# STA6171: Statistical Computing for DS 1
# EM Algorithm

Ick Hoon Jin

Yonsei University, Department of Statistics and Data Science

2020.10.07

# Introduction

- Develop for handling missing outcomes.
- Use a lot in optimization

- The expectation-maximization (EM) algorithm is an iterative optimization strategy motivated by a notion of missingness and by consideration of the conditional distribution of what is missing given what is observed.

- Popularity of the EM algorithm
  - Simple to implement
  - Reliable to find the global optimum.

# Introduction

*Observed Data: X*
*Missing Data: Z* ) *Complete Data: Y = (X, Z)*

- Frequentist Setting

  - Observed data from $X$ along with missing data from $Z$.
  - Complete data $Y = (X, Z)$. *We want to maximize $L(\theta|x)$*
  - Given observed data $x$, maximize a likelihood $L(\theta|x)$. Difficult to work with this likelihood *Difficult to work with this likelihood (Not able to apply Newton method)*
  - A easier way is working with the density $Y|\theta$ and $Z|x, \theta$. *Use EM algorithm*
    *$\varepsilon$     $f(Y|\theta)$   $f(Z|X, \theta)$*
- Bayesian Setting: Interest Often focuses on estimating the mode of a posterior distribution $f(\theta|x)$. *→ maximize $L(\theta|x)$*

- Missing data may not truly be missing: they may be only a conceptual ploy that simplifies the problem. In this case, $Z$ is often referred to as *latent*. *Sometimes we want to maximize $L(\theta|x)$*
  *$Z$ : latent variable  $Z|X, \theta$*
  *$(Y|\theta)$*

# Missing Data and Marginalization - Frequentist

observed
$$Y = (X, Z)$$
missing

- In the presence of missing data, only a function of the complete-data $y$ is observed.

$$f(y \mid \theta) = f(x, z \mid \theta) = f_X(x \mid \theta) f_{Z \mid X}(z \mid x, \theta)$$

$$l_y(\theta) = l_X(\theta) + \log f_{Z \mid X}(z \mid x, \theta)$$

- Assume that the missing data are random, so that

$$f(y|\theta) = f(x, z|\theta) = f_X(x|\theta) \cdot f_{Z|X}(z|x, \theta).$$

Thus, it follows that

$$l_X(\theta) = l_y(\theta) - \log f_{Z \mid X}(z \mid X, \theta)$$

maximization diff.

maximization is straightforward.

$$l_X(\theta) = l_Y(\theta) - \log f_{Z|X}(Z|X, \theta).$$

- Useful when maximizing $l_X(\theta)$ can be difficult but maximizing the complete log-likelihood $l$ is simple.

# Missing Data and Marginalization - Bayesian

Complete data likelihood $\quad L(\theta|y) = L(\theta|x, z)$

$L(\theta|x)$ : marginalization of $L(\theta|y)$

- View the likelihood $L(\theta|x)$ as a marginalization of the complete-data likelihood $L(\theta|y) = L(\theta|x,z)$.

- Consider there to be missing parameter $\psi$, whose inclusion simplifies Bayesian calculations even though $\psi$ is of no interest itself. Since $Z$ and $\psi$ are both missing random quantities, it *nuisance parameter (latent)* matters little whether we use notation that suggests the missing variables to be unobserved data or parameters.

Introduction
EM Algorithm
EM Variants

EM Algorithm
Examples
Variance Estimation

# EM Algorithm

- EM algorithm ~~iteratively seeks to maximize $L(\theta|x)$~~ with respect to $\theta$.

- Let $\theta^{(t)}$ ⟶ *estimate.* denote the estimated maximizer at iteration $t$, for $t = 0, 1, \cdots$.

- Define $Q(\theta|\theta^{(t)})$ to be the ~~expectation of the joint log-likelihood for the complete data, conditional on the observed data $X = x$.~~

$$Q\left(\theta|\theta^{(t)}\right) = E\left\{\log L(\theta|Y)\Big|x, \theta^{(t)}\right\}.$$

$$Q(\theta|\theta^{(t)}) = E\{\log L(\theta|Y) | x, \theta^{(t)}\}$$

- Then, $Q\left(\theta|\theta^{(t)}\right)$ is maximized w.r.t $\theta$, that is $\theta^{(t+1)}$ is found such that

$$Q\left(\theta^{(t+1)}|\theta^{(t)}\right) \geq Q\left(\theta|\theta^{(t)}\right)$$

for all $\theta \in \Theta$.

Introduction
EM Algorithm
EM Variants

EM Algorithm
Examples
Variance Estimation

# Example: Simple Exponential Distribution

$$f(y|\theta) = \prod_{i=1}^{2} \theta e^{-\theta y_i} = \theta^2 e^{-\theta \sum_{i=1}^{2} y_i}$$

- Suppose $Y_1, Y_2 \overset{iid}{\sim} exp(\theta)$ and $y_1 = 5$ is observed but the value $y_2$ is missing.

  $Y = (Y_1, Y_2)$ → observed, missing

  $\ell_y(\theta) = 2\log\theta - \theta y_1 - \theta y_2$

- The complete-data log likelihood function is

$$\log L(\theta|y) = 2\log\theta - \theta y_1 - \theta y_2.$$

  $E(Y_2|Y_1, \theta^{(t)}) = E(Y_2|\theta^{(t)})$

  make $y_2$ using conditional expectation

  $= 1/\theta^{(t)}$

- Because

$$E(Y_2|y_1, \theta^{(t)}) = E(Y_2|\theta^{(t)}) = \frac{1}{\theta^{(t)}},$$

  the conditional expectation of $\log L(\theta|Y)$ yields

  $Q(\theta|\theta^{(t)}) = E[\ell_y(\theta)|y_1, \theta^{(t)}] = E[2\log\theta - \theta y_1 - \theta y_2 | y_1, \theta^{(t)}]$

$$Q(\theta|\theta^{(t)}) = 2\log\theta - 5\theta - \theta/\theta^{(t)}.$$

  $= 2\log\theta - 5\theta - \theta E[Y_2|y_1, \theta^{(t)}] = 2\log\theta - 5\theta - \theta/\theta^{(t)}$

- The maximizer of $Q(\theta|\theta^{(t)})$ with respect to $\theta$ is easily to found to be the root of $2/\theta - 5 - 1/\theta^{(t)} = 0$.

  $\frac{2}{\theta} = \frac{5\theta^{(t)} + 1}{\theta^{(t)}}$   $\theta^{(t+1)} = \frac{2\theta^{(t)}}{5\theta^{(t)} + 1}$