

Computational Statistics

Chapter 3: EM algorithm

1. Let $\mathbf{Y} = (Y_1, \dots, Y_n)$ be an i.i.d. sample from a mixture of a normal distribution $\mathcal{N}(\mu, \sigma)$ and a uniform distribution $\mathcal{U}([-a, a])$, with pdf

$$g(y; \theta) = \pi \phi(y; \mu, \sigma) + (1 - \pi)c, \quad (1)$$

where $\phi(\cdot; \mu, \sigma)$ is the normal pdf, $c = (2a)^{-1}$, π is the proportion of the normal distribution in the mixture and $\theta = (\mu, \sigma, \pi)^T$ is the vector of parameters. Typically, the uniform distribution corresponds to outliers in the data. The proportion of outliers in the population is then $1 - \pi$. We want to estimate parameter θ using the EM algorithm

- (a) Using the functions `sample`, `rnorm` and `runif`, generate a sample of size $n = 100$. Draw a box plot of the data.
 - (b) Write an EM algorithm for this problem.
 - (c) Apply the EM algorithm to the data, with different initializations. Draw the estimated probabilities $1 - z_i^{(t)}$ of being an outlier, as a function of y_i . Does it make sense?
 - (d) Compare the estimates with those computed using the `optim` function.
2. We will now apply the same idea to linear regression. We assume that we have an independent sample $\mathbf{Y} = (Y_1, \dots, Y_n)$, where the distribution of each observation Y_i is a mixture of a normal distribution $\mathcal{N}(v_i^T \beta, \sigma^2)$ and a uniform distribution $\mathcal{U}([-a, a])$, v_i being a vector of covariates and β a vector of coefficients. The pdf of Y_i is

$$g(y; \theta) = \pi \phi(y; v_i^T \beta, \sigma) + (1 - \pi)c, \quad (2)$$

with $\theta = (\beta^T, \sigma, \pi)^T$ and $c = (2a)^{-1}$.

- (a) Using the functions `rnorm` and `runif`, generate a sample of size $n = 100$, with $v_i \sim \mathcal{U}([-6, 6])$, $\beta = (1, 2)^T$, $\sigma = 2$, $a = 20$ and $\pi = 0.5$. Draw a scatter plot of the data.
- (b) Compute the ordinary least squares estimates (OLS) of the coefficients, and draw the corresponding line.

- (c) Write an EM algorithm for this problem. (In the M-step, you will have to solve a weighted least-squares problem. You can use the `lm` functions with input parameter `weights`).
- (d) Apply the EM algorithm to the data, taking the OLS estimates as initial values. Draw the line with coefficients equal to the MLEs. Plot the points (v_i, y_i) such that $z_i^{(t)} < 0.5$ as filled circles. Does it make sense?
3. We consider again the problem of estimating the parameters in a mixture of a normal distribution $\mathcal{N}(\mu, \sigma)$ and a uniform distribution $\mathcal{U}([-a, a])$, where a is a known constant. The observed data are an iid sample w_1, \dots, w_n from W with pdf

$$g(w; \theta) = \pi \phi(w; \mu, \sigma) + (1 - \pi)c, \quad (3)$$

where $\phi(\cdot; \mu, \sigma)$ is the normal pdf, $c = (2a)^{-1}$, π is the proportion of the normal distribution in the mixture and $\theta = (\mu, \sigma, \pi)^T$ is the vector of parameters. Typically, the uniform distribution corresponds to outliers in the data. The proportion of outliers in the population is then $1 - \pi$.

We have seen how to find the MLE $\hat{\theta}$ of θ using the EM algorithm. We now want to estimate the variance of $\hat{\theta}$.

- (a) Set $\theta = (0, 1, 0.9)$ and $a = 5$. Generate $N = 1000$ samples of size $n = 100$. For each sample, compute $\hat{\theta}$ using the EM algorithm. Estimate the variance of $\hat{\theta}$.
- (b) We now consider one sample w_1, \dots, w_n and we wish to estimate $\text{Var}(\hat{\theta})$ from that sample, without knowing the true value of θ . We will use two methods.
- i. Louis' method: compute $\hat{i}_{\mathbf{x}}(\hat{\theta})$ and estimate $\hat{i}_{\mathbf{z}|\mathbf{y}}(\hat{\theta})$ by Monte Carlo simulation; compute an estimate of $\hat{i}_{\mathbf{y}}(\hat{\theta})$ using the missing information principle equation, and its inverse $\hat{i}_{\mathbf{y}}(\hat{\theta})^{-1}$.
 - ii. Bootstrapping: generate $B = 1000$ bootstrap samples. Estimate $\text{Var}(\hat{\theta})$ by the sample variance of the bootstrap estimates of $\hat{\theta}$.
- (c) Compare the estimates of $\text{Var}(\hat{\theta})$ obtained by the different methods.