# ST340 Lab 4: Expectation Maximization

*2019–20*

## 1: Expectation Maximization—mixture of Gaussians

```r
if (!require(mvtnorm)) {
  install.packages("mvtnorm")
  library(mvtnorm)
}
```

```
## Loading required package: mvtnorm
```

```
## Warning: package 'mvtnorm' was built under R version 3.5.2
```

(a) Define parameters for $K = 3$ multivariate normal distributions.

```r
K <- 3
mus.actual <- matrix(0,3,2)
mus.actual[1,] <- c(0,0)
mus.actual[2,] <- c(4,4)
mus.actual[3,] <- c(-4,4)
```

(b) Generate the covariance matrices randomly.

```r
Sigmas.actual <- list()
for (k in 1:K) {
  mtx <- matrix(1,2,2)
  mtx[1,2] <- runif(1)*2-1
  mtx[2,1] <- mtx[1,2]
  Sigmas.actual[[k]] <- mtx*exp(rnorm(1))
}
```

(c) Generate some random mixture weights.

```r
ws <- runif(K)
ws <- ws/sum(ws)
```

(d) Generate 1000 data points in $\mathbb{R}^2$. Hint: look at `?rmvnorm`.

```r
n <- 1000
p <- 2
xs <- matrix(0,n,p)
cols <- rep(0,n)
for (i in 1:n) {
  # sample from the mixture by sampling a mixture component k...
  # YOUR CODE HERE

  # ...and then sampling from that mixture component
  # YOUR CODE HERE

  cols[i] <- k
}
```

(e) Plot the data points coloured by cluster.

```R
plot(xs,col=cols,pch=20)
```

(f) Plot the data points without the colours.

```R
plot(xs,pch=20)
```

(g) Run the EM algorithm on your generated data. You can try seeing what happens if $K = 2$ or $K = 4$ as well. . .

```R
source("em_mixture_gaussians.R")
print(system.time(out <- em_mix_gaussian(xs,K=3)))

my.colors <- rep(0,n)
for (i in 1:n) {
  my.colors[i] <- rgb(out$gammas[i,1],out$gammas[i,2],out$gammas[i,3])
}

plot(xs,col=my.colors,pch=20)
```

## 2: Expectation Maximization—mixture of Bernoullis

(a) Create a file called `em_mixture_bernoullis.R` which contains a function called `em_mix_bernoulli` that is the analogue of `em_mix_gaussian`. You could use `em_mixture_gaussians.R` as a template.

Hint: do not initialize any of the cluster `mus` to be either 0 or 1. (Do you know why?)

Hint: if, by numerical error, any of the parameters $\mu_{kj}$ are greater than 1, the algorithm will most likely fail. To avoid this, you can, after the $M$ step, perform the update

```R
mus[which(mus > 1,arr.ind=TRUE)] <- 1 - 1e-15
```

where `mus` is a $K \times p$ matrix.

(b) Test your code.

```R
n <- 500; p <- 50
K.actual <- 2
mix <- runif(K.actual); mix <- mix / sum(mix)
mus.actual <- matrix(runif(K.actual*p),K.actual,p)
zs.actual <- rep(0,n)

xs <- matrix(0,n,p)
for (i in 1:n) {
  cl <- sample(K.actual,size=1,prob=mix)
  zs.actual[i] <- cl
  xs[i,] <- (runif(p) < mus.actual[cl,])
}
```

(c) Calculate the mixture parameters.

```R
source("em_mixture_bernoullis.R")
print(system.time(out <- em_mix_bernoulli(xs,K.actual)))
```

(d) Check if the learned parameters are close to the truth.

```R
v1 <- sum(abs(out$mus-mus.actual))
v2 <- sum(abs(out$mus[2:1,]-mus.actual))
```

```
vm <- min(v1,v2)/p/2
print(vm)
if (vm > .3) print("probably not working") else print("might be working")
```