

# Reporte de estancia

Alejandro Nieto Ramos, Marc Lavielle (asesor)

1 de agosto de 2016

## El algoritmo EM

### El algoritmo EM en general

Sea  $\mathbf{Y}$  un vector aleatorio que corresponde con los datos observados  $\mathbf{y}$  y cuya función de densidad de probabilidad es  $g(\mathbf{y}; \Psi)$ , donde  $\Psi = (\Psi_1, \dots, \Psi_d)$  es un vector de parámetros desconocidos con espacio parametral  $\Omega$ .

El algoritmo Esperanza Maximización (**EM**) permite calcular estimadores de máxima verosimilitud (**EML**) a través de un proceso iterativo, cuando por falta de datos adicionales, la estimación por máxima verosimilitud sería directa. Desde este punto de vista, al vector de datos observados  $\mathbf{y}$  se le considera como incompleto y a su vez, se toma como una función observable de los datos completos. Ahora bien, cuando se habla de datos incompletos, no solo se hace referencia a la falta de datos en sí, sino también a situaciones donde los datos incompletos representen lo que se podría observar en un experimento hipotético. En este último contexto, se define a  $\mathbf{y}_c$  como el vector de los datos completos y a  $\mathbf{z}$  como al vector de los datos adicionales, es decir, los no observables o faltantes.

Sea entonces  $g_c(\mathbf{x}; \Psi)$  la función de densidad de probabilidad del vector aleatorio  $\mathbf{Y}_c$  que corresponde a los datos completos. Por lo tanto, la función de verosimilitud de  $\Psi$  si se tuvieran los datos completos, estaría dada por  $\log L_c(\Psi) = \log g_c(\mathbf{y}_c; \Psi)$ .

El algoritmo EM resuelve indirectamente el problema de hallar la solución a la ecuación

$$\frac{\partial}{\partial \Psi} \log L(\Psi) = 0, \quad (1)$$

iterando en términos de la función de verosimilitud de los datos completos  $\log g_c(\mathbf{y}_c)$ . Como no se observa, se cambia por su esperanza condicional dado  $\mathbf{y}$ , usando la aproximación actual de  $\Psi$ .

Es decir, sea  $\Psi^{(0)}$  un valor inicial de  $\Psi$ . En la primera iteración, el paso E requiere del cálculo de

$$Q(\Psi; \Psi^{(0)}) = E_{\Psi^{(0)}}[\log L_c(\Psi) | \mathbf{y}].$$

El paso M requiere de la maximización de  $Q(\Psi; \Psi^{(0)})$  con respecto a  $\Psi$  sobre el espacio parametral  $\Omega$ . Es decir, se toma a  $\Psi^{(1)}$  tal que  $Q(\Psi^{(1)}; \Psi^{(0)}) \geq Q(\Psi; \Psi^{(0)})$  para todos los  $\Psi \in \Omega$ . Los pasos E y M se aplican nuevamente cambiando ahora la aproximación actual por  $\Psi^{(1)}$ .

En general, en la iteración  $k + 1$ , el algoritmo es el siguiente:

**Paso E.** Se calcula  $Q(\Psi; \Psi^{(k)})$ , donde  $Q(\Psi; \Psi^{(k)}) = E_{\Psi^{(k)}}[\log L_c(\Psi) | \mathbf{y}]$ .

**Paso M.** Se elige  $\Psi^{k+1}$  como cualquier valor  $\Psi \in \Omega$  que maximice  $Q(\Psi; \Psi^{(k)})$ , es decir,  $Q(\Psi^{(k+1)}; \Psi^{(k)}) \geq Q(\Psi; \Psi^{(k)})$  para todo  $\Psi \in \Omega$ . Se alternan repetidamente ambos pasos hasta que  $|L(\Psi^{(k+1)}) - L(\Psi^{(k)})| < \epsilon$ , para un  $\epsilon > 0$  previamente dado, en caso de que la sucesión  $L(\Psi)$  converja. Para asegurar que la sucesión sea convergente, basta escoger que esté acotada superiormente<sup>1</sup>.

---

<sup>1</sup>Teorema de Bolzano.

## El algoritmo EM para modelos de mezclas

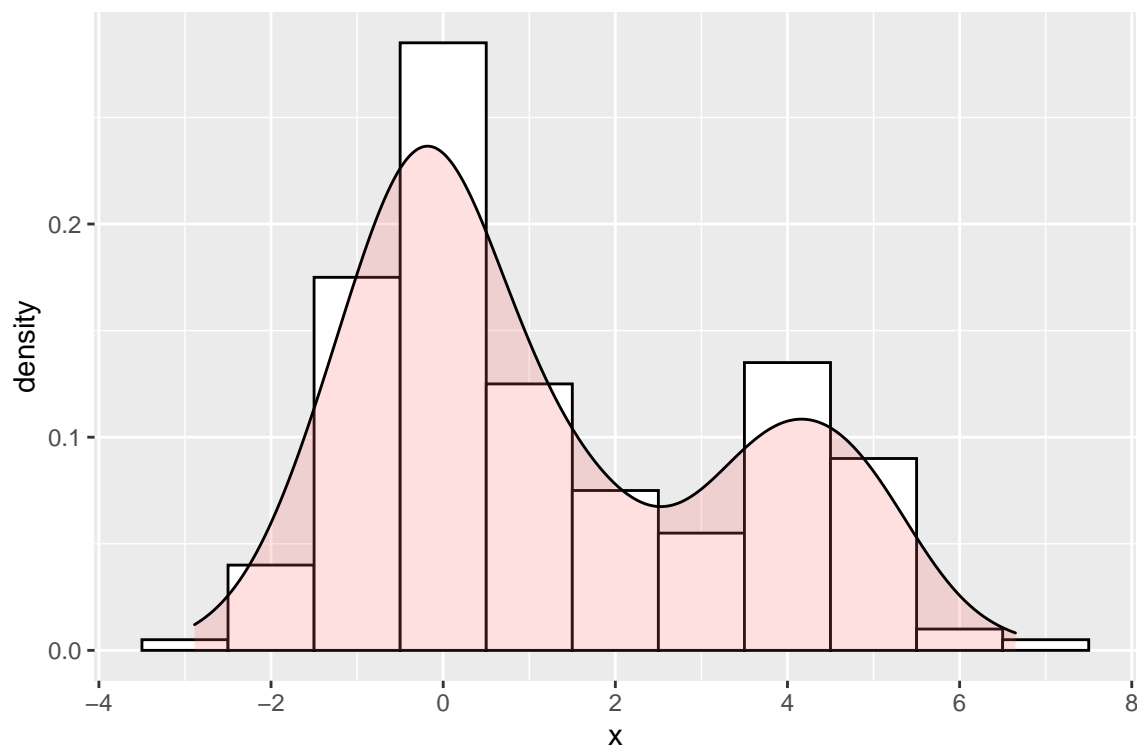
Sea  $(Y_1, \dots, Y_n)$  una muestra aleatoria de tamaño  $n$ , donde  $Y_j$  es un vector de dimensión  $p$  con función de densidad de probabilidad  $f(y_j)$  en  $\mathbb{R}^p$ . Sea  $\mathbf{Y} = (Y_1^T, \dots, Y_n^T)^T$  y defínase una observación de  $\mathbf{Y}$  como  $\mathbf{y} = (y_1^T, \dots, y_n^T)^T$ . Supóngase que que la densidad  $f(y_j)$  de  $Y_j$  se puede escribir en la forma

$$f(y_j) = \sum_{i=1}^g \pi_i f_i(y_j), \quad (2)$$

donde las  $f_i$  son las densidades y los valores  $\pi_i$  son tales que  $\sum_{i=1}^g \pi_i = 1$  con  $0 \leq \pi_i \leq 1$  para  $i = 1, \dots, g$ .

A las cantidades  $\pi_1, \dots, \pi_g$  se les conoce como proporciones de la mezcla o pesos. Dado que las funciones  $f_1(y_j), \dots, f_g(y_j)$  son densidades, la ecuación (??) define una densidad. A las  $f_i(y_j)$  se les conoce como densidades componentes de la mezcla. Cabe señalar que en esta formulación, el número de componentes se considera fijo, aunque en muchas aplicaciones pueda ser desconocido.<sup>2</sup> A la densidad (??) se le conoce como densidad de una mezcla finita de  $g$  componentes.

```
mezcla<-function(n,pesos,mu,sigma)
{
  G <- length(mu)
  Z <- sample(1:G, n, prob=pesos, replace=T)
  x<-NULL
  for (g in 1:G)
  {
    x<-c(x,rnorm(length(which(Z==g)),mu[g],sigma[g]))
  }
  return(x)
}
```



Gráfica de la mezcla de dos normales, la primera con  $\mu = 0, \sigma = 1$  y una proporción del 70%, y la segunda con  $\mu = 4, \sigma = 1$  (y una proporción del 30%).

<sup>2</sup>En un caso así, el número de componentes se tendría que aproximar.

## Interpretación

Una manera clara de generar un vector aleatorio  $Y_{ij}$  con la mezcla de  $g$  densidades componente  $f(y_j)$  es la siguiente. Defínase a  $Z_j$  como la variable aleatoria categórica que toma valores en  $1, \dots, g$  con probabilidades respectivas  $\pi_1, \dots, \pi_g$ , y supóngase que la densidad condicional de  $Y_j$  dado  $Z_j = i$  es  $f_i(y_j)$  para  $i = 1, \dots, g$ . Entonces la densidad marginal de  $Y_j$  está dada por  $f(y_j)$ . Considerando este enfoque, la variable  $Z_j$  se puede pensar como la componente de pertenencia del vector  $Y_j$ . Se define entonces un vector de componentes de pertenencia de dimensión  $g$ ,  $\mathbf{Z}_j$ , en lugar de la variable categórica única  $Z_j$ , donde el  $i$ -ésimo elemento de  $\mathbf{Z}_j$ ,  $Z_{ij} = (\mathbf{Z}_j)_i$  es 0 o 1 de acuerdo a si la componente de origen de  $Y_j$  en la mezcla es igual o no a  $i$ ,  $i = 1, \dots, g$ . Por lo tanto,  $\mathbf{Z}_j$  se distribuye de acuerdo a una distribución multinomial que consiste de un experimento en  $g$  categorías con probabilidades  $\pi_1, \dots, \pi_g$ , es decir,  $\Pr(\mathbf{Z}_j = \mathbf{z}_j) = \pi_1^{z_{1j}} \dots \pi_g^{z_{gj}}$ ; ésta se representa como  $Mult_g(1, \boldsymbol{\pi})$  donde  $\boldsymbol{\pi} = (\pi_1, \dots, \pi_g)$ .

## Enfoque directo

Sea

$$f(y_j; \boldsymbol{\Psi}) = \sum_{i=1}^g \pi_i f_i(y_j; \theta_i) \quad (3)$$

un modelo de mezcla paramétrico de una muestra aleatoria observada  $\mathbf{y} = (y_1^T, \dots, y_n^T)^T$ , donde  $\boldsymbol{\Psi} = (\pi_1, \dots, \pi_{g-1}, \xi^T)^T$  es el vector que contiene todos los parámetros desconocidos de la mezcla del modelo y  $\xi$  es el vector de los parámetros  $\theta_1, \dots, \theta_g$  a priori tomados como distintos. La log verosimilitud de  $\boldsymbol{\Psi}$  que puede obtenerse de los datos observados está dada por

$$\log L(\boldsymbol{\Psi}) = \sum_{j=1}^n \log(f(y_j; \boldsymbol{\Psi})) \quad (4)$$

$$= \sum_{j=1}^n \log \left\{ \sum_{i=1}^g \pi_i f_i(y_j; \theta_i) \right\} \quad (5)$$

Si se calcula el EMV de  $\boldsymbol{\Psi}$  tomando esta log verosimilitud, es preciso resolver  $\frac{\partial}{\partial \boldsymbol{\Psi}} \log L(\boldsymbol{\Psi}) = 0$ . Esta ecuación se puede manipular de modo que el EMV,  $\hat{\boldsymbol{\Psi}}$ , satisfaga

$$\hat{\pi}_i = \sum_{j=1}^n \frac{1}{n} \tau_i(y_j; \hat{\boldsymbol{\Psi}}), \quad (6)$$

para todo  $i = 1, \dots, g$ , y

$$\sum_{i=1}^g \sum_{j=1}^n \tau_i(y_j; \hat{\boldsymbol{\Psi}}) \frac{\partial}{\partial \xi} \log f_i(y_j; \theta_i) = 0, \quad (7)$$

donde

$$\tau_i(y_j; \boldsymbol{\Psi}) = \frac{\pi_i f_i(y_j; \theta_i)}{\sum_{h=1}^g \pi_h f_h(y_j; \theta_h)} \quad (8)$$

es la probabilidad posterior de que  $y_j$  pertenezca a la componente  $i$  de la mezcla.

## Enfoque donde se ve a los datos como incompletos

Retómese el mismo enfoque del que se habló en el algoritmo EM, donde el vector  $\mathbf{y} = (y_1^T, \dots, y_n^T)^T$  se toma como incompleto, ya que los vectores de pertenencia a cierta componente asociados,  $z_1, \dots, z_n$ , no están disponibles. Con este enfoque, en donde cada  $y_j$  se ve como que viene de una de los componentes del modelo de mezcla que se desea ajustar,  $\mathbf{z}_j$  es un vector de dimensión  $g$  con  $z_{ij} = (z_j)_i$  igual a 1 o 0 de de acuerdo a si  $y_j$  vino o no del componente  $i$ -ésimo de la mezcla, para todo  $i = 1, \dots, g$  y  $j = 1, \dots, n$ . El vector de datos completo es por lo tanto,  $\mathbf{y}_c = (y^T, \mathbf{z}^T)^T$ , donde  $\mathbf{z} = (z_1^T, \dots, z_n^T)^T$ .

Los vectores de pertenencia a cierta componente  $z_1, \dots, z_n$  se toman como realizaciones de los vectores aleatorios  $Z_1, \dots, Z_n$ , donde por la característica de independencia de los datos, se puede asumir que se distribuyen independientes con distribución  $Mult_g(1, \pi)$ . Asumir esto significa que la distribución del vector de los datos completos  $\mathbf{Y}_c$  da origen a la distribución apropiada para los datos incompletos  $\mathbf{Y}$ . La log verosimilitud de los datos completos para  $\Psi$ ,  $\log L_c(\Psi)$  es

$$\log L_c(\Psi) = \sum_{i=1}^g \sum_{j=1}^n z_{ij} \{\log \pi_i + \log f_i(y_j; \theta_i)\}. \quad (9)$$

### Paso E

Se aplica el algoritmo EM a este problema tratando a los  $z_{ij}$  como datos faltantes. Como la log verosimilitud de los datos completos,  $L_c$ , es lineal con respecto a los datos no observados  $z_{ij}$ , el paso E en la iteración  $k+1$  requiere únicamente del cálculo de la esperanza condicional actual de  $Z_{ij}$  dados los datos observados  $y$ , donde  $Z_{ij}$  es la variable aleatoria correspondiente a  $z_{ij}$ . Ahora bien,

$$\mathbb{E}_{\Psi^{(k)}}(Z_{ij}|\mathbf{y}) = \Pr_{\Psi^{(k)}}(Z_{ij} = 1|\mathbf{y}) \quad (10)$$

$$= \tau_i(y_j; \Psi^{(k)}), \quad (11)$$

y de acuerdo a la ecuación (??),

$$\tau_i(y_j; \Psi^{(k)}) = \frac{\pi_i^{(k)} f_i(y_j; \theta_i^{(k)})}{f(y_j; \Psi^{(k)})} \quad (12)$$

$$= \frac{\pi_i^{(k)} f_i(y_j; \theta_i^{(k)})}{\sum_{h=1}^g \pi_h^{(k)} f_h(y_j; \theta_h^{(k)})} \quad (13)$$

para todo  $i = 1, \dots, g$ ,  $j = 1, \dots, n$ . La cantidad  $\tau_i(y_j; \Psi^{(k)})$  es la probabilidad posterior de que el  $j$ -ésimo miembro de la muestra con valor observado  $y_j$  pertenezca al componente  $i$ -ésimo de la mezcla. Aplicando los cálculos hechos en el párrafo anterior, tomando la esperanza condicional de (??) dado  $y$  se tiene que

$$Q(\Psi; \Psi^{(k)}) = \sum_{i=1}^g \sum_{j=1}^n \tau_i(y_j; \Psi^{(k)}) \{\log \pi_i + \log f_i(y_j; \theta_i)\}. \quad (14)$$

### Paso M

El paso M en la iteración  $k+1$ , requiere que se maximice de manera global  $Q(\Psi; \Psi^{(k)})$  con respecto a  $\Psi$  sobre el espacio parametral  $\Omega$  para obtener la estimación actualizada  $\Psi^{(k+1)}$ . Para el modelo de mezcla finita, los estimadores actualizados  $\pi_i^{(k+1)}$  de las proporciones de la mezcla  $\pi_i$  se calculan independientemente de los estimadores actualizados  $\xi^{(k+1)}$  del vector de parámetros  $\xi$  que contiene a los parámetros desconocidos en las densidades componentes.

Si los  $z_{ij}$  fueran observables, entonces el EVM de  $\pi_i$  estaría dado por la ecuación (??), es decir,  $\hat{\pi}_i = \sum_{j=1}^n \frac{1}{n} z_{ij}$  para  $i = 1, \dots, g$ .

Como el paso E solo reemplaza cada  $z_{ij}$  por su esperanza condicional  $\tau_i(y_j; \Psi^{(k)})$  en la log verosimilitud de los datos completos, el estimador actualizado de  $\pi_i$  está dado al reemplazar cada  $z_{ij}$  en la ecuación anterior por  $\tau_i(y_j; \Psi^{(k)})$ , es decir

$$\pi_i^{(k+1)} = \sum_{j=1}^n \frac{1}{n} \tau_i(y_j; \Psi^{(k)}) \quad (15)$$

para todo  $i = 1, \dots, g$ .

Con respecto a la actualización de  $\xi$  en el paso M de la iteración  $k + 1$ , se puede ver de la ecuación (??) que  $\xi^{(k+1)}$  se obtiene como una raíz apropiada de

$$\sum_{i=1}^g \sum_{j=1}^n \tau_i(y_j; \Psi^{(k)}) \frac{\partial}{\partial \xi} \log f_i(y_j; \theta_i) = 0, \quad (16)$$

En ciertos casos, la solución se puede encontrar en forma analítica, como se verá en la siguiente sección. Finalmente, el algoritmo termina al igual que como se explicó al final de la sección ??.

## Mezcla de normales

Sea ahora

$$f(y_j; \Psi) = \sum_{i=1}^g \pi_i \phi_i(y_j; \mu_i, \Sigma_i), \quad (17)$$

donde  $\phi(y_j; \mu_i, \Sigma_i) = \frac{1}{(\sqrt{2\pi})^n \sqrt{|\Sigma_i|}} \exp(-\frac{1}{2}(y_j - \mu_i)^T \Sigma_i^{-1} (y_j - \mu_i))$ .

En esta caso,  $\Psi = (\pi_1, \dots, \pi_g, \xi^T)^T$ , donde  $\xi$  contiene a los elementos de las componente medias  $\mu_i$  y las distintas matrices componentes de varianza-covarianza  $\Sigma_i$ , ( $i = 1, \dots, g$ ).

Se ha visto que en la iteración  $k + 1$  del paso E, se reemplazan las variables de pertenencia a cierta componente, es decir, los  $z_{ij}$ , por sus esperanzas condicionales dadas por sus probabilidades posteriores de pertenecer a los datos observados  $y_j$ ,  $\tau_i(y_j; \Psi^{(k)})$ . Por lo tanto,

$$\tau_i(y_j; \Psi) = \frac{\pi_i \phi_i(y_j; \mu_i, \Sigma_i)}{\sum_{h=1}^g \pi_h \phi_h(y_j; \mu_h, \Sigma_h)} \quad (18)$$

para  $i = 1, \dots, g$ ,  $j = 1, \dots, n$ .

A continuación se presenta el código en R para calcular  $\tau_i$  en el caso de una de mezcla de funciones qe se distribuyen normales.

```
calculotau<-function(x,teta)
{
  n<-length(x)
  G<-length(teta$p)
  tau<-matrix(NA,n,G)
  for (g in 1:G)
    tau[,g]<-teta$p[g]*dnorm(x,teta$mu[g],teta$sigma[g])
  tau=tau/matrix(rep(rowSums(tau),G),nrow=n)
  return(tau)
}
```

El paso M para componentes normales existe en forma analítica. Las actualizaciones de las medias componentes  $\mu_i$  y las matrices de varianza covarianza  $\Sigma_i$  son:

$$\mu_i^{(k+1)} = \frac{\sum_{j=1}^n \tau_{ij}^{(k)} y_j}{\sum_{j=1}^n \tau_{ij}^{(k)}} \quad (19)$$

$$\Sigma_i^{(k+1)} = \frac{\sum_{j=1}^n \tau_{ij}^{(k)} (y_j - \mu_i^{(k+1)})(y_j - \mu_i^{(k+1)})^T}{\sum_{j=1}^n \tau_{ij}^{(k)}} \quad (20)$$

para  $i = 1, \dots, g$ , donde  $\tau_{ij}^{(k)} = \tau_i(y_j; \Psi^{(k)})$  con  $i = 1, \dots, g$  y  $j = 1, \dots, n$ .

Para simplificar los cálculos computacionales, es conveniente escribir las actualizaciones anteriores (??) y (??) en términos de las esperanzas condicionales actuales de los estadísticos suficientes  $S_{i1}$ ,  $S_{i2}$  y  $S_{i3}$  para  $\Psi$  para el enfoque de los datos completos. Es decir,

$$S_{i1}^{(k)} = \sum_{j=1}^n \tau_{ij}^{(k)}$$

$$S_{i2}^{(k)} = \sum_{j=1}^n \tau_{ij}^{(k)} y_j$$

$$S_{i3}^{(k)} = \sum_{j=1}^n \tau_{ij}^{(k)} y_j^2.$$

De este modo, el paso E implementado en R queda como sigue:

```
estadistica<-function(x,Z)
{
  G<-dim(Z)[2]
  M<-dim(Z)[3]
  if (is.na(M))
  {
    M <- 1
    dim(Z) <- c(dim(Z),1)
  }
  s1 <- 0
  s2 <- 0
  s3 <- 0
  for (m in 1:M)
  {
    Z.m <- Z[, ,m]
    s1 <- s1 + colSums(Z.m)
    s2 <- s2 + x %*% Z.m
    s3 <- s3 + (x^2) %*% Z.m
  }
  est<-list(s1=s1/M,s2=s2/M,s3=s3/M)
  return(est)
}

pasoE<-function(x,teta)
{
  tau <- calculotau(x,teta)
  s <- estadistica(x,tau)
  return(s)
}
```

Ahora bien, el código en R para calcular el paso M es

```
pasoM<-function(S,n)
{
  p<-S$s1/n
  mu<-S$s2/S$s1
  sigma<-sqrt(S$s3/S$s1-(mu)^2)
  teta<-list(p=p,mu=mu,sigma=sigma)
  return(teta)
}
```

donde la varianza está dada por

```
logLikelihood <- function(x,df)
{
  K <- dim(df)[1]
  G <- (dim(df)[2]-1)/3
  ll <- NULL
  for (k in (1:K))
  {
    lk <- 0
    for (g in (1:G))
    {
      pi.k <- df[k,g+1]
      mu.k <- df[k,G+g+1]
      sigma.k <- df[k,2*G+g+1]
      lk <- lk + pi.k * dnorm(x, mean=mu.k, sd=sigma.k)
    }
    ll <- c(ll, sum(log(lk)))
  }
  df$deviance <- -2*ll
  return(df)
}
```

El estimador de la  $i$ -ésima proporción  $\pi_i$  está dado por la ecuación (??).

Finalmente, el programa iterativo que hace uso de las funciones anteriores se presenta a continuación:

```
mixtl.em <- function(x, teta0, K)
{
  G<-length(mu)
  col.names <- c("iteration", paste0("p",1:G), paste0("mu",1:G), paste0("sigma",1:G))

  teta.est <- matrix(NA,K+1,3*G+1)
  teta.est[1,] <- c(0, teta0$p, teta0$mu, teta0$sigma)

  teta<-teta0
  for (k in 1:K)
  {
    s<-pasoE(x,teta)
    teta<-pasoM(s,n)
    teta.est[k+1,] <- c(k, teta$p, teta$mu, teta$sigma)
  }
  df.em <- as.data.frame(teta.est)
  names(df.em) <- col.names
  return(df.em)
}

# valores iniciales para el algoritmo
pesos0<-c(.5,.5)
mu0<-c(1,2)
sigma0<-c(.5,2)

K <- 250
M <- 1

G<-length(mu)
```

```

col.names <- c("iteration", paste0("p",1:G), paste0("mu",1:G), paste0("sigma",1:G))
teta0<-list(p=pesos0,mu=mu0,sigma=sigma0)
s0<-list(s1=0,s2=0,s3=0)

## EM
df.em <- mixt1.em(x, teta0, K)
df.em <- logLikelihood(x,df.em)
print(round(df.em[nrow(df.em),],3))

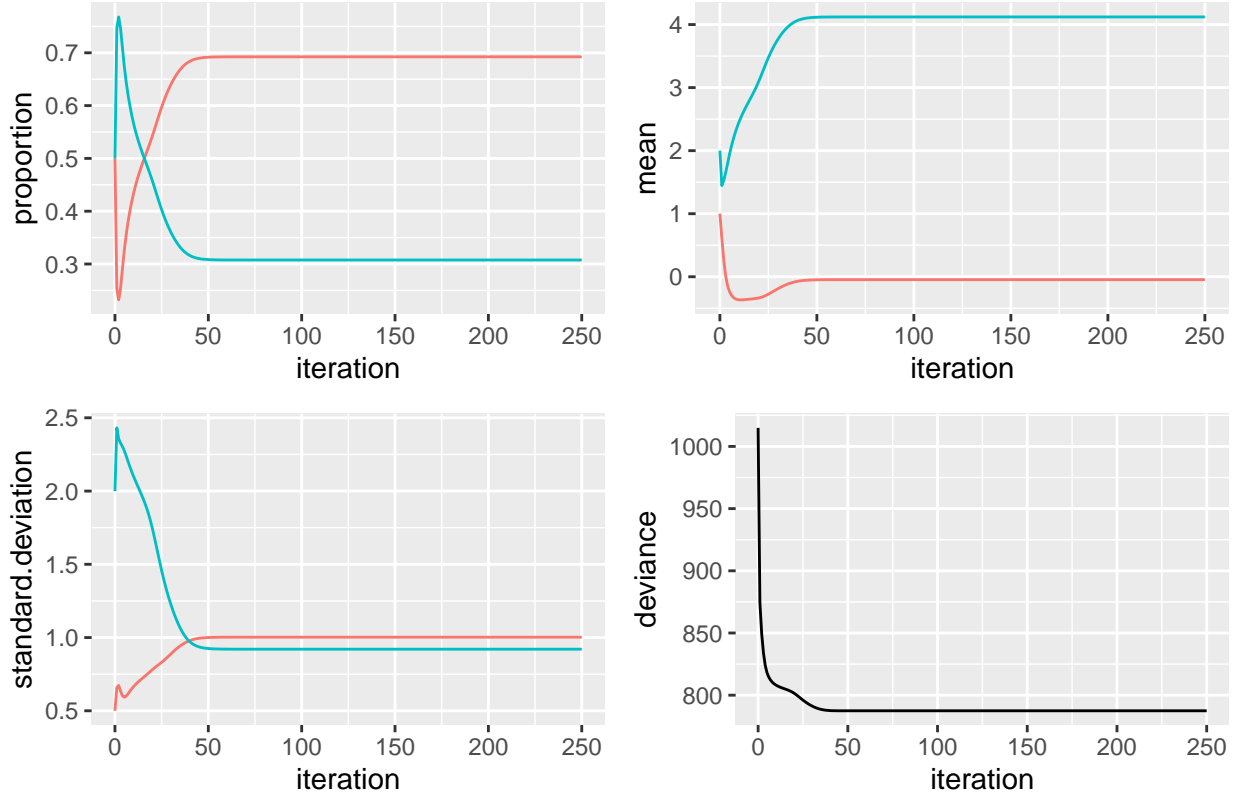
##      iteration    p1    p2    mu1    mu2 sigma1 sigma2 deviance
## 251          250 0.692 0.308 -0.046 4.121  1.002  0.921  787.436

graphConvergence <- function(df, title=NULL)
{
  G <- (dim(df)[2]-1)/3
  df1<-melt(df[,c(1,2,3)],id.vars="iteration",variable="proportion")
  graf1 <- ggplot(df1,aes(iteration,value))+geom_line(aes(color=proportion)) +
    ylab("proportion") + theme(legend.position="none")
  df2<-melt(df[,c(1,4,5)],id.vars="iteration",variable="mean")
  graf2 <- ggplot(df2,aes(iteration,value))+geom_line(aes(color=mean)) +
    ylab("mean") + theme(legend.position="none")
  df3<-melt(df[,c(1,6,7)],id.vars="iteration",variable="standard.deviation")
  graf3 <- ggplot(df3,aes(iteration,value))+geom_line(aes(color=standard.deviation)) +
    ylab("standard.deviation") + theme(legend.position="none")
  graf4 <- ggplot(df,aes(iteration,deviance))+geom_line()
  grid.arrange(graf1,graf2,graf3,graf4,nrow=2, top=title)
}
graphConvergence(df.em, title="EM")

```



## EM



Ejemplo donde se calculan los parámetros de la mezcla del ejemplo presentado al principio ( $pesos = (70, 30)$ ,  $\mu = (0, 4)$ ,  $\sigma = (1, 1)$ ). Los valores iniciales dados son  $\pi_0 = (.5, .5)$ ,  $\mu = (1, 2)$ ,  $\sigma = (.5, 2)$

## El algoritmo EM estocástico

En el paso E del algoritmo EM, cada vector  $z_j$  se reemplaza por su esperanza condicional dados los datos observados:  $z_j^{(k)} = (z_{1j}^{(k)}, \dots, z_{gj}^{(k)})^T$ , donde  $z_{ij}^{(k)}$  se define como en la ecuación (??) y es la probabilidad posterior actual de que la  $j$ -ésima observación venga del componente  $i$ -ésimo de la mezcla ( $i = 1, \dots, g; j = 1, \dots, n$ ). Existen varias versiones del algoritmo EM estocástico y se verán dos de ellas en este documento.

### El algoritmo EM estocástico de Monte Carlo (MCEM)

Para este caso, la probabilidad posterior actual se usa en un paso estocástico E, donde se obtiene una sola muestra<sup>3</sup> de la distribución condicional actual de  $z$  dados los datos observados  $y$ . Por la independencia de los datos completos, esto se lleva a cabo obteniendo una muestra para cada  $j$  con  $j = 1, \dots, n$ ; es decir, se obtiene una muestra  $z_j^{(1k)}$  de una distribución multinomial con  $g$  categorías que tienen probabilidad  $z_{ij}^{(k)}$ .

El paso M consiste en encontrar el EMV de  $\Psi$  como si los  $y_1, \dots, y_n$  fueran clasificados de manera determinista de acuerdo a  $z_1^{(1k)}, \dots, z_n^{(1k)}$ .

En resumen, el algoritmo estocástico empieza con probabilidades multinomiales arbitrarias  $\tau_{ij}^{(0)}$ , con  $i = 1, \dots, g$ , para cada observación  $j$ . Luego, si  $\Psi^{(k)}$  es el valor de  $\Psi$  en la iteración  $k$ , se llevan a cabo los siguientes pasos:

#### Paso estocástico E.

<sup>3</sup>Se pueden obtener varias, y de hecho, en el algoritmo programado para este documento, se tiene esta opción.

Para cada  $j$ , se obtiene una muestra<sup>4</sup> de la distribución multinomial con probabilidades  $\tau_{ij}^{(k)}$  con  $i = 1, \dots, g$ , las cuales pueden ser arbitrarias cuando  $k = 0$  (por simplicidad, se pueden tomar todas iguales) o están dadas por el paso M previo cuando  $k > 0$ . En la  $k$ -ésima iteración, se define  $z_j^{(1_k)} = (z_1^{(1_k)}, \dots, z_n^{(1_k)})$  la muestra de la  $j$ -ésima observación  $y_j$ . Esto resulta en una partición de las  $n$  observaciones  $y_1, \dots, y_n$  con respecto a las  $g$  componentes de la mezcla, porque para cada  $j$ , exactamente una de estas  $z_{ij}^{(1_k)}$  es 1 y las otras 0.

#### Paso M.

Calcular

$$\pi_i^{(k+1)} = \sum_{j=1}^n \frac{1}{n} z_{ij}^{(1_k)}, \quad (21)$$

con  $i = 1, \dots, g$ .

Si las  $f_i$  son normales de dimensión  $p$ ,  $\mu_i^{(k+1)}$  y  $\Sigma_i^{(k+1)}$  están dadas, respectivamente, por las ecuaciones (??) y (??). Las probabilidades posteriores para el paso E estocástico en la siguiente iteración se actualizan de la siguiente forma:

$$\tau_{ij}^{(k+1)} = \tau_i(y_j; \Psi^{(k+1)}) \quad (22)$$

$$= \frac{\pi_i^{(k+1)} f_i(y_j; \theta^{(k+1)})}{f(y_j; \Psi^{(k+1)})} \quad (23)$$

para  $i = 1, \dots, g$ .

Se presenta el código en R del paso estocástico para el caso de una mezcla de distribuciones normales escalares:

```
pasoS<-function(x,teta,M)
{
  n<-length(x)
  G<-length(teta$p)
  Z<-array(NA,c(n,G,M))
  tau<-calculotau(x,teta)
  for (i in 1:n)
  {
    Z[i,,]<-rmultinom(n=M,size=1,prob=tau[i,])
  }
  if (M==1) Z <- Z[, ,1]
  return(Z)
}
```

## El algoritmo EM de aproximación estocástica (SAEM)

Como se observó en la sección anterior, el algoritmo de EM de Monte Carlo reemplaza el paso E por una aproximación de Monte Carlo basada en un gran número de simulaciones independientes de los parámetros individuales no observados  $z$ . En el algoritmo de EM de aproximación estocástica (**SAEM**) se reemplaza el paso E por una aproximación estocástica basada en una sola simulación de  $z$ . Este algoritmo requiere de un valor inicial  $\theta_0$ . En la  $k$ -ésima iteración se tienen tres pasos: un paso de simulación que llamaremos paso S, un paso de aproximación estocástica y finalmente el paso M de maximización.

#### Paso S.

Se obtiene una muestra  $z_i^{(k)}$  de una distribución condicional  $p(z_i | y_i; \Psi^{(k-1)})$ .

<sup>4</sup>Este es un caso particular del algoritmo EM de Monte Carlo, en el cual, se obtienen  $M$  muestras y la función  $Q$  se aproxima por medio de la suma  $Q(\Psi, \Psi^{(k)}) = \frac{1}{M} \sum_{m=1}^M \log p(\Psi | z^{(m_k)}; y)$ .

### Paso SA.

Se actualiza  $Q_{k-1}(\Psi)$  de acuerdo a

$$Q_k(\Psi) = Q_{k-1}(\Psi) + \gamma_k(\log p(z_i|y_i; \Psi) - Q_{k-1}(\Psi)),$$

donde  $\{\gamma_k\}_{k=1}^{\infty}$  es una sucesión decreciente tal que  $\gamma_k \in \mathbb{R}$  con  $\gamma_1 = 1$ .

### Paso M

Se actualiza  $\Psi^{(k+1)}$  igual que antes.

Se presenta el código en R para el caso del paso de aproximación estocástica.

```
pasoSA<-function(x,Z,s.old,gamma)
{
  S<-estadistica(x,Z)
  s11<-s.old$s1+gamma*(S$s1-s.old$s1)
  s12<-s.old$s2+gamma*(S$s2-s.old$s2)
  s13<-s.old$s3+gamma*(S$s3-s.old$s3)
  s.new<-list(s1=s11,s2=s12,s3=s13)
  return(s.new)
}
```

Se tiene entonces el algoritmo SAEM para la mezcla presentada para el algoritmo determinista, donde  $\{\gamma_k\}_{k=1}^{\infty} = \{1\}_{k=1}^{\infty}$  (es decir, éste es un algoritmo MCEM puro) y  $M = 1$ .

```
mixt1.saem <- function(x, teta0, K, K1, M)
{
  G<-length(mu)
  col.names <- c("iteration", paste0("p",1:G), paste0("mu",1:G), paste0("sigma",1:G))

  K2 <- K - K1
  gamma<-c(rep(1,K1),1/(1:K2))

  teta.est <- matrix(NA,K+1,3*G+1)
  teta.est[1,] <- c(0, teta0$p, teta0$mu, teta0$sigma)

  teta<-teta0
  s<-s0
  for (k in 1:K)
  {
    Z<-pasoS(x,teta,M)
    s<-pasoSA(x,Z,s,gamma[k])
    teta<-pasoM(s,n)
    teta.est[k+1,] <- c(k, teta$p, teta$mu, teta$sigma)
  }
  df.saem <- as.data.frame(teta.est)
  names(df.saem) <- col.names
  return(df.saem)
}

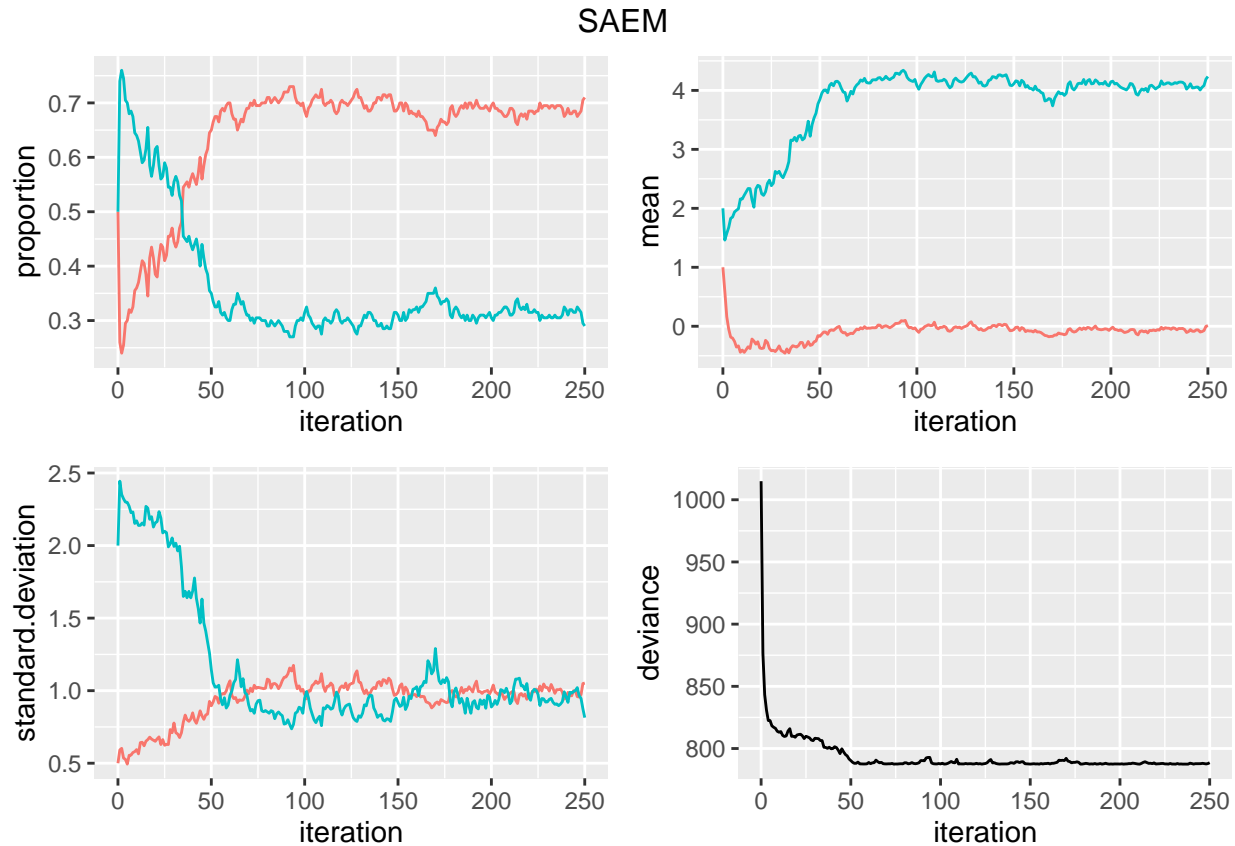
K <- 250
K1 <- 250
M <- 1

## SAEM
df.saem <- mixt1.saem(x, teta0, K, K1, M)
```

```
df.saem <- logLikelihood(x,df.saem)
print(round(df.saem[nrow(df.saem),],3))
```

```
##      iteration  p1  p2  mu1  mu2 sigma1 sigma2 deviance
## 251          250 0.71 0.29 0.011 4.236  1.053  0.814  788.594
```

```
graphConvergence(df.saem, title="SAEM")
```



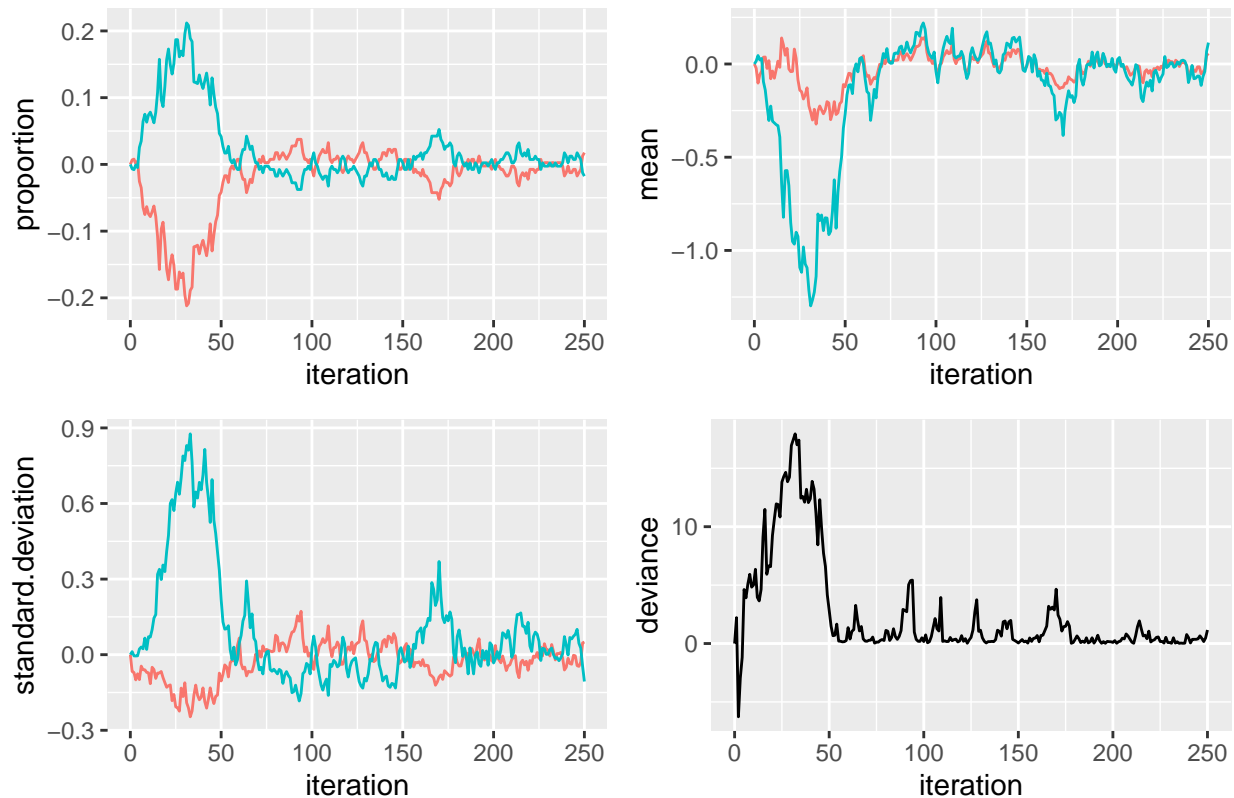
Haciendo la diferencia entre el método SAEM y EM se tiene:

```
## SAEM - EM
df.diff <- df.saem - df.em
df.diff$iteration <- df.em$iteration
print(round(df.diff[nrow(df.saem),],3))
```

```
##      iteration  p1    p2  mu1  mu2 sigma1 sigma2 deviance
## 251          250 0.018 -0.018 0.057 0.115  0.05 -0.106   1.159
```

```
graphConvergence(df.diff, title="SAEM - EM")
```

## SAEM – EM



Si ahora  $M = 10$ , se puede ver que los valores del algoritmo SAEM, el cual sigue siendo un algoritmo MCEM puro, se aproximan a los del algoritmo EM (las gráficas de SAEM se suavizan):

```
K <- 250
M <- 10

## EM
df.em <- mixt1.em(x, teta0, K)
df.em <- logLikelihood(x, df.em)

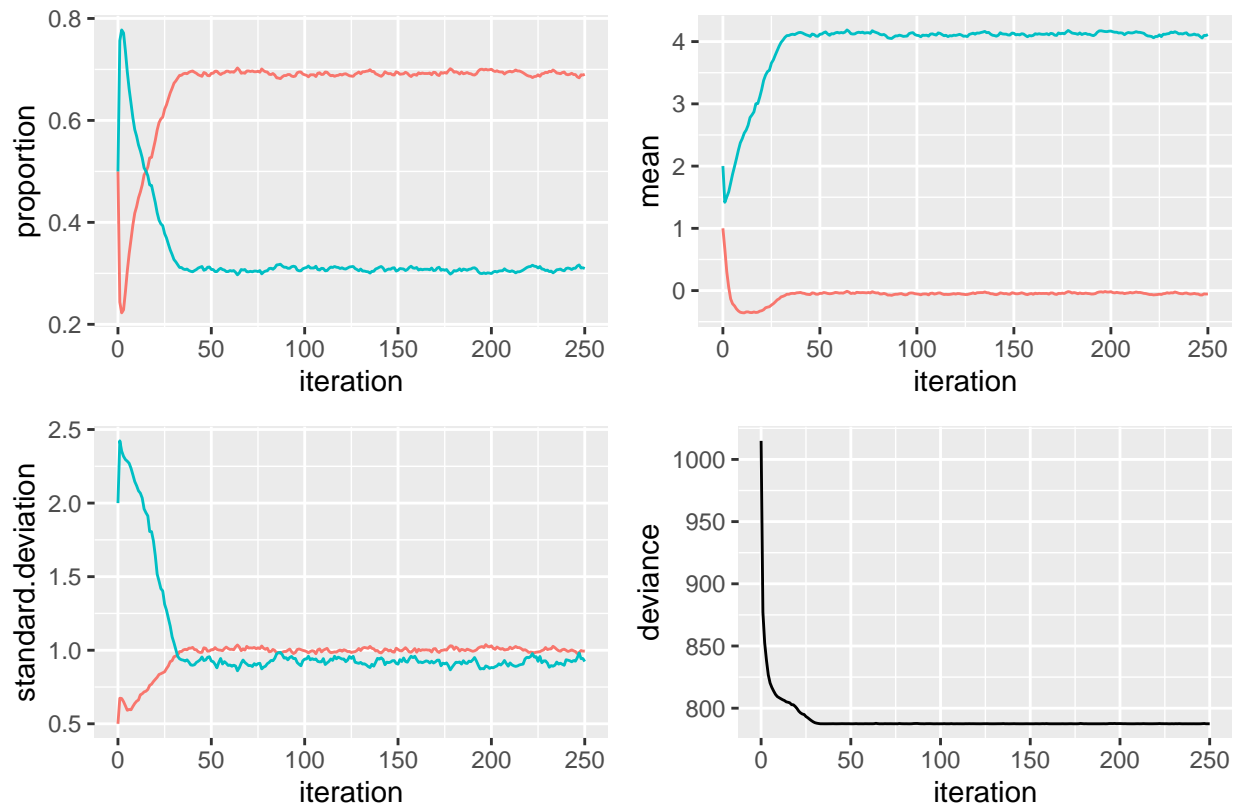
K <- 250
K1 <- 250
M <- 10

## SAEM
df.saem <- mixt1.saem(x, teta0, K, K1, M)
df.saem <- logLikelihood(x, df.saem)
print(round(df.saem[nrow(df.saem), ], 3))

##      iteration    p1    p2    mu1    mu2 sigma1 sigma2 deviance
## 251         250 0.691 0.31 -0.053 4.112  0.994  0.923  787.453

graphConvergence(df.saem, title="SAEM")
```

## SAEM



Si ahora se toma la sucesión decreciente  $\{\gamma_k\}_{k=1}^{\infty} = \{\frac{1}{k}\}_{k=1}^{\infty}$ , la gráfica de los valores del algoritmo SAEM se suaviza aún más.

```
K <- 250
M <- 10

## EM
df.em <- mixt1.em(x, teta0, K)
df.em <- logLikelihood(x, df.em)

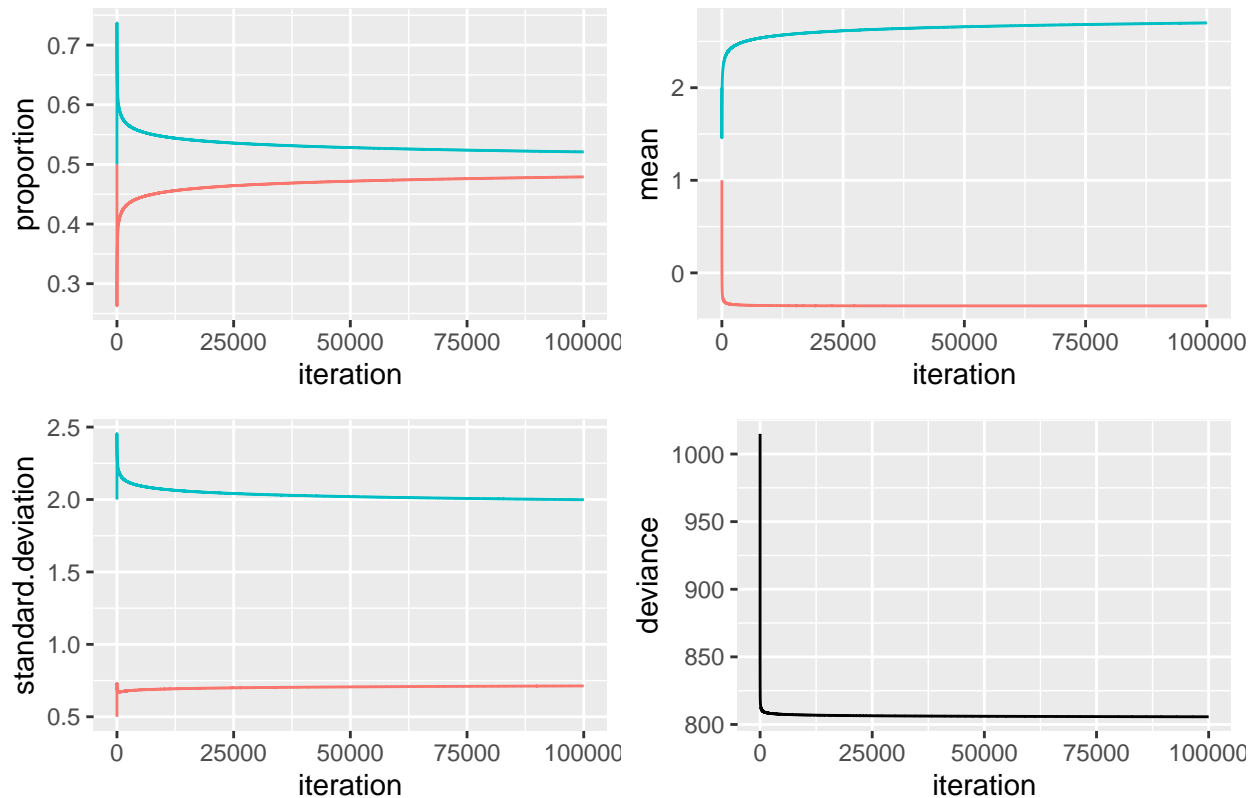
K <- 100000
K1 <- 0
M <- 10

## SAEM
df.saem <- mixt1.saem(x, teta0, K, K1, M)
df.saem <- logLikelihood(x, df.saem)
print(round(df.saem[nrow(df.saem), ], 3))

##      iteration    p1    p2    mu1 mu2 sigma1 sigma2 deviance
## 100001      1e+05 0.479 0.521 -0.356 2.7   0.713  1.999  805.696

graphConvergence(df.saem, title="SAEM")
```

## SAEM



Lo cierto es que los valores convergen muy lentamente, además de que no convergen a los valores del algoritmo EM. En este caso entonces, se puede formar una sucesión para tener una combinación de los algoritmos MCEM puro y SAEM, para poder encontrar una mejor aproximación al valor verdadero de los parámetros. Es decir, los primeros  $k_1$  valores de la sucesión podrían ser todos 1 y a partir del valor del valor  $k_1 + 1$  podría tenerse propiamente una sucesión decreciente de números reales.

En el siguiente ejemplo se hace uso de ambos de ambos algoritmos estocásticos con  $k_1=100$  y  $\{\gamma_k\}_{k=1}^{\infty} = \{\frac{1}{k}\}_{k=101}^{\infty}$ .

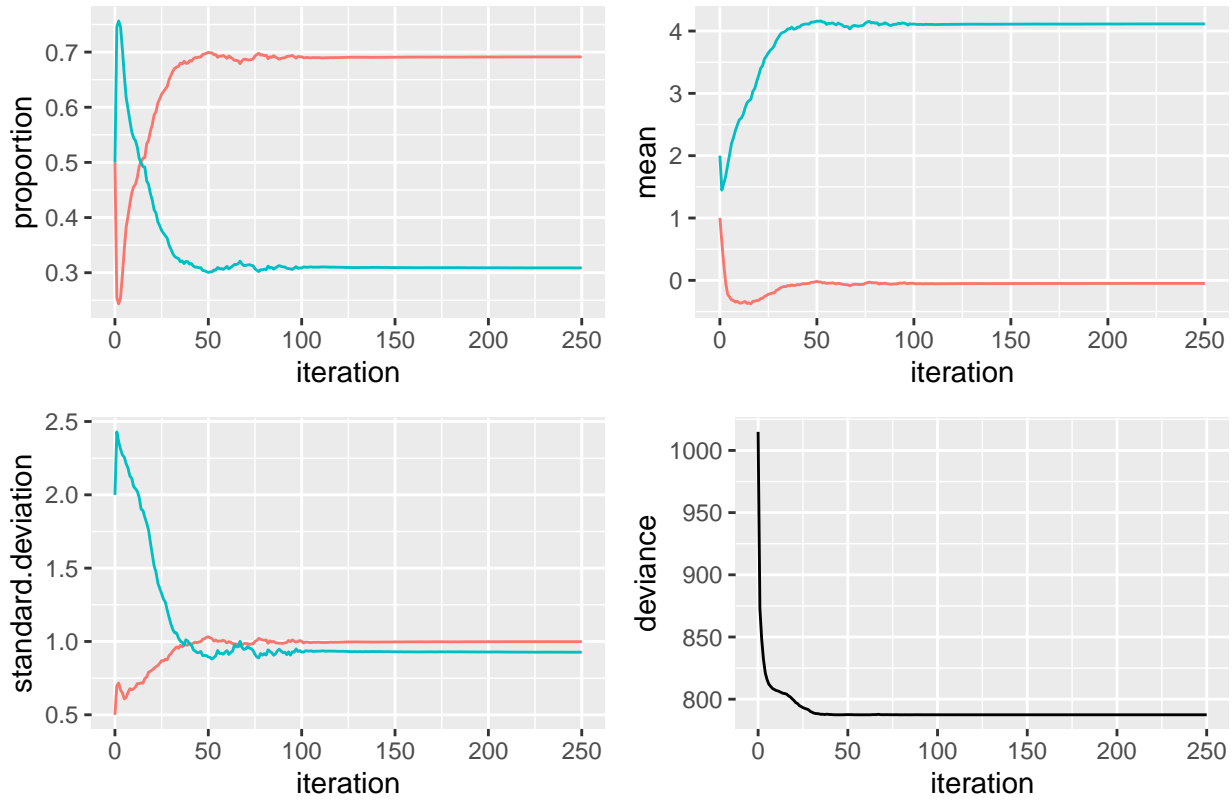
```
K <- 250
K1 <- 100
M <- 10

## SAEM
df.saem <- mixt1.saem(x, teta0, K, K1, M)
df.saem <- logLikelihood(x, df.saem)
print(round(df.saem[nrow(df.saem), ], 3))

##      iteration    p1    p2    mu1    mu2 sigma1 sigma2 deviance
## 251         250 0.691 0.309 -0.049 4.114  0.999  0.927  787.439

graphConverge(df.saem, title="SAEM")
```

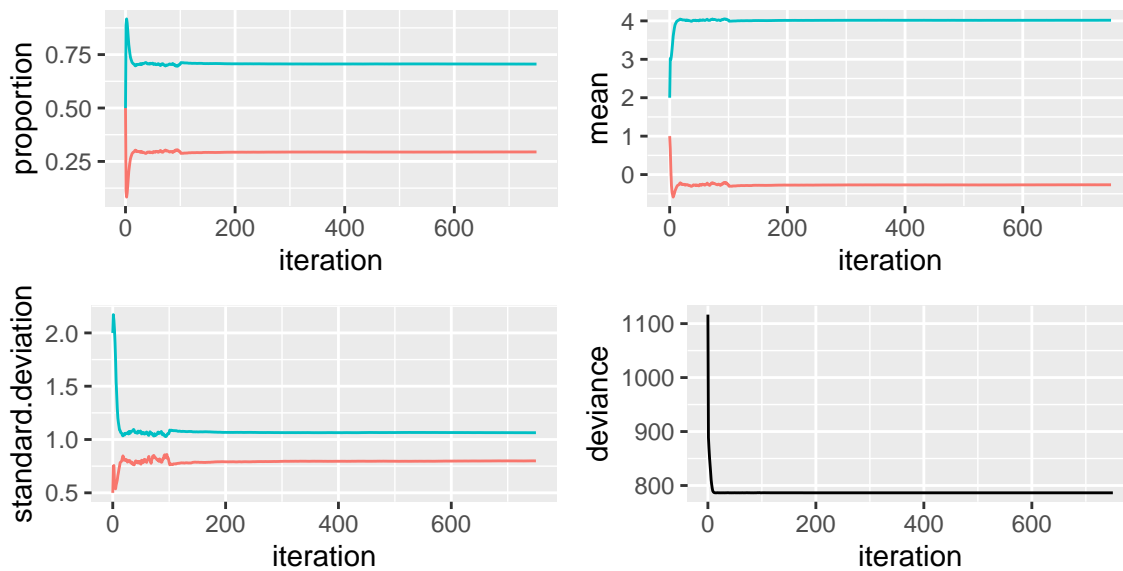
### SAEM



Por último, se puede modificar el algoritmo SAEM para que en vez de tomar la sucesión  $\{\gamma_k\}_{k=1}^{\infty}$  se tome la sucesión  $\{\gamma_k^{\alpha}\}_{k=1}^{\infty}$ , donde  $0.5 < \alpha < 1$ . Con  $\gamma = .8$  se observa que la aproximación a los valores de los parámetros es aún mejor, además de que las oscilaciones se vuelven más pequeñas a medida que se tienen más iteraciones.

```
##      iteration    p1    p2    mu1    mu2 sigma1 sigma2 deviance
## 751          750 0.294 0.706 -0.266 4.016   0.8   1.064   786.444
```

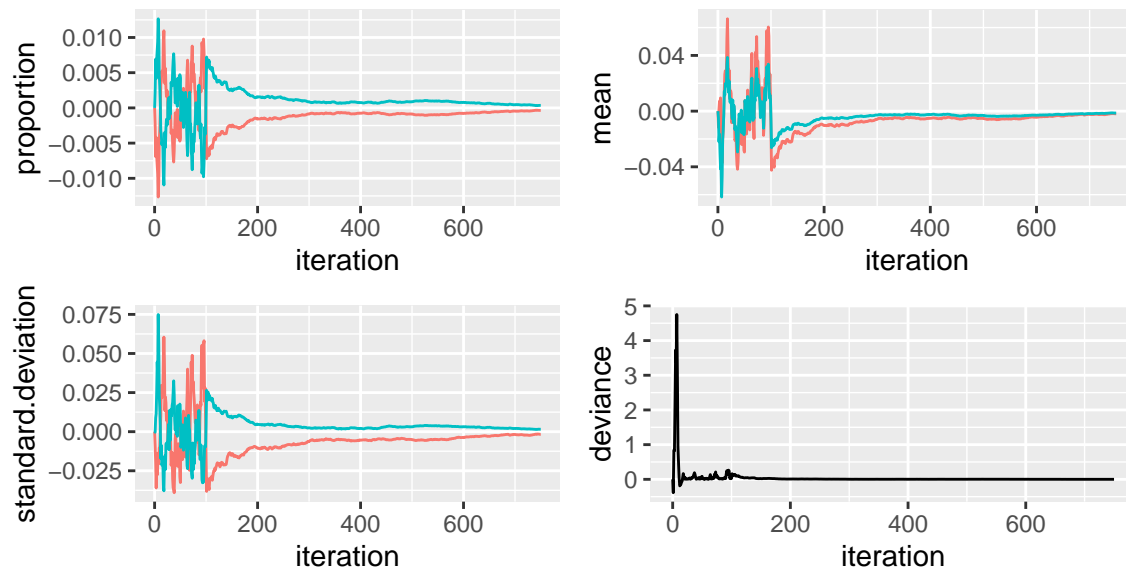
### SAEM





```
##      iteration p1 p2      mu1      mu2 sigma1 sigma2 deviance
## 751          750  0  0 -0.002 -0.001 -0.002  0.002          0
```

### SAEM – EM



## References

- [1] Lavielle, M. (2014). Mixed effects models for the population approach: models, tasks, methods and tools. CRC press.
- [2] McLachlan, G. J., & Peel, D. (2000). Finite Mixture Models. Wiley series in probability and statistics.
- [3] McLachlan, G. J., & Krishnan, T. (1997). The EM algorithm and extensions. Wiley series in probability and statistics.