k, \sigma_k)} \alpha_k N(x_i; \mu_k, \sigma_k) \frac{(x_i - \mu_k)}{\sigma_k^2} \right\} = 0 \quad (3)$$

by (2) and (3)

$$\begin{aligned} \sum_{i=1}^n \Lambda_{i,k} \frac{(x_i - \mu_k)}{\sigma_k^2} &= 0 \\ \Rightarrow \hat{\mu}_k &= \frac{\sum_{i=1}^n x_i \Lambda_{i,k}}{\sum_{i=1}^n \Lambda_{i,k}} \end{aligned}$$

Similarly, let the partial derivative of (1) about  $\sigma_k^2$ , we have

$$\sum_{i=1}^n \left\{ \frac{1}{\sum_{k=1}^m \alpha_k N(x_i; \mu_k, \sigma_k)} \alpha_k N(x_i; \mu_k, \sigma_k) \left( \frac{(x_i - \mu_k)^2}{2\sigma_k^4} - \frac{1}{2\sigma_k^2} \right) \right\} = 0 \quad (4)$$

by (2) and (4)

$$\hat{\sigma}_k^2 = \frac{\sum_{i=1}^n \Lambda_{i,k} \|x_i - \hat{\mu}_k\|^2}{\sum_{i=1}^n \Lambda_{i,k}}$$

According to the definition of  $\Lambda_{i,k}$ , the estimate of  $\alpha_k$  is obviously that

$$\hat{\alpha}_k = \frac{1}{n} \sum_{i=1}^n \Lambda_{i,k}$$

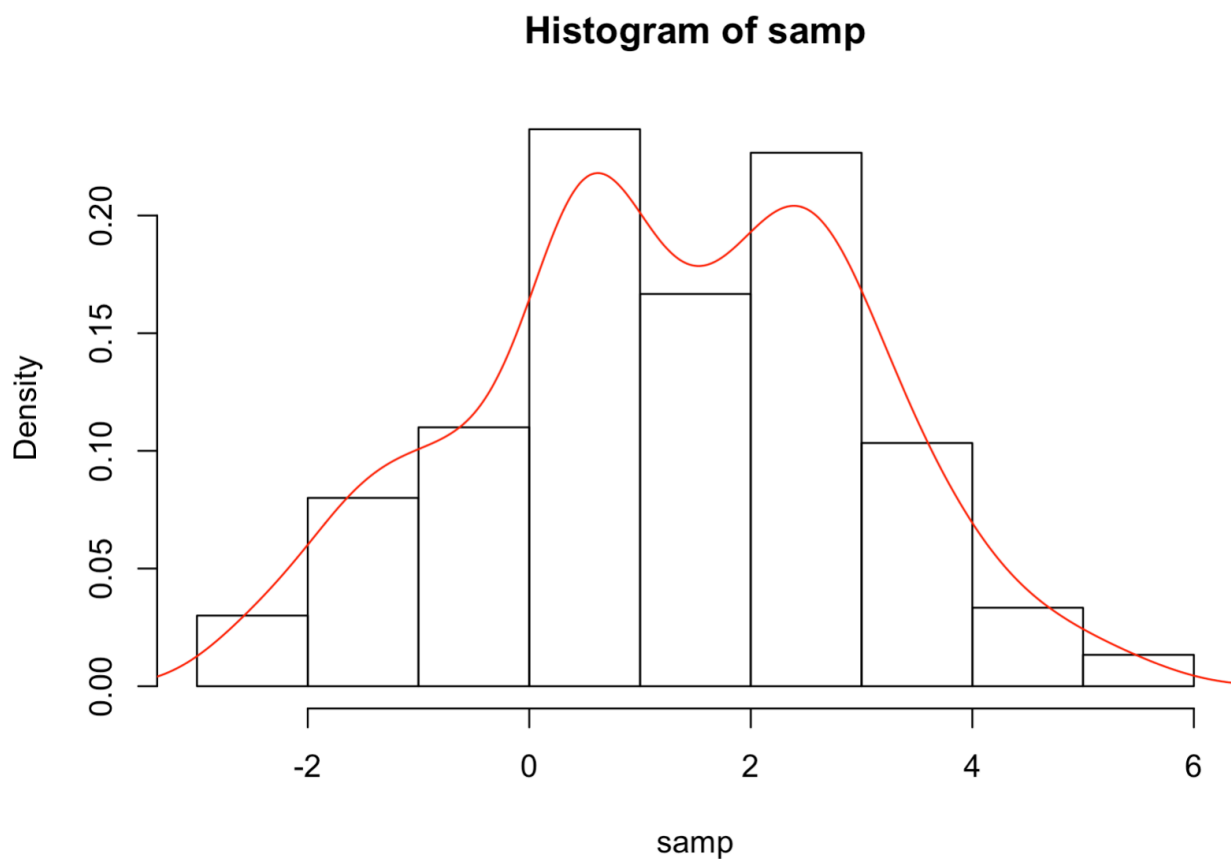
Q.E.D

**Prob 2 . For  $k = 3$ , simulate 100 data points from such a GMM, and use the EM algorithm to estimate the parameters. Check how well the estimates are.**

```

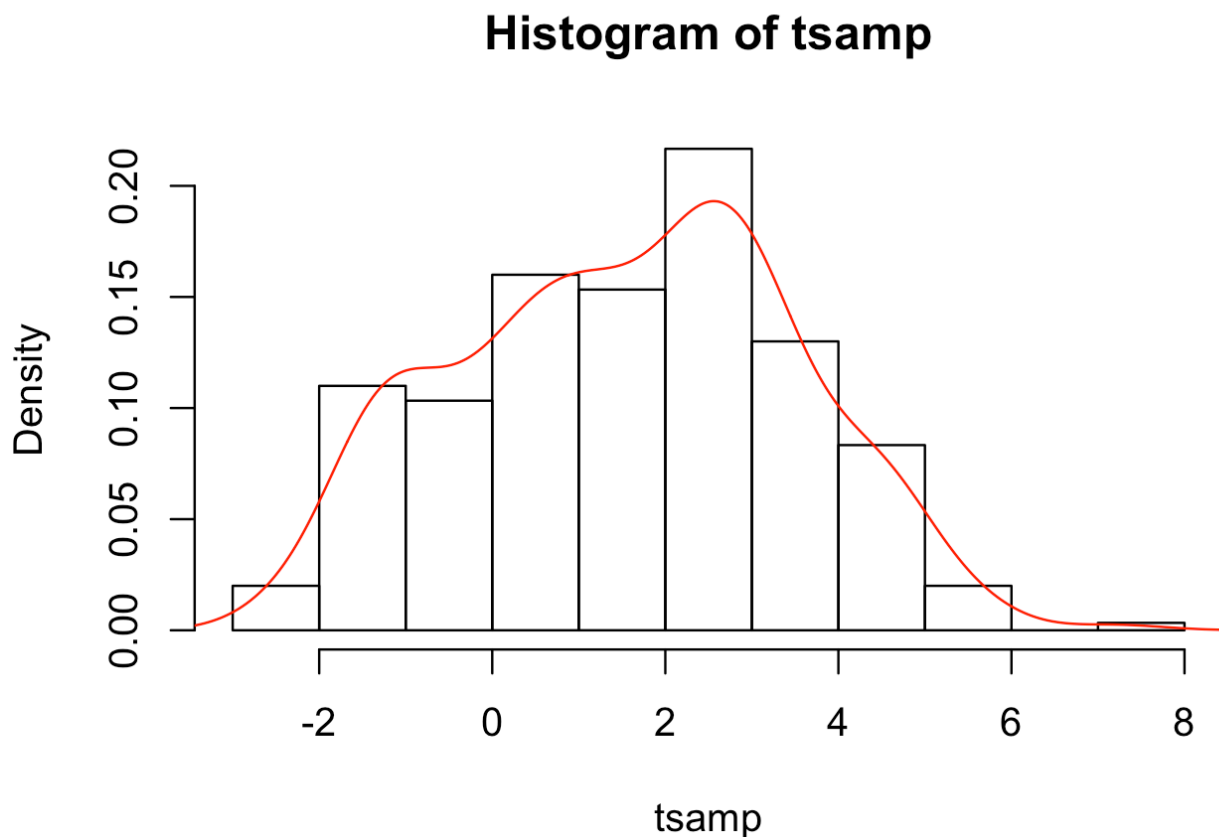
1  set.seed(1234)
2  n <- 300
3  alpha1 <- 0.4
4  mu1 <- 3
5  sigma1 <- 1
6  alpha2 <- 0.4
7  mu2 <- -2
8  sigma2 <- 2
9  alpha3 <- 0.2
10 mu3 <- 0
11 sigma3 <- 2^0.5
12 n1 <- floor(n*alpha1)
13 n2 <- floor(n*(alpha2+alpha1))-n1
14 n3 <- n-n2-n1
15 samp <- numeric(n)
16 samp[1:n1] <- rnorm(n1, mu1, sigma1)
17 samp[(n1+1):n1+n2] <- rnorm(n2, mu2, sigma2)
18 samp[(n2+1):n] <- rnorm(n3, mu3, sigma3)
19 hist(samp, freq = FALSE)
20 lines(density(samp), col = 'red')
21 #EM
22 library(mclust)
23 em <- Mclust(samp)
24 summary(em, parameters = T)

```



```
1 Gaussian finite mixture model fitted by EM algorithm
2 -----
3
4 Mclust V (univariate, unequal variance) model with 3 components:
5
6 log-likelihood   n df      BIC      ICL
7      -566.8767 300  8 -1179.384 -1242.349
8
9 Clustering table:
10    1  2  3
11   63 52 185
12
13 Mixing probabilities:
14           1           2           3
15  0.2089691 0.1350616 0.6559693
16
17 Means:
18           1           2           3
19  -1.1335438  0.4727417  2.1804613
20
21 Variances:
22           1           2           3
23  0.60241205 0.01478454 1.57299074
```

The result is quite different from the initial set parameters. This is because the three distributions are dense and coupled together, and the estimated value may not be unique. If the distribution is sparse and  $n$  is larger, the effect is better. The following figure is a histogram drawn according to the estimated value, and you can see that the fit is well.



### Prob 3 . Finish all the exercises on pages entitled "Exercise".

The law school data set `law` in the `bootstrap` package contains average LSAT and average GPA for 15 law schools. This data set is a random sample from the universe of 82 law schools in `law82`.

1. Estimate the correlation between LSAT and GPA scores, and compute the bootstrap estimate of the standard error of the sample correlation.
2. Use the `boot` function from package `boot`.

```

1 > print(cor(law$LSAT, law$GPA)) # sample correlation coef
2 [1] 0.7763745
3 > print(cor(law82$LSAT, law82$GPA)) # population correlation coef
4 [1] 0.7599979
5 # bootstrap estimate of standard error of
6 B <- 200 #number of replicates
7 n <- nrow(law) #sample size
8 R <- numeric(B) #storage for replicates
9 #bootstrap estimate of standard error of R
10 for (b in 1:B) {

```

```

11 #randomly select the indices
12 i <- sample(1:n, size = n, replace = TRUE)
13 LSAT <- law$LSAT[i] # i -- vector of indices
14 GPA <- law$GPA[i]
15 R[b] <- cor(LSAT, GPA)
16 }
17 #output
18 > print(se.R <- sd(R))
19 [1] 0.1297349

```

Compute the bootstrap estimation of bias for the law school sample correlation problem

```

1 library(bootstrap)
2 data(law)
3 theta.hat <- cor(law$LSAT, law$GPA)
4 # bootstrap estimate of bias
5 B <- 2000 # number of bootstrap replicates
6 n <- nrow(law) # sample size
7 theta.b <- numeric(B)
8 for (b in 1:B) { # randomly select the indices
9   i <- sample(1:n, size = n, replace = TRUE) # i is a vector of indices
10  LSAT <- law$LSAT[i]
11  GPA <- law$GPA[i]
12  theta.b[b] <- cor(LSAT, GPA)
13 }
14 bias <- mean(theta.b) - theta.hat
15 bias
16 > bias
17 [1] 0.007822616

```

The patch (bootstrap) data contains measurements of a certain hormone in the bloodstream of eight subjects after wearing a medical patch. The parameter of interest is

$$\theta = \frac{E(new) - E(old)}{E(old) - E(placebo)}$$

If  $|\theta| \leq 0.2$ , this indicates bioequivalence of the old and new patches. The statistic is  $\bar{Y}/\bar{Z}$ . Compute a bootstrap estimate of bias in the bioequivalence ratio statistic.

```

1 data(patch, package = "bootstrap")
2 n <- nrow(patch)
3 y <- patch$y
4 z <- patch$z
5 theta.hat <- mean(y) / mean(z)

```

```

6 > print (theta.hat)
7 [1] -0.0713061
8 theta.jack <- numeric(n)
9 for (i in 1:n){
10 theta.jack[i] <- mean(y[-i]) / mean(z[-i])
11 }
12 bias <- (n - 1) * (mean(theta.jack) - theta.hat)
13 > print(bias)
14 [1] 0.008002488

```

Compute the confidence intervals for the patch ratio statistic.

```

1 library(boot)          #for boot and boot.ci
2 data(patch, package = "bootstrap")
3 theta.boot <- function(dat, ind) {
4     #function to compute the statistic
5     y <- dat[ind, 1]
6     z <- dat[ind, 2]
7     mean(y) / mean(z)
8 }
9 y <- patch$y
10 z <- patch$z
11 dat <- cbind(y, z)
12 boot.obj <- boot(dat, statistic = theta.boot, R = 2000)
13 print(boot.obj)
14 print(boot.ci(boot.obj,
15               type = c("basic", "norm", "perc")))
16 #calculations for bootstrap confidence intervals
17 alpha <- c(.025, .975)
18 #normal
19 print(boot.obj$t0 + qnorm(alpha) * sd(boot.obj$t))
20 #basic
21 print(2*boot.obj$t0 - quantile(boot.obj$t, rev(alpha), type=1))
22 #percentile
23 print(quantile(boot.obj$t, alpha, type=6))

```

```

1 ORDINARY NONPARAMETRIC BOOTSTRAP
2
3
4 Call:
5 boot(data = dat, statistic = theta.boot, R = 2000)
6
7
8 Bootstrap Statistics :
9     original      bias    std. error

```

```

10 t1* -0.0713061 0.007706594 0.1009925
11
12 BOOTSTRAP CONFIDENCE INTERVAL CALCULATIONS
13 Based on 2000 bootstrap replicates
14
15 CALL :
16 boot.ci(boot.out = boot.obj, type = c("basic", "norm", "perc"))
17
18 Intervals :
19 Level      Normal      Basic      Percentile
20 95%  (-0.2770, 0.1189 )  (-0.3126, 0.0863 )  (-0.2289, 0.1700 )
21 Calculations and Intervals on Original Scale
22 >
23 >
24 > #calculations for bootstrap confidence intervals
25 > alpha <- c(.025, .975)
26 >
27 > #normal
28 > print(boot.obj$t0 + qnorm(alpha) * sd(boot.obj$t))
29 [1] -0.2692477 0.1266355
30 >
31 > #basic
32 > print(2*boot.obj$t0 -
33 +       quantile(boot.obj$t, rev(alpha), type=1))
34       97.5%      2.5%
35 -0.31247046 0.08631008
36 >
37 > #percentile
38 > print(quantile(boot.obj$t, alpha, type=6))
39       2.5%      97.5%
40 -0.2288723 0.1700097
41

```

**Compute the confidence intervals using the boot. ci function.**

```

1 library(boot)
2 data(law, package = "bootstrap")
3 boot.obj <- boot(law, R = 2000,
4       statistic = function(x, i){cor(x[i,1], x[i,2])})
5 print(boot.ci(boot.obj, type=c("basic", "norm", "perc")))

```

```

1 BOOTSTRAP CONFIDENCE INTERVAL CALCULATIONS
2 Based on 2000 bootstrap replicates
3
4 CALL :
5 boot.ci(boot.out = boot.obj, type = c("basic", "norm", "perc"))
6
7 Intervals :
8 Level      Normal      Basic      Percentile
9 95%      ( 0.5223, 1.0388 ) ( 0.5902, 1.0886 ) ( 0.4641, 0.9625 )
10 Calculations and Intervals on Original Scale

```

---

**Use linear regression model to predict cat heart weights with cat body weights. Compute the confidence intervals of the regression coefficients with bootstrap resampling.**

```

1 library(MASS)
2 data(cats)
3 resample <- function(x) {
4   sample(x,size=length(x),replace=TRUE)
5 }
6 coefs.cats.lm <- function(subset) {
7   fit <- lm(Hwt~Bwt,data=cats,subset=subset)
8   return(coefficients(fit))
9 }
10 cats.lm.sampling.dist <- replicate(1000,
11   coefs.cats.lm(resample(1:nrow(cats))))
12 (limits <- apply(cats.lm.sampling.dist,1,quantile,c(0.025,0.975)))

```

```

1      (Intercept)      Bwt
2 2.5%      -1.948358 3.442391
3 97.5%      1.210477 4.636989

```