

Methods of Maximizing the Likelihood

Maximum likelihood estimation requires maximization of the log likelihood $\ell(\boldsymbol{\theta}) = \log L(\boldsymbol{\theta}|\mathbf{Y})$.

In most cases, this means taking derivatives and solving likelihood equations

$$\underline{\underline{\ell}}(\boldsymbol{\theta}) = \frac{\partial}{\partial \boldsymbol{\theta}^T} \ell(\boldsymbol{\theta}) = 0.$$

Sometimes we can do this analytically (Yay!).

When an analytical solution doesn't exist, we have options:

- Standard optimization methods like Newton-Raphson (or fancy ones like gradient descent).
- Profile likelihood. (later).
- EM algorithm (today).

"Expectation"
"maximization"

1 EM Algorithm

Approach solving the likelihood equation via viewing the observed data \mathbf{Y} as incomplete and that there is missing data \mathbf{Z} that would make the problem simpler if we had it.

↳ sometimes it is actually missing data

↳ others just additional data we wish we had.

Intuition for
"wish we had"
is how we could
proceed.

Example (Two-Component Mixture): Suppose Y_1, \dots, Y_n are iid from the mixture density

$$f(y; \theta) = p f_1(y; \mu_1, \Sigma_1) + (1-p) f_2(y; \mu_2, \Sigma_2),$$

where f_1 and f_2 are bivariate normal densities with mean vectors μ_1 and μ_2 and variance matrices Σ_1 and Σ_2 , respectively. Thus, the parameter vector $\theta = (p, \mu_1, \mu_2, \Sigma_1, \Sigma_2)$ and the likelihood is

$$L(p, \mu_1, \mu_2, \Sigma_1, \Sigma_2) = \prod_{i=1}^n [p f_1(y_i; \mu_1, \Sigma_1) + (1-p) f_2(y_i; \mu_2, \Sigma_2)]$$

$$\Rightarrow \ell(p, \mu_1, \mu_2, \Sigma_1, \Sigma_2) = \sum_{i=1}^n \log \{ p f_1(y_i; \mu_1, \Sigma_1) + (1-p) f_2(y_i; \mu_2, \Sigma_2) \}$$

... and we're stuck.

We cannot get nice expressions for $\hat{\mu}_{k,MLE}$ or $\hat{\Sigma}_{k,MLE}$ $k=1,2$.

Actually this log likelihood has maxima on boundary of the parameter space \Rightarrow not well-behaved.

```

library(mvtnorm) ## multivariate normal

p = .6
mu1 <- c(0, 0)
sig1 <- matrix(c(1, 0, 0, 1), ncol = 2)
mu2 <- c(1.5, 1.5)
sig2 <- matrix(c(1, .6, .6, 1), ncol = 2)

## sample from the mixture
n <- 50
z <- rbinom(n, 1, p)

y1 <- rmvnorm(sum(z), mean = mu1, sigma = sig1)
y2 <- rmvnorm(n - sum(z), mean = mu2, sigma = sig2)
y <- matrix(NA, nrow = n, ncol = 2) ## observed data
y[z == 1, ] <- y1
y[z == 0, ] <- y2

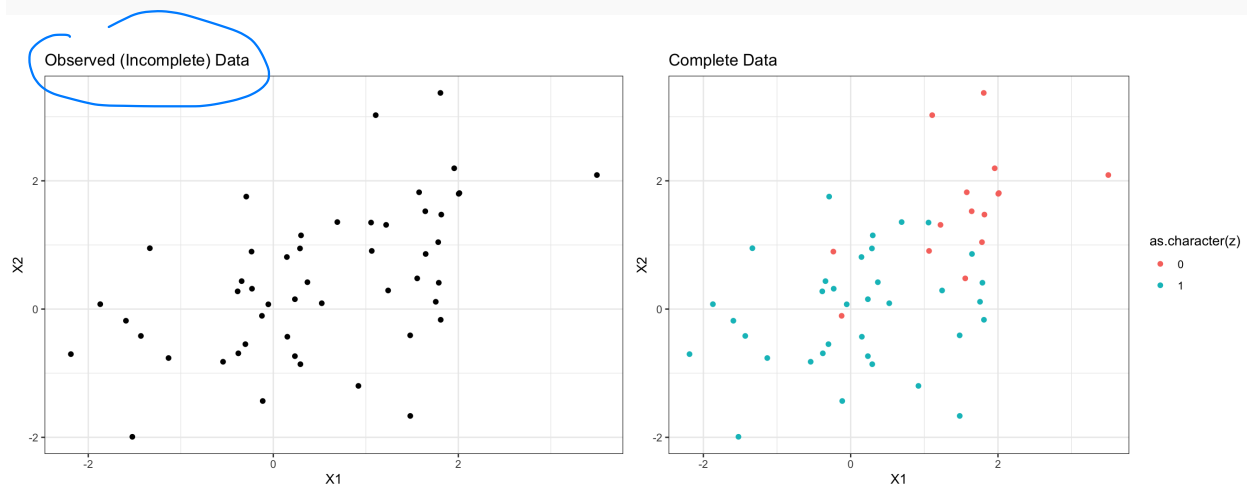
df <- data.frame(y, z)

## plot data
ggplot(df) +
  geom_point(aes(X1, X2)) +
  ggtitle("Observed (Incomplete) Data")

ggplot(df) +
  geom_point(aes(X1, X2, colour = as.character(z))) +
  ggtitle("Complete Data")

```

simulate
data.



Let's try to maximize the likelihood

```

# loglikelihood of incomplete data--no knowledge of z
loglik_mixture <- function(par, data) {
  p <- plogis(par[1]) # p guaranteed to be in [0,1]
  mu1 <- c(par[2], par[3])
  sig1 <- matrix(c(exp(par[4]), par[5], par[5], exp(par[4])), nrow
    = 2)
  mu2 <- c(par[6], par[7])
  sig2 <- matrix(c(exp(par[8]), par[9], par[9], exp(par[8])), nrow
    = 2)
  # note: exponential guarantees the diagonal elements are
    positive, but
  # nothing to guarantee matrices are positive definite. (Could do
    square root)

  out <- log(p * dmvnorm(data, mean = mu1, sigma = sig1) +
    (1-p) * dmvnorm(data, mean = mu2, sigma = sig2))
  return(sum(out))
}

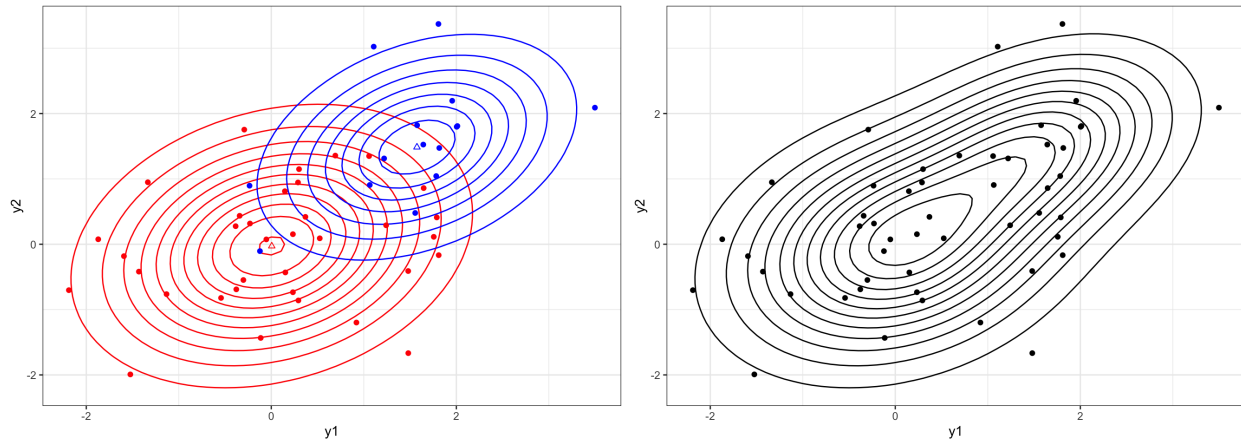
## optimize from different starting values
mle1 <- optim(c(0, -.2, -.2, .5, 0, 2, 2, .5, 0), loglik_mixture,
  data = y, control = list(fnscale = -1))
mle2 <- optim(c(.405, 0, 0, 0, 0, 1.5, 1.5, 0, .6), loglik_mixture,
  data = y, control = list(fnscale = -1))

```

Parameter	Truth	MLE1	MLE2
p	0.6	0.6771	0.6634
μ_{11}	0.0	0.0307	0.0050
μ_{12}	0.0	-0.0512	-0.0281
Σ_{111}	1.0	0.9757	0.9757
Σ_{112}	0.0	0.2178	0.2267
μ_{21}	1.5	1.5597	1.5744
μ_{22}	1.5	1.4815	1.4859
Σ_{211}	1.0	0.7161	0.7220
Σ_{212}	0.6	0.2679	0.2436

OK, not great not terrible.

Fitted results:



This seems pretty good... can we break this with initialization?

```
# Centered the second mixture component at a data point, and shrink
# variance, so normal is super-concentrated around that point.
loglik_mixture(c(.6, 0, 0, 0, 0, y[30, 1], y[30, 2], -50, 0), data =
  y)
```

```
## [1] -137.7964
```

```
mle3 <- optim(c(.6, 0, 0, 0, 0, y[30, 1], y[30, 2], -50, 0),
  loglik_mixture, data = y, control = list(fnscale = -1))
```

Parameter	Truth	MLE3
p	0.6	0.9873 <i>Yikes.</i>
μ_{11}	0.0	0.0000
μ_{12}	0.0	0.0000
Σ_{111}	1.0	1.0000
Σ_{112}	0.0	0.0000
μ_{21}	1.5	1.8067
μ_{22}	1.5	3.3712
Σ_{211}	1.0	0.0000 <i>← that's bad.</i>
Σ_{212}	0.6	0.0000

What would change if we were given the complete data, where $Z_i \stackrel{iid}{\sim} \text{Bern}(p)$?
now we know cluster assignment!

$$f_{Y,Z}(y,z;\theta) = (p f_1(y; \mu_1, \Sigma_1))^z ((1-p) f_2(y; \mu_2, \Sigma_2))^{1-z}$$

$$\Rightarrow \ell(p, \mu_1, \mu_2, \Sigma_1, \Sigma_2 | Y, Z) = \sum_{i=1}^n \{ z_i \log f_1(y_i; \mu_1, \Sigma_1) + (1-z_i) \log f_2(y_i; \mu_2, \Sigma_2) + z_i \log p + (1-z_i) \log(1-p) \}$$

$$\frac{\partial \ell(\theta | Y, Z)}{\partial \mu_1} = \sum_{i=1}^n z_i \frac{\partial \log f_1(y_i; \mu_1, \Sigma_1)}{\partial \mu_1} \quad \log f_1(y_i; \mu_1, \Sigma_1) = -\log 2\pi - \frac{1}{2} \log \det \Sigma_1 - \frac{1}{2} (y_i - \mu_1)^T \Sigma_1^{-1} (y_i - \mu_1)$$

$$\Rightarrow \frac{\partial \log f_1(y_i; \mu_1, \Sigma_1)}{\partial \mu_1} = -\Sigma_1^{-1} (y_i - \mu_1)$$

plugging in:

$$\frac{\partial \ell(\theta | Y, Z)}{\partial \mu_1} = -\sum_{i=1}^n z_i \Sigma_1^{-1} (y_i - \mu_1) \stackrel{\text{set}}{=} 0 \quad \text{solve} \Rightarrow \hat{\mu}_{1, \text{MLE}} = \frac{1}{n_{\{Z_i=1\}}} \sum_{i=1}^n z_i y_i$$

So MLE is the sample mean of the observations from the first density (DMM).

The other gaussian parameter estimates are also exactly what you think:

$$\hat{\mu}_{2, \text{MLE}} = \frac{1}{n_{\{Z_i=0\}}} \sum_{i=1}^n (1-z_i) y_i, \quad \hat{\Sigma}_{1, \text{MLE}} = \frac{1}{n_{\{Z_i=1\}}} \sum_{i=1}^n z_i (y_i - \hat{\mu}_{1, \text{MLE}})^T (y_i - \hat{\mu}_{1, \text{MLE}}), \quad \text{similar for } \hat{\Sigma}_{2, \text{MLE}}.$$

Now p :

$$\frac{\partial \ell(\theta | Y, Z)}{\partial p} = \frac{1}{p} \sum_{i=1}^n z_i - \frac{1}{(1-p)} \sum_{i=1}^n (1-z_i) \stackrel{\text{set}}{=} 0$$

$$\sum_{i=1}^n z_i - p \sum_{i=1}^n z_i = np - p \sum_{i=1}^n z_i \Rightarrow \hat{p}_{\text{MLE}} = \frac{\sum_{i=1}^n z_i}{n}$$

also exactly what we expect!

So, if we knew which mixture component the data came from, our life would be easy...

Consider the complete log-likelihood:

$$\ell(p, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \Sigma_1, \Sigma_2 | \mathbf{Y}, \mathbf{Z}) = \sum_{i=1}^n \{ Z_i \log f_1(Y_i; \boldsymbol{\mu}_1, \Sigma_1) + (1 - Z_i) \log f_2(Y_i; \boldsymbol{\mu}_2, \Sigma_2) + Z_i \log p + (1 - Z_i) \log(1 - p) \}.$$

We could consider the Z_i 's as "weights" which represent our *current* belief in which density each datum come from.

If we were given this information (the complete data) our belief weights are 0s and 1s.

Instead have current belief based on model parameters.

↙ based on current knowledge of model parameters

Given what our belief is in the weights of the data, what is our estimate of the model parameters?

This seems circular. (and it is \rightarrow iterative procedure).

$$\hat{\boldsymbol{\theta}}^{(k)} \xrightarrow{\text{belief weights}} \hat{\boldsymbol{\theta}}^{(k+1)}$$

"E"
"M"

$$\hat{\mu}_1^{(k+1)} = \frac{\sum_{i=1}^n w_i^{(k)} \psi_i}{\sum_{i=1}^n w_i^{(k)}}$$

$$\hat{\mu}_2^{(k+1)} = \frac{\sum_{i=1}^n (1 - w_i^{(k)}) \psi_i}{\sum_{i=1}^n (1 - w_i^{(k)})}$$

$$\hat{\sigma}_1^{(k+1)} = \frac{1}{\sum_{i=1}^n w_i^{(k)}} \sum_{i=1}^n w_i^{(k)} (\psi_i - \hat{\mu}_1^{(k)}) (\psi_i - \hat{\mu}_1^{(k)})$$

$$\hat{\sigma}_2^{(k+1)} = \frac{1}{\sum_{i=1}^n (1 - w_i^{(k)})} \sum_{i=1}^n (1 - w_i^{(k)}) (\psi_i - \hat{\mu}_2^{(k)}) (\psi_i - \hat{\mu}_2^{(k)})$$

$$\hat{p}^{(k+1)} = \frac{\sum_{i=1}^n w_i^{(k)}}{n}$$

This is the basic intuition for the EM algorithm. We will view our data \mathbf{Y} as incomplete and imagine there is missing data \mathbf{Z} that would make the problem simpler if we had it.

The EM algorithm then follows:

① Write down joint likelihood of the "complete" data (\mathbf{y}, \mathbf{z}) , $L_c(\underline{\theta} | \mathbf{y}, \mathbf{z})$ *don't actually have \mathbf{z} !*
need to maximize something that is only a function of $\underline{\theta}, \mathbf{y}$.
 \Rightarrow conditional

② E-step: compute conditional expectation of $\log L_c(\underline{\theta} | \mathbf{y}, \mathbf{z})$ given \mathbf{y} , assuming parameter is $\underline{\theta}^{(v)}$ *expectation!*
 $Q(\underline{\theta}, \underline{\theta}^{(v)} | \mathbf{y}) = E_{\underline{\theta}^{(v)}} [\log L_c(\underline{\theta} | \mathbf{y}, \mathbf{z}) | \mathbf{y}] = \int \log L_c(\underline{\theta} | \mathbf{y}, \mathbf{z}) f_{\mathbf{z} | \mathbf{y}}(\mathbf{z} | \mathbf{y}, \underline{\theta}^{(v)}) d\mathbf{z}$ *actually current value in the iteration.*

③ M-step: Maximize $Q(\underline{\theta}, \underline{\theta}^{(v)} | \mathbf{y})$ wrt $\underline{\theta}$ ($\underline{\theta}^{(v)}$ fixed).

$$\text{i.e. } \underline{\theta}^{(v+1)} = \underset{\underline{\theta}}{\operatorname{argmax}} Q(\underline{\theta}, \underline{\theta}^{(v)} | \mathbf{y}).$$

repeat ② & ③ until convergence (values of $\underline{\theta}^{(v)}$ and $\underline{\theta}^{(v+1)}$ not changing much).

Example (Two-Component Mixture, Cont'd): The EM algorithm for the two-component Gaussian mixture model is

(initialized) start with $\hat{\underline{\theta}}^{(0)}$, for $v = 0, 1, 2, \dots$

$$\text{① E step: } Q(\underline{\theta}, \hat{\underline{\theta}}^{(v)} | \mathbf{y}) = E_{\hat{\underline{\theta}}^{(v)}} [\log L_c(\underline{\theta} | \mathbf{y}, \mathbf{z}) | \mathbf{y}] = \sum_{i=1}^n \{ w_i^{(v)} \log f_1(y_i; \mu_1^{(v)}, \Sigma_1^{(v)}) + (1-w_i^{(v)}) \log f_2(y_i; \mu_2^{(v)}, \Sigma_2^{(v)}) + w_i^{(v)} \log p^{(v)} + (1-w_i^{(v)}) \log (1-p^{(v)}) \}$$

$$w_i^{(v)} = E_{\hat{\underline{\theta}}^{(v)}}(z_i | y_i) = \frac{p^{(v)} f_1(y_i; \mu_1^{(v)}, \Sigma_1^{(v)})}{p^{(v)} f_1(y_i; \mu_1^{(v)}, \Sigma_1^{(v)}) + (1-p^{(v)}) f_2(y_i; \mu_2^{(v)}, \Sigma_2^{(v)})}$$

② M-step: see page 7.

Your Turn: Implement the EM algorithm for the two-component mixture model on our example data.

1.1 Convergence of the EM algorithm

We will show that $l(\hat{\theta}^{(v+1)}) \geq l(\hat{\theta}^{(v)})$.

In other words, each step of the EM algorithm leads to an improvement in the log-likelihood value.

If the likelihood is well-behaved, it will achieve the MLE, otherwise it will achieve a local maximum (if there is one).
 ↳ banded, unimodal

y = observed data
 z = hidden

We know $f_{Z|Y}(z|y; \theta) = \frac{f_{YZ}(y, z; \theta)}{f_Y(y; \theta)}$. def'n of conditional density. true for any y, z

$\Rightarrow f_y(y_j; \theta) = \frac{f_{yz}(y, z; \theta)}{f_{z|y}(z|y; \theta)}$ just rewritten (not clear why yet).

Assume we observe $\mathbf{y} = (y_1, \dots, y_n)$, then

$L(\theta | \mathbf{y}) = f_y(\mathbf{y}; \theta) = \frac{f_{yz}(\mathbf{y}, \mathbf{z}; \theta)}{f_{z|y}(\mathbf{z} | \mathbf{y}; \theta)}$ (if iid, product of univariate densities)

$\Rightarrow l(\theta | \mathbf{y}) = \log f_{yz}(\mathbf{y}, \mathbf{z}; \theta) - \log f_{z|y}(\mathbf{z} | \mathbf{y}; \theta) = \underbrace{\log \text{likelihood of complete data } \mathbf{y}, \mathbf{z}}_{l_c(\theta | \mathbf{y}, \mathbf{z})} - \underbrace{l(\theta | \{\mathbf{z} | \mathbf{y}\})}_{\text{conditional log-likelihood}}$ holds for any $\mathbf{z}!$
 \Rightarrow take expected values wrt $\mathbf{z} | \mathbf{y}; \hat{\theta}^{(v)}$

$\int l(\theta | \mathbf{y}) f_{z|y}(\mathbf{z} | \mathbf{y}; \hat{\theta}^{(v)}) d\mathbf{z} = \int \log f_{yz}(\mathbf{y}, \mathbf{z}; \theta) f_{z|y}(\mathbf{z} | \mathbf{y}; \hat{\theta}^{(v)}) d\mathbf{z} - \int \log f_{z|y}(\mathbf{z} | \mathbf{y}; \theta) f_{z|y}(\mathbf{z} | \mathbf{y}; \hat{\theta}^{(v)}) d\mathbf{z}$

$l(\theta | \mathbf{y}) \int f_{z|y}(\mathbf{z} | \mathbf{y}; \hat{\theta}^{(v)}) d\mathbf{z} = \dots$
 $l(\theta | \mathbf{y}) = Q(\theta, \hat{\theta}^{(v)}) - H(\theta, \hat{\theta}^{(v)})$

So, in order to show that $l(\hat{\theta}^{(v+1)}) \geq l(\hat{\theta}^{(v)})$, this is the same as

$Q(\hat{\theta}^{(v+1)}, \hat{\theta}^{(v)}) - H(\hat{\theta}^{(v+1)}, \hat{\theta}^{(v)}) \geq Q(\hat{\theta}^{(v)}, \hat{\theta}^{(v)}) - H(\hat{\theta}^{(v)}, \hat{\theta}^{(v)})$

Step 1: Show that $H(\theta, \hat{\theta}^{(v)})$ is maximized when $\theta = \hat{\theta}^{(v)}$.

$$\text{i.e. } H(\hat{\theta}^{(v)}, \hat{\theta}^{(v)}) \geq H(\theta, \hat{\theta}^{(v)}) \text{ for any } \theta.$$

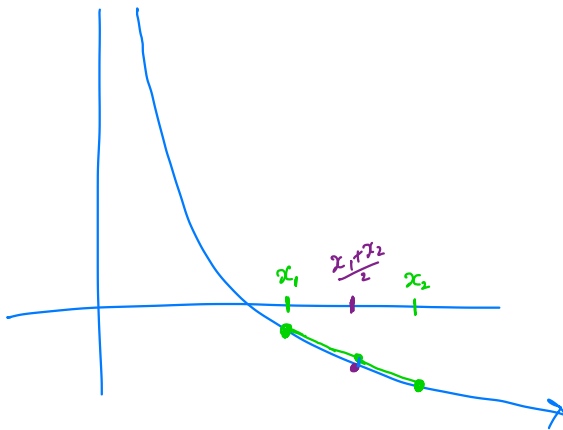
Recall: Jensen's Inequality. A function Φ is convex if $\Phi\left(\frac{x_1+x_2}{2}\right) \leq \frac{1}{2}\Phi(x_1) + \frac{1}{2}\Phi(x_2)$. Then

$$\Phi(\mathbb{E}[g(X)]) \leq \mathbb{E}[\Phi(g(X))], \Leftrightarrow \Phi\left(\int g(x)f(x)dx\right) \leq \int \Phi(g(x))f(x)dx$$

where g is a real-valued integrable function.

where $X \sim f$.

Fact: $-\log$ is convex



Consider $H(\hat{\theta}^{(v)}, \hat{\theta}^{(v)}) - H(\theta, \hat{\theta}^{(v)})$ WTS this is positive $\forall \theta$.

$$H(\theta, \hat{\theta}^{(v)}) = \int \log(f_{z|y}(z|y; \theta)) f_{z|y}(z|y; \hat{\theta}^{(v)}) dz$$

$$\Rightarrow H(\hat{\theta}^{(v)}, \hat{\theta}^{(v)}) - H(\theta, \hat{\theta}^{(v)}) = \int (\log f_{z|y}(z|y; \hat{\theta}^{(v)}) - \log f_{z|y}(z|y; \theta)) f_{z|y}(z|y; \hat{\theta}^{(v)}) dz$$

$$= \int -\log\left(\frac{f_{z|y}(z|y; \theta)}{f_{z|y}(z|y; \hat{\theta}^{(v)})}\right) f_{z|y}(z|y; \hat{\theta}^{(v)}) dz.$$

$$\geq -\log\left\{ \int \frac{f_{z|y}(z|y; \theta)}{f_{z|y}(z|y; \hat{\theta}^{(v)})} f_{z|y}(z|y; \hat{\theta}^{(v)}) dz \right\}$$

$$= -\log\left\{ \int f_{z|y}(z|y; \theta) dz \right\}$$

$$= 0$$

$$\Rightarrow H(\hat{\theta}^{(v)}, \hat{\theta}^{(v)}) \geq H(\theta, \hat{\theta}^{(v)}) \quad \forall \theta //$$

This allows us to only focus on maximizing Q when performing the EM algorithm!

Step 2: Find a $\hat{\theta}^{(u+1)}$ that will optimize Q .

Recall goal is to $\hat{\theta}^{(u+1)}$ s.t. $l(\hat{\theta}^{(u+1)} | y) \geq l(\hat{\theta}^{(u)} | y)$.

And $l(\theta | y) = Q(\theta, \hat{\theta}^{(u)}) - H(\theta, \hat{\theta}^{(u)})$.

Let $\hat{\theta}^{(u+1)} = \underset{\theta}{\operatorname{argmax}} Q(\theta, \hat{\theta}^{(u)})$.

We know $H(\hat{\theta}^{(u+1)}, \hat{\theta}^{(u)}) \leq H(\hat{\theta}^{(u)}, \hat{\theta}^{(u)})$ because true all θ .

+ $Q(\hat{\theta}^{(u+1)}, \hat{\theta}^{(u)}) \geq Q(\hat{\theta}^{(u)}, \hat{\theta}^{(u)})$ by optimization.

So $l(\hat{\theta}^{(u+1)} | y) = Q(\hat{\theta}^{(u+1)}, \hat{\theta}^{(u)}) - H(\hat{\theta}^{(u+1)}, \hat{\theta}^{(u)})$
 $\leq Q(\hat{\theta}^{(u)}, \hat{\theta}^{(u)}) - H(\hat{\theta}^{(u)}, \hat{\theta}^{(u)}) = l(\hat{\theta}^{(u)} | y)$ ✓

Example (Two-Component Mixture, Cont'd):

$$Q(\theta, \hat{\theta}^{(v)}) = \int \log f_{Y,Z}(y, z; \theta) f_{Z|Y}(z | y; \hat{\theta}^{(v)}) dz.$$

For the Gaussian mixture, we know the complete log-likelihood:

$$\log f_{Y,Z}(Y, Z; \theta) = \sum_{i=1}^n \{ z_i \log f_1(y_i; \mu_1, \Sigma_1) + (1-z_i) \log f_2(y_i; \mu_2, \Sigma_2) + z_i \log p + (1-z_i) \log(1-p) \}$$

To get the conditional density,

$$f_{Z|Y}(z | y; \hat{\theta}^{(v)}) = \prod_{i=1}^n f_{Z_i|Y_i}(z_i | y_i; \hat{\theta}^{(v)}).$$

$$f_{Z_i|Y_i}(z_i | y_i; \hat{\theta}^{(v)}) = \frac{f_{Y,Z}(y_i, z_i; \hat{\theta}^{(v)})}{f_Y(y_i; \hat{\theta}^{(v)})} \quad \begin{array}{l} \text{complete density} \\ \text{observed density.} \end{array}$$

$$= \frac{[\hat{p}^{(v)} f_1(y_i; \hat{\mu}_1^{(v)}, \hat{\Sigma}_1^{(v)})]^{z_i} [(1-\hat{p}^{(v)}) f_2(y_i; \hat{\mu}_2^{(v)}, \hat{\Sigma}_2^{(v)})]^{1-z_i}}{\hat{p}^{(v)} f_1(y_i; \hat{\mu}_1^{(v)}, \hat{\Sigma}_1^{(v)}) + (1-\hat{p}^{(v)}) f_2(y_i; \hat{\mu}_2^{(v)}, \hat{\Sigma}_2^{(v)})}$$

$$P(Z_i = 1 | Y = y_i) = \frac{\hat{p}^{(v)} f_1(y_i; \hat{\mu}_1^{(v)}, \hat{\Sigma}_1^{(v)})}{\hat{p}^{(v)} f_1(y_i; \hat{\mu}_1^{(v)}, \hat{\Sigma}_1^{(v)}) + (1-\hat{p}^{(v)}) f_2(y_i; \hat{\mu}_2^{(v)}, \hat{\Sigma}_2^{(v)})} \quad \stackrel{\text{define}}{=} \hat{w}_i^{(v)}; \text{ Then}$$

$$\Rightarrow Z_i | Y_i, \hat{\theta}^{(v)} \sim \text{Bern}(\hat{w}_i^{(v)}).$$

$$\Rightarrow Q(\theta, \hat{\theta}^{(v)}) = \sum_{i=1}^n E_{Z_i|Y_i} \left[\log f_{Y,Z}(Y_i, Z_i; \theta) \right] \quad \text{and} \quad E_{Z_i|Y_i} [g(Z_i)] = g(1) \hat{w}_i^{(v)} + g(0) (1 - \hat{w}_i^{(v)}).$$

$$= \sum_{i=1}^n \hat{w}_i^{(v)} \left[\log f_1(Y_i; \mu_1, \Sigma_1) + \log p \right] + (1 - \hat{w}_i^{(v)}) \left[\log f_2(Y_i; \mu_2, \Sigma_2) + \log(1-p) \right]$$

take derivatives,
set = 0, solve.

↳ Which yields the intuitive update formulas from before!

So "plugging in the weights" makes sense from an optimization standpoint in this example.

In general can't always separate E + M in this way for Q.

The EM algorithm allows us to obtain $\hat{\theta}_{EM}$, the parameter estimate which optimizes the algorithm.

If likelihood is "nice" = $\hat{\theta}_{MLE}$

1.2 Variance Estimation for EM estimates

The EM algorithm find the MLE, but it does not automatically produce an estimate of the covariance matrix. Why not?
→ likelihood is well-behaved.

Could'n't we just look at the curvature of the optimized surface?

No. We are optimizing Q , which is not the log-likelihood! ($Q - H$).

⇒ we can't rely on information results without computing the Hessian of the log-likelihood. (Which may be hard)

There are several options to estimate the variance.

1. **Bootstrapping** (we will discuss bootstrapping in depth later in the course).

A simple procedure is easy to imagine:

1. Find $\hat{\theta}_{EM}$ from the original data $\mathcal{Y} = \mathcal{Y}_1, \dots, \mathcal{Y}_n$

2. Resample data with replacement to obtain bootstrap samples (sample from empirical dsn).

$$\mathcal{Y}^{*b} = \mathcal{Y}_{-i}^{*b}, \dots, \mathcal{Y}_{-n}^{*b} \quad \text{for } b = 1, \dots, B$$

← # of bootstrap samples

Find $\hat{\theta}_{EM}^{*b}$ for each \mathcal{Y}^{*b} .

3. Use $\hat{\theta}_{EM}^{*b}$'s to get bootstrap CI for $\hat{\theta}_{EM}$.

↳ percentile bootstrap.

↳ t-based

↳ $B \subset A$ bootstrap.

Straight forward, easy.

Can be computationally expensive (need EM $b=1, \dots, B$ times).

↳ could be done in parallel to ease the burden.

2. Louis's Method

$$l(\underline{\theta}|\underline{y}) = Q(\underline{\theta}, \underline{\theta}^{(v)}) - H(\underline{\theta}, \underline{\theta}^{(v)}) \quad \text{true for any } \underline{\theta}^{(v)}$$

$$\text{So in particular } l(\underline{\theta}|\underline{y}) = Q(\underline{\theta}, \underline{\theta}) - H(\underline{\theta}, \underline{\theta})$$

$$\Rightarrow -l''(\underline{\theta}|\underline{y}) = -Q''(\underline{\theta}, \underline{\theta}) + H''(\underline{\theta}, \underline{\theta}). \quad (\text{derivatives wrt to } \underline{\theta} \text{ first argument})$$

$$\uparrow$$

$$n\bar{I}(\underline{y}, \underline{\theta})$$

from before.

observed information = "complete information" - "missing information"

$$-Q''(\underline{\theta}, \underline{\theta}) = -E_{\underline{z}|\underline{y}} [l''(\underline{\theta}|\underline{y}, \underline{z})] = -\int l''(\underline{\theta}|\underline{y}, \underline{z}) f_{\underline{z}|\underline{y}}(\underline{z}|\underline{y}; \underline{\theta}) d\underline{z}$$

Somewhat similar in form to Fisher Information (except for ds'n expectation is wrt).
similar statement about $H''(\underline{\theta}, \underline{\theta})$.

So to get CIs, need $Q''(\underline{\theta}, \underline{\theta})$ and $H''(\underline{\theta}, \underline{\theta})$:

$Q''(\underline{\theta}, \underline{\theta})$ is the curvature of the optimization surface \Rightarrow can often get numerically.

$$H''(\underline{\theta}, \underline{\theta}) = \text{Var} \left[\frac{\partial \log f_{\underline{z}|\underline{y}}(\underline{z}|\underline{y}; \underline{\theta})}{\partial \underline{\theta}} \right] \text{ wrt } f_{\underline{z}|\underline{y}} \quad (\text{again use fact (2)})$$

$$\text{because } H''(\underline{\theta}, \underline{\theta}) = \int \frac{\partial^2 \log f_{\underline{z}|\underline{y}}(\underline{z}|\underline{y}; \underline{\theta})}{\partial \underline{\theta}} f_{\underline{z}|\underline{y}}(\underline{z}|\underline{y}; \underline{\theta}) d\underline{z}$$

A MC estimate of $H''(\underline{\theta}, \underline{\theta})$ is thus the sample variance of $\frac{\partial \log f_{\underline{z}|\underline{y}}(\underline{z}|\underline{y}; \underline{\theta})}{\partial \underline{\theta}}$

of a sample of \underline{z} 's imputed from $f_{\underline{z}|\underline{y}}$

Other options see Guinness & Glesby Section 4.2.3

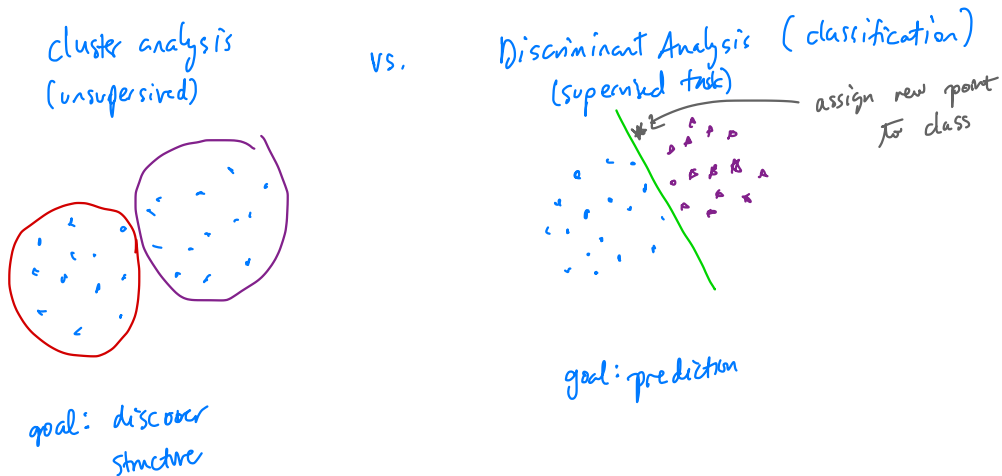
- Supplemental EM
- Empirical Information
- Numerical differentiation to get Hessian.

assuming integration
and differentiation order
are interchangeable

1.3 Another way to cluster: K-means

Goal of clustering:

Find an optimal grouping for which the observations within each group are "similar" but clusters are "dissimilar" to each other.



Methods for clustering include hierarchical and non-hierarchical, algorithmic and model-based.

hierarchical: proceed by going from $n \rightarrow 1$ clusters (or $1 \rightarrow n$)

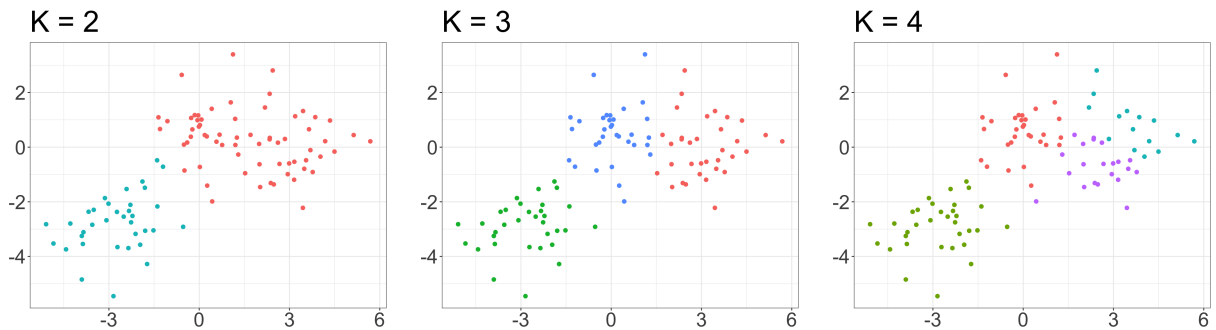
non-hierarchical: model-based (finite mixture of Normals) } both work for set # of clusters
 algorithmic (k-means) } but can be applied to different #s.

Not an exhaustive list.

K-means is a simple and elegant approach to partition a data set into K distinct, non-overlapping clusters.

First specify how many clusters (K), then K-means assigns each observation to one of K clusters.

Eg. clustering $n=100$ observations into K clusters using $p=2$ features.



The K -means clustering procedure results from a simple and intuitive mathematical problem. Let C_1, \dots, C_K denote sets containing the indices of observations in each cluster. These satisfy two properties:

e.g. if obs i is in cluster k , $i \in C_k$

$$1. C_1 \cup C_2 \cup \dots \cup C_K = \{1, \dots, n\}$$

each obs. belongs to a cluster

$$2. C_k \cap C_{k'} = \emptyset \quad \forall k \neq k'$$

clusters are non-overlapping.

Idea: "Good" clustering is one for which within-cluster-variation is as small as possible.
 \hookrightarrow clustered data are "similar"

The within-cluster variation for cluster C_k is a measure of the amount by which the observations within a cluster differ from each other. call this $W(C_k)$.

The we want to

$$\text{minimize } \left\{ \sum_{k=1}^K W(C_k) \right\}$$

C_1, \dots, C_K

i.e. partition data into K clusters s.t. total within cluster variation is minimized.

To solve this, we need to define within-cluster variation.

many ways.

Most common: squared Euclidean distance.

$$W(C_k) = \frac{1}{|C_k|} \sum_{i, i' \in C_k} \sum_{j=1}^p (x_{ij} - x_{i'j})^2$$

This results in the following optimization problem that defines K-means clustering:

$$\text{minimize } \left\{ \sum_{k=1}^K \frac{1}{|C_k|} \sum_{i, i' \in C_k} \sum_{j=1}^p (x_{ij} - x_{i'j})^2 \right\}$$

objective function.

this is very hard to solve exactly, $\approx K^n$ ways to partition n obs into K clusters!

A very simple algorithm has been shown to find a local optimum to this problem:

↳ "pretty good solution"

1. Randomly assign a number from 1 to K to each observation (initialize cluster assignments).

2. Iterate until cluster assignments stop changing:

(a) for each of the K clusters, compute cluster centroid

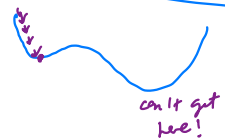
→ vector of the p feature means for each cluster.

(b) Assign each obs to closest centroid cluster.

↑ euclidean distance.

This algorithm is guaranteed to decrease the value of objective function at each step.

Clustering depends on initial (random) cluster assignment.



⇒ run the algorithm multiple times from different initial configurations and choose clustering w/ smallest objective function.

Still need K ...

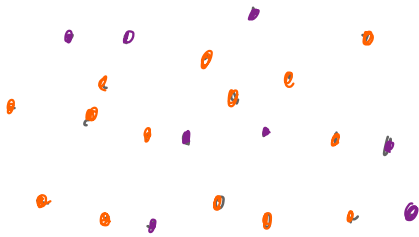
k-means illustration:

①

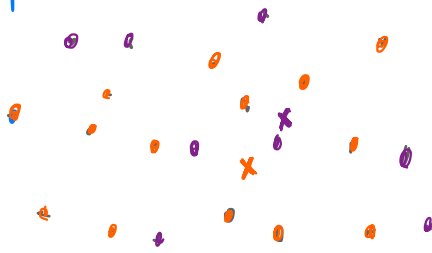


k=2

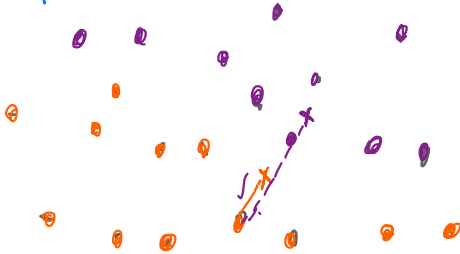
② Randomly assign



③ compute medoids.



④ Update assignments



repeat ③ & ④ until clusters stable.

Questions about the algorithm:

1. How do we define distance?

Normally Euclidean.

Could choose

- Minkowski $d(x_1, x_2) = \left(\sum_{i=1}^p |x_{1i} - x_{2i}|^d \right)^{1/d}$

= Mahalanobis

2. How do we choose starting values?

randomly usually (shotgun approach).

Maybe using another method? like hierarchical?

3. How do we choose K ?

look at between SS vs. within SS? (elbow plot of ratios)

Another way: Dunn index \rightarrow compares to a "null" clustering.

No one right way.

Compared to the Gaussian mixture problem,

finite.



also fixed # clusters

"soft assignment"

"hard assignment"

In the Gaussian mixture, we get probabilities of assignment, not just assignment.



model-based (making assumptions!)

K-means sensitive to starting values.

In R, kmeans function will fit this algorithm.

4. Should we scale the data?
Most times yes, unless no.