

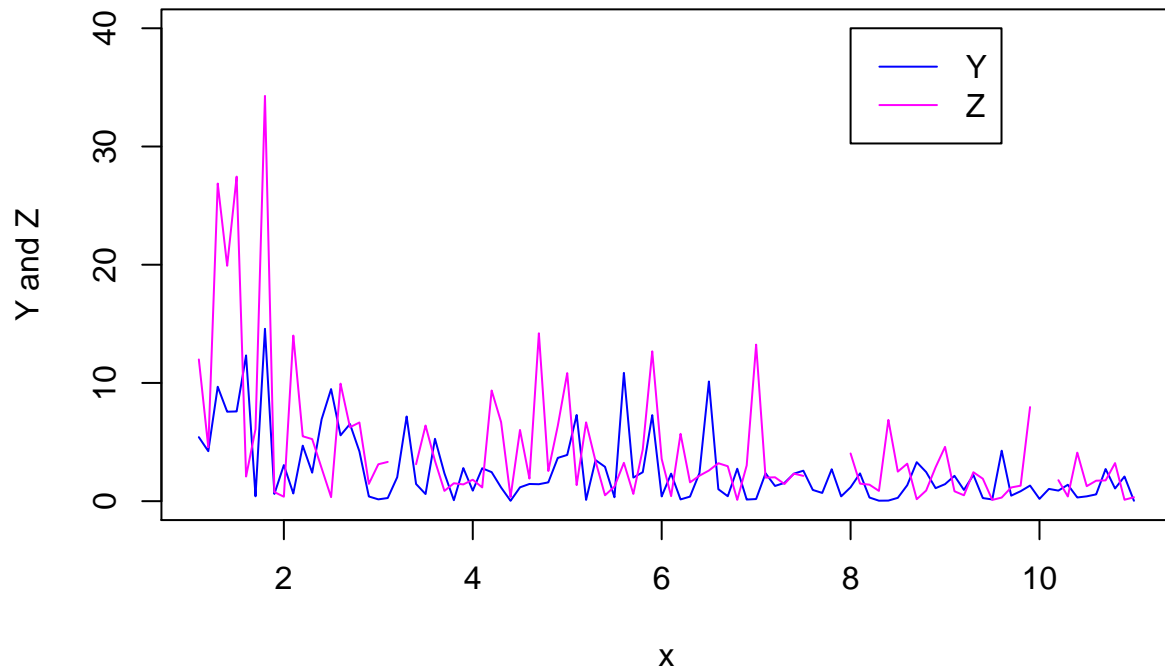
## lab6-2

Shwetha

12/10/2020

### Question 2: EM algorithm

1.



From the above plot we can observe that Z has some missing values , hence the discontinuity in plot. Z has a lot of peaks , intensity of peaks are higher at the beginning ie when x values are small and with increase in X the peak intensity reduces. Y also has a similar pattern as Z. For both Y and Z we can say that the variation when compared to X is high for smaller values of X and low for bigger values of X.

2.

There are missing values in Z. The model of Y and Z are as follows

$$Y_i \sim \exp(\frac{X_i}{\lambda})$$

$$Z_i \sim \exp(\frac{X_i}{2\lambda})$$

The goal is to derive an EM algorithm that estimates.

Since Y and Z follow exponential distribution, pdf of these can be stated as follows

$$f(Y|X, \lambda) = \frac{X_i}{\lambda} \exp(-\frac{X_i Y_i}{\lambda})$$

$$f(Z|X, \lambda) = \frac{X_i}{2\lambda} \exp(-\frac{X_i Z_i}{2\lambda})$$

Likelihood of lambda :

$$L(\lambda) = \prod_{i=1}^n \frac{X_i}{\lambda} \exp(-\frac{X_i Y_i}{\lambda}) \prod_{i=1}^n \frac{X_i}{2\lambda} \exp(-\frac{X_i Z_i}{2\lambda})$$

$$L(\lambda) = \prod_{i=1}^n \frac{X_i^2}{2\lambda^2} \exp(-\frac{X_i Y_i}{\lambda}) \exp(-\frac{X_i Z_i}{2\lambda})$$

$$L(\lambda) = (\frac{1}{2\lambda^2})^n \exp(-\sum_{i=1}^n \frac{X_i Y_i}{\lambda} + \frac{X_i Z_i}{2\lambda}) \prod_{i=1}^n X_i^2$$

Log likelihood :

$$\log(L(\lambda)) = -n \log(2) - 2n \log(\lambda) - \sum_{i=1}^n \frac{X_i Y_i}{\lambda} - \sum_{i=1}^n \frac{X_i Z_i}{2\lambda} + \sum_{i=1}^n X_i^2$$

Terms which do not depend on lambda are can be written as C(constants).

$$\log(L(\lambda)) = -2n \log(\lambda) - \sum_{i=1}^n \frac{X_i Y_i}{\lambda} - \sum_{i=1}^n \frac{X_i Z_i}{2\lambda} + C$$

We know that Z has missing values , so this can be split in to two parts. Z with known values :

$$\sum_{known} \frac{X_i Z_i}{2\lambda}$$

Z with unknown values :

$$\sum_{unknown} \frac{X_j Z_j}{2\lambda}$$

$$\sum_{i=1}^n \frac{X_i Z_i}{2\lambda} = \sum_{known} \frac{X_i Z_i}{2\lambda} + \sum_{unknown} \frac{X_j Z_j}{2\lambda}$$

Substituting this to log likelihood :

$$\log(L(\lambda)) = -2n \log(\lambda) - \sum_{i=1}^n \frac{X_i Y_i}{\lambda} - \sum_{known} \frac{X_i Z_i}{2\lambda} - \sum_{unknown} \frac{X_j Z_j}{2\lambda} + C$$

Expected log likelihood :

$$E[\log(L(\lambda))] = E \left[ -2n \log(\lambda) - \sum_{i=1}^n \frac{X_i Y_i}{\lambda} - \sum_{known} \frac{X_i Z_i}{2\lambda} - \sum_{unknown} \frac{X_j Z_j}{2\lambda} + C \right]$$

All terms other than the summation of unknown Z are on observed values. So the uncertainty is only in this term, so we need to calculate expected value for this term and rest are constant values hence expected value is the same. Z<sub>j</sub> represents the missing Z values, this can be estimated by using exponential distribution as follows:

$$E[Z_j] = \frac{2\lambda_{prev}}{X_j}$$

$$E\left[\sum_{unknown} \frac{X_j Z_j}{2\lambda}\right] = \sum_{unknown} \frac{X_j}{2\lambda} \frac{2\lambda_{prev}}{X_j} = \frac{U\lambda_{prev}}{\lambda}$$

Where U is the number of unknown values in Z.

E-Step is given by :

$$E[\log(L(\lambda))] = -2n \log(\lambda) - \sum_{i=1}^n \frac{X_i Y_i}{\lambda} - \sum_{known} \frac{X_i Z_i}{2\lambda} - \frac{U\lambda_{prev}}{\lambda} + C$$

For M step we just need to differentiate expected log likelihood wrt to lambda and equate it to 0 to find the estimate of lambda value, this will be the lambda value for next iteration.

$$\frac{\partial E[\log(L(\lambda))]}{\partial \lambda} = \frac{-2n}{\lambda} + \frac{1}{\lambda^2} \sum_{i=1}^n X_i Y_i + \frac{1}{2\lambda^2} \sum_{known} X_i Z_i + \frac{U\lambda_{prev}}{\lambda^2} = 0$$

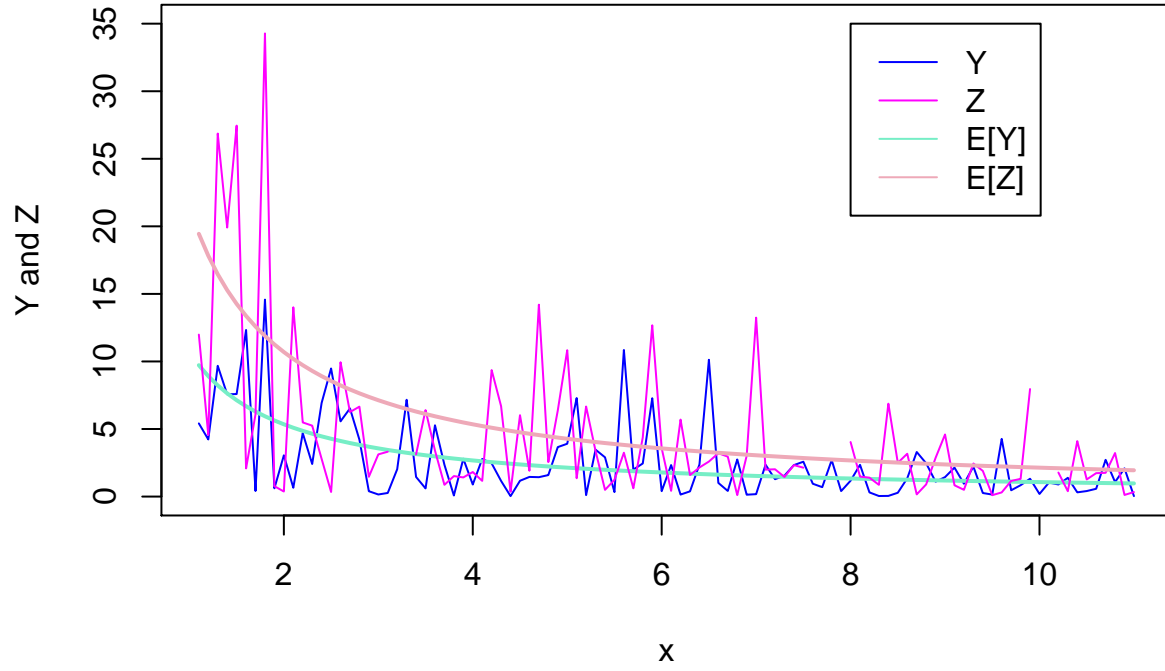
$$\lambda_{mle} = \frac{1}{2n} \sum_{i=1}^n X_i Y_i + \frac{1}{4n} \sum_{known} X_i Z_i + \frac{U\lambda_{prev}}{2n}$$

3.

```
n = nrow(physical)
U = length(which(is.na(physical$Z)))
known_z = which(!is.na(physical$Z))
count = 1
lambda_diff = Inf
eps = 0.001
lambda = 100
repeat{
  lambda_new = (sum(physical$X * physical$Y) + (0.5 * sum(physical$X[known_z] * physical$Z[known_z]))) +
  lambda_diff = abs(lambda - lambda_new)
  lambda = lambda_new
  count = count + 1
  if(lambda_diff < eps){break}
}
```

Optimal lambda is 10.6956555 Number of iterations are 6

4.



As Y and Z follow exponential distribution , their expected value will be as follows.

$$E[Y] = \frac{\lambda}{X_i}$$

$$E[Z] = \frac{2\lambda}{X_i}$$

From the graph we can say that the computed lambda seems to be reasonable as both E[Y] and E[Z] lines seem to represent the mean of the actual plot pretty well. The variation with respect to X is higher in Z when compared to Y and even this is captured as the mean line of Z is higher than that of Y