# UiO : Matematisk institutt

## Det matematisk-naturvitenskapelige fakultet

**STK-4051/9051  Computational Statistics  Spring 2022**
**Chaper 4 (part 1)**

Instructor: Odd Kolbjørnsen, oddkol@math.uio.no

# Missing data

- Data often (partly) missing
- Censored data (ex 2.3): Time to event not completely known
- Classification of images: Classes to some pixels known, unknown for most of the pixels
- Clustering: Data to be allocated to groups, group membership unknown
- If complete data, Likelihood "often easy"
- Likelihood becomes complicated when data are missing
- Notation:
    - $Y = (X, Z)$ are complete data
    - $X$ observed,
    - $Z$ missing
    - $X = M(Y)$ is observed part
    - Have $f_Y(y|\theta)$
    - Want $\max_{\theta} f_X(x|\theta)$

$$f_X(x|\theta) = \int_{y:M(y)=x} f_Y(y|\theta)\,dy = \int_z f_Y(x,z|\theta)\,dz$$

$$f_X(x|\theta) = \frac{f_Y(y|\theta)}{f_{z|x}(z|x,\theta)}$$

# EM algorithm

- Main idea: Iterate between
  - Estimate $Z$ given $X, \theta$ (E-step)
  - Estimate $\theta$ given $(X, Z)$ (M-step)
- Formally a bit more complicated
  - If complete data, we want to maximize $\log L(\theta|Y)$
  - $\log L(\theta|Y)$ unknown, but given a current value $\theta^{(t)}$ we can estimate it by

$$Q\big(\boldsymbol{\theta}\big|\boldsymbol{\theta}^{(t)}\big) = E\big[\log L(\boldsymbol{\theta}|\boldsymbol{Y}) \mid \boldsymbol{x}, \boldsymbol{\theta}^{(t)}\big]$$
$$= E\big[\log f_Y(y|\theta)\big|\boldsymbol{x}, \boldsymbol{\theta}^{(t)}\big]$$
$$= \int_z \log f_Y(y|\theta)]\, f_{z|x}(z|x, \theta^t)dz$$

- Algorithm:
  1. E-step: Compute $Q\big(\boldsymbol{\theta}\big|\boldsymbol{\theta}^{(t)}\big)$
  2. M-step: Maximize $Q\big(\boldsymbol{\theta}\big|\boldsymbol{\theta}^{(t)}\big)$ wrt $\boldsymbol{\theta}$ to obtain $\boldsymbol{\theta}^{(t+1)}$.
  3. Return to E-step unless a stopping criterion has been met

# The Q function

- $Q\big(\boldsymbol{\theta}\big|\boldsymbol{\theta}^{(t)}\big) = E\big[\log L(\boldsymbol{\theta}|\boldsymbol{Y}) \quad \big| \, \boldsymbol{x}, \boldsymbol{\theta}^{(t)} \quad \big]$

$$= E\big[\log f_Y(\boldsymbol{y}|\boldsymbol{\theta})\big|\boldsymbol{x}, \boldsymbol{\theta}^{(t)}\big]$$

$$= \int_Z \log f_Y(\boldsymbol{y}|\boldsymbol{\theta})]\, f_{Z|x}(\boldsymbol{z}|\boldsymbol{x}, \boldsymbol{\theta}^{(t)})d\boldsymbol{z}$$

The expected value of the complete log likelihood (contains $\boldsymbol{\theta}$ )
given the observed data and the curent estimate of parameter ($\boldsymbol{\theta}^{(t)}$)

# Peppered Moths - Example

- Color based on one single gene
- Three different allels (C,I,T)
- C is dominant to I, and I is dominant to T

| | |
|---|---|
| TT | Light-colored |
| II,IT | Intermediate |
| CC,CI,CT | Black coloring |

- Observing color, interest in frequency of $C, I, T$

- 

- Assume frequencies $p_C, p_I, p_T, p_C + p_I + p_T = 1$

| Color | Probability |
|---|---|
| White | $p_T^2$ |
| Intermediate | $p_I^2 + 2p_I p_T$ |
| Black | $p_C^2 + 2p_C p_I + 2p_C p_T$ |

- Observed $(n_T, n_I, n_C) = (341, 196, 85)$
- Complete $(n_{CC}, n_{CI}, n_{CT}, n_{II}, n_{IT}, n_{TT})$

| Options | Count |
|---|---|
| CC | 1 |
| CI (IC) | 2 |
| CT (TC) | 2 |
| II | 1 |
| IT (TI) | 2 |
| TT | 1 |

# Peppered Moths - Likelihood

- Complete data $(n_{CC}, n_{CI}, n_{CT}, n_{II}, n_{IT}, n_{TT})$
- Complete likelihood (multinomial distribution)

$$f_{\mathbf{Y}}(\mathbf{y}|\mathbf{p}) = \frac{n!}{n_{CC}!n_{CI}!n_{CT}!n_{II}!n_{IT}!n_{TT}!} p_C^{2n_{CC}} (2p_C p_I)^{n_{CI}} (2p_C p_T)^{n_{CT}} p_I^{2n_{II}} (2p_I p_T)^{n_{IT}} p_T^{2n_{TT}}$$

$$= \frac{n!}{n_{CC}!n_{CI}!n_{CT}!n_{II}!n_{IT}!n_{TT}!} 2^{n_{CI}+n_{CT}+n_{IT}} \times$$

$$p_C^{2n_{CC}+n_{CI}+n_{CT}} p_I^{2n_{II}+n_{CI}+n_{IT}} p_T^{2n_{TT}+n_{CT}+n_{IT}}$$

- Complete log-likelihood

$$\log\{f_{\mathbf{Y}}(\mathbf{y}|\mathbf{p})\} = \log\left(\frac{n!}{n_{CC}!n_{CI}!n_{CT}!n_{II}!n_{IT}!n_{TT}!}\right) +$$

$$[n_{CI} + n_{CT} + n_{IT}]\log(2) + [2n_{CC} + n_{CI} + n_{CT}]\log(p_C) +$$

$$[2n_{II} + n_{CI} + n_{IT}]\log(p_I) + [2n_{TT} + n_{CT} + n_{IT}]\log(p_T)$$

- $Q(\mathbf{p}|\mathbf{p}^{(t)}) = E[\log\{f_{\mathbf{Y}}(\mathbf{y}|\mathbf{p})\}|n_C, n_I, n_T, \mathbf{p}^{(t)}]$
- Note: First term do not depend on $\mathbf{p} = (p_C, p_I, p_T)$, not needed in the optimization step!

# **Peppered Moths – updating E & M**

- Complete log-likelihood

$$Q(\mathbf{p}|\mathbf{p}^{(t)}) = \text{Const} + E[n_{CI} + n_{CT} + n_{IT}|\mathbf{p}^{(t)}]\log(2)+$$
$$E[2n_{CC} + n_{CI} + n_{CT}|\mathbf{p}^{(t)}]\log(p_C)+$$
$$E[2n_{II} + n_{CI} + n_{IT}|\mathbf{p}^{(t)}]\log(p_I)+$$
$$E[2n_{TT} + n_{CT} + n_{IT}|\mathbf{p}^{(t)}]\log(p_T)$$

$$\mathrm{E}[N_{CC}|n_C, n_I, n_T, \mathbf{p}^{(t)}] = n_{CC}^{(t)} = \frac{n_C(p_C^{(t)})^2}{(p_C^{(t)})^2 + 2p_C^{(t)}p_I^{(t)} + 2p_C^{(t)}p_T^{(t)}}$$

Expectation

- Updating:

$$p_C^{(t+1)} = \frac{2n_{CC}^{(t)} + n_{CI}^{(t)} + n_{CT}^{(t)}}{2n}$$
$$p_T^{(t+1)} = \frac{2n_{TT}^{(t)} + n_{CT}^{(t)} + n_{IT}^{(t)}}{2n},$$

$$p_I^{(t+1)} = \frac{2n_{II}^{(t)} + n_{IT}^{(t)} + n_{CI}^{(t)}}{2n}$$

Maximization

- `Moth_EM.R`

8

$$E[N_{CC}|n_C, n_I, n_T, \boldsymbol{p}^{(t)}] = n_{CC}^{(t)} = \frac{n_C (p_C^{(t)})^2}{(p_C^{(t)})^2 + 2p_C^{(t)} p_I^{(t)} + 2p_C^{(t)} p_T^{(t)}}$$

X could be either C,T, or I

$$E\left(N_{CC}|n_c, n_I, n_T, \boldsymbol{p}^{(t)}\right) = n_c \cdot P(CC|\, CX, \boldsymbol{p}^{(t)})$$

$$= n_c \cdot \frac{P(CC \,\&\, CX|\boldsymbol{p}^{(t)})}{P(CX|\boldsymbol{p}^{(t)})}$$

$$= n_c \frac{P(CC|\boldsymbol{p}^{(t)})}{P(CC|\boldsymbol{p}^{(t)}) + P(CI|\boldsymbol{p}^{(t)}) + P(CT|\boldsymbol{p}^{(t)})}$$

$$P(CC) = p_c^2$$
$$P(CI) = 2p_c p_I$$
$$P(CI) = 2p_c p_T$$

Insert to get result

# Moths in R    Data = (85, 196, 341)

```
> show(c(p.old,l.old,NA))
[1] 0.3333333 0.3333333 0.3333333 0.0000000          NA
> more = TRUE
> while(more){
+     n = allele.e(x,p)
+     p = allele.m(x,n)
+     l = loglik(p,n)
+     more = abs(l-l.old)>eps
+     R = sum((p-p.old)^2)/sum(p.old^2)
+     more = R > eps
+     show(c(p,l,R))
+     l.old = l
+     p.old = p
+ }
[1]    0.08199357    0.23740622    0.68060021 -90.55303903    0.57890393
[1]    0.071248952   0.197869614   0.730881433 -68.467059735   0.007993122
[1]    7.085204e-02  1.903604e-01  7.387876e-01 -6.526257e+01  2.058264e-04
[1]    7.083746e-02  1.890227e-01  7.401398e-01 -6.474409e+01  6.163093e-06
[1]    7.083693e-02  1.887869e-01  7.403762e-01 -6.465487e+01  1.894317e-07
[1]    7.083691e-02  1.887454e-01  7.404177e-01 -6.463926e+01  5.851928e-09
>
> ## OUTPUT
> p    # FINAL ESTIMATE FOR ALLELE PROBABILITIES (p.c, p.i, p.t)
[1] 0.07083691 0.18874537 0.74041772
```

# Convergence EM

- Iterations increases log likelihood
  - Jensen's inequality $\qquad \ell(\boldsymbol{\theta}|\boldsymbol{x}) = \log f_X(\boldsymbol{x}|\boldsymbol{\theta})$
  - For convex $f(x)$, we have:

  $$f\big(E(X)\big) \leq E(f(X))$$

- Convergence order $\beta > 0$:

| How fast iteration $x^{(t)}$ approaches the true solution $x^*$ |
|---|

$$\epsilon^{(t)} = x^{(t)} - x^*$$
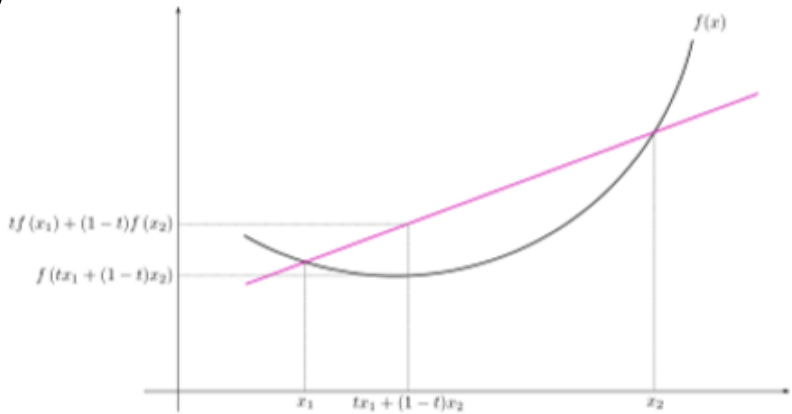
$$\lim_{t \to \infty} |\epsilon^{(t)}| \to 0$$

$$\lim_{t \to \infty} \frac{|\epsilon^{(t+1)}|}{|\epsilon^{(t)}|^{\beta}} = c$$

# Jensen's inequality



- Convex functions

$$f(tx_1 + (1-t)x_2) \leq tf(x_1) + (1-t)f(x_2)$$
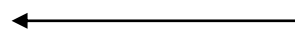
for $t \in [0, 1]$

- Finite form

$$f\left(\sum t_i x_i\right) \leq \sum t_i f(x_i), \quad t_i \text{ positive}, \sum t_i = 1$$

- Infinite form ($g(\cdot)$ non-negative, integrable):

$$f\left(\frac{1}{b-a}\int_a^b g(x)dx\right) \leq \frac{1}{b-a}\int_a^b f(g(x))dx$$

- Probabilistic form ($g(\cdot)$ density):

$$f(E^g[X]) \leq E^g[f(X)]$$

Ill prove this next week in exercise

# Iterations increase the value of, $\ell(\boldsymbol{\theta}|x) = \log(f_X(x|\theta))$

$$f_{z|x}(\mathbf{z}|\mathbf{x}, \theta) = \frac{f_y(\mathbf{y}|\boldsymbol{\theta})}{f_x(\mathbf{x}|\boldsymbol{\theta})}$$

$$\Downarrow$$

$$\log f_x(\mathbf{x}|\boldsymbol{\theta}) = \log f_y(\mathbf{y}|\boldsymbol{\theta}) - \log f_{z|x}(\mathbf{z}|\mathbf{x}, \theta)$$

$$\Downarrow$$

$$E[\log f_x(\mathbf{x}|\boldsymbol{\theta})] = E[\log f_y(\mathbf{y}|\boldsymbol{\theta})] - E[\log f_{z|x}(\mathbf{z}|\mathbf{x}, \theta)]$$

Any expectation

- If expectation with respect to $\mathbf{Z}|(\mathbf{x}, \boldsymbol{\theta}^{(t)})$,

$$\log f_x(\mathbf{x}|\boldsymbol{\theta}) = Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)}) - E[\log f_{z|x}(\mathbf{z}|\mathbf{x}, \theta)|\mathbf{x}, \boldsymbol{\theta}^{(t)}]$$

$$= Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)}) - H(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)})$$

$$Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)}) = E[\log f_Y(\boldsymbol{y}|\boldsymbol{\theta})|\, \boldsymbol{x}, \boldsymbol{\theta}^{(t)}\,]$$

$$H(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)}) = E[\log f_{z|x}(\boldsymbol{z}|\boldsymbol{x}, \boldsymbol{\theta})|\, \boldsymbol{x}, \boldsymbol{\theta}^{(t)}\,]$$

Select :
Expectation with
respect to the
distribution the
missing data  have
under the
under the current
estimate of the
parameter

**Proof:** $H\left(\theta^{(t)}\middle|\theta^{(t)}\right) \geq H\left(\theta\middle|\theta^{(t)}\right)$ **for any** $\theta$

$$H\left(\boldsymbol{\theta}\middle|\boldsymbol{\theta}^{(t)}\right) = E\left[\log f_{z|x}(\boldsymbol{z}|x,\theta)\middle|\, x, \theta^{(t)}\right]$$

$$H\left(\theta^{(t)}\middle|\theta^{(t)}\right) - H\left(\theta\middle|\theta^{(t)}\right) = E\left\{\log f_{z|x}\left(z\middle|x,\theta^{(t)}\right) - \log f_{z|x}(z|x,\theta)\right\}$$

$$= E\left\{-\log\frac{f_{z|x}(z|x,\theta)}{f_{z|x}\left(z\middle|x,\theta^{(t)}\right)}\right\} \quad \geq \quad -\log E\left\{\frac{f_{z|x}(z|x,\theta)}{f_{z|x}\left(z\middle|x,\theta^{(t)}\right)}\right\}$$

Jensen's

$$= -\log\int\frac{f_{z|x}(z|x,\theta)}{f_{z|x}\left(z\middle|x,\theta^{(t)}\right)}f_{z|x}\left(z\middle|x,\theta^{(t)}\right)dz$$

$$= -\log\int f_{z|x}(z|x,\theta)\,dz = -\log E\{1\} = 0$$

# Proof of increasing likelihood

$\log f_x(\mathbf{x}|\boldsymbol{\theta}^{(t+1)}) - \log f_x(\mathbf{x}|\boldsymbol{\theta}^{(t)})$

$= \underbrace{Q(\boldsymbol{\theta}^{(t+1)}, \boldsymbol{\theta}^{(t)}) - Q(\boldsymbol{\theta}^{(t)}, \boldsymbol{\theta}^{(t)})} - [\underbrace{H(\boldsymbol{\theta}^{(t+1)}, \boldsymbol{\theta}^{(t)}) - H(\boldsymbol{\theta}^{(t)}, \boldsymbol{\theta}^{(t)})}]$

In maximization step choose the $\theta^{(t+1)}$ such that it improves the old $\theta^{(t)}$
$> 0$.
If you are not able to improve Q, you have converged

Select $\theta = \theta^{(t+1)}$ and apply result from previous page. Result holds for any $\theta$ in particular for $\theta = \theta^{(t+1)}$
$\geq 0$

$$H(\theta^{(t)}|\theta^{(t)}) \geq H(\theta^{(t+1)}|\theta^{(t)})$$
$$\Updownarrow$$
$$-H(\theta^{(t+1)}|\theta^{(t)}) + H(\theta^{(t)}|\theta^{(t)}) \geq 0$$

$$\log f_x(x|\theta^{(t+1)}) > \log f_x(x|\theta^{(t)})$$

# Convergence order (good to know, but need not derive)

- The EM algorithm defines a mapping $\theta^{(t+1)} = \Psi(\theta^{(t)})$
- When the EM algorithm converges, $\hat{\theta} = \Psi(\hat{\theta})$
- Tayler expansion:

$$
\begin{aligned}
\varepsilon^{(t+1)} &\equiv \theta^{(t+1)} - \hat{\theta} \\
&= \Psi(\theta^{(t)}) - \Psi(\hat{\theta}) \\
&\approx \Psi(\theta^{(t)}) - [\Psi(\theta^{(t)}) + \Psi'(\theta^{(t)})(\hat{\theta} - \theta^{(t)})] \\
&= \Psi'(\theta^{(t)})(\theta^{(t)} - \hat{\theta}) \\
&= \Psi'(\theta^{(t)})\varepsilon^{(t)}
\end{aligned}
$$

- Convergence order $\beta$ if $\lim_{t\to\infty} \frac{\|\varepsilon^{(t+1)}\|}{|\varepsilon^{(t)}|^\beta} = \rho$

- $p = 1$: $\lim_{t\to\infty} \frac{|\varepsilon^{(t+1)}|}{|\varepsilon^{(t)}|} = \Psi'(\hat{\theta})$, linear convergence

- $p > 1$: Still linear if $-\ell''(\hat{\theta}|\mathbf{x})$ is positive definite
- (Newton's method has convergence order $\beta = 2$)

# Example: Mixture Gaussian clustering

- Assume $\mathbf{Y}_i = (X_i, C_i)$ are distributed according to

$$\Pr(C_i = k) = \pi_k, \quad k = 1, \ldots, K$$
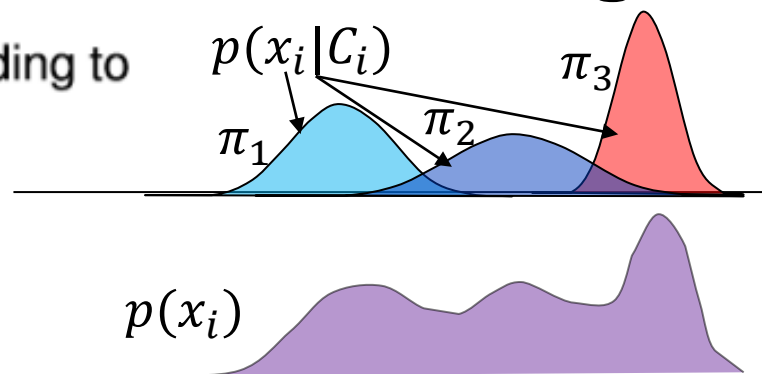
$$X_i | C_i = k \sim N(\mu_k, \sigma_k)$$

$p(x_i|C_i)$   $\pi_3$

$\pi_1$   $\pi_2$

$p(x_i)$

- The $C_i$'s are missing
- Complete log-density:

$$\log f(\mathbf{y}_i) = \log(\pi_{c_i}) + \log[\phi(x_i; \mu_{c_i}, \sigma_{c_i})]$$

$$= \sum_{k=1}^{K} I(c_i = k)[\log(\pi_k) + \log[\phi(x_i; \mu_k, \sigma_k)]]$$

- Complete log-likelihood:

$$\log f_Y(\mathbf{y}|\theta) = \sum_{i=1}^{n} \sum_{k=1}^{K} I(c_i = k)[\log(\pi_k) + \log[\phi(x_i; \mu_k, \sigma_k^2)]]$$

17

# E-step- Mixture Gaussian

- Complete log-likelihood:

$$\log f_Y(\mathbf{y}|\boldsymbol{\theta}) = \sum_{i=1}^{n} \sum_{k=1}^{K} I(c_i = k)[\log(\pi_k) + \log[\phi(x_i; \mu_k, \sigma_k^2)]]$$

- E-step (the $C_i$'s the only stochastic part)

$$Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)}) = E\left[\sum_{i=1}^{n} \sum_{k=1}^{K} I(C_i = k)[\log(\pi_k) + \log[\phi(x_i; \mu_k, \sigma_k^2)]]|\mathbf{x}, \boldsymbol{\theta}^{(t)}\right]$$

$$= \sum_{i=1}^{n} \sum_{k=1}^{K} E[I(C_i = k|\mathbf{x}, \boldsymbol{\theta}^{(t)})][\log(\pi_k) + \log[\phi(x_i; \mu_k, \sigma_k)]]$$

$$= \sum_{i=1}^{n} \sum_{k=1}^{K} \Pr(C_i = k|\mathbf{x}, \boldsymbol{\theta}^{(t)})[\log(\pi_k) + \log[\phi(x_i; \mu_k, \sigma_k)]]$$

$$\Pr(C_i = k|\boldsymbol{x}, \boldsymbol{\theta}^{(t)}) = \frac{\pi_k^{(t)}\phi(x_i, \mu_k^{(t)}, \sigma_k^{(t)})}{\sum_l \pi_l^{(t)}\phi(x_i, \mu_l^{(t)}, \sigma_l^{(t)})}$$

$$\Pr(C_i = k | \boldsymbol{x}, \boldsymbol{\theta}^{(t)}) = \frac{\pi_k^{(t)} \phi(x_i, \mu_k^{(t)}, \sigma_k^{(t)})}{\sum_l \pi_l^{(t)} \phi(x_i, \mu_l^{(t)}, \sigma_l^{(t)})}$$
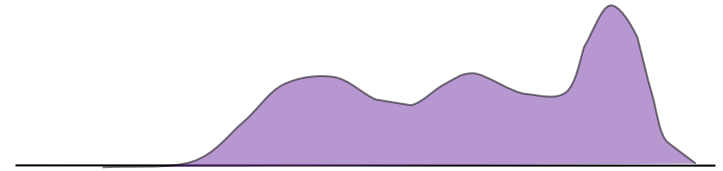
$$P(C_i = k) = \pi_k$$

$$p(x_i | C_i = k) = \phi(x_i; \mu_k, \sigma_k)$$

$p(x_i | C_i)$

$\pi_3$

$\pi_1$ $\pi_2$

$$p(x_i) = \sum_l \pi_l \phi(x_i; \mu_l, \sigma_l)$$

$$P(C_i = k \,| X_i = x_i) = \frac{p(C_i = k \,\&\, X_i = x_i)}{p(X_i = x_i)} = \frac{P(C_i = k)\, p(x_i | C_i = k)}{p(x_i)}$$

# M-step- Mixture Gaussian

- **M-step**: Taking into account $\sum_{k=1}^{K} \pi_k = 1$:

$$Q_{lagr}(\theta|\theta^{(t)}) = \sum_{i=1}^{n}\sum_{k=1}^{K} \Pr(C_i = k|\mathbf{x}, \theta^{(t)})[\log(\pi_k) + \log[\phi(x_i; \mu_k, \sigma_k^2)] + \lambda(1 - \sum_{k=1}^{K}\pi_k)$$

$$\frac{\partial}{\partial \pi_k} Q_{lagr}(\theta|\theta^{(t)}) = \sum_{i=1}^{n} \Pr(C_i = k|\mathbf{x}, \theta^{(t)})\pi_k^{-1} - \lambda$$

$$\Downarrow$$

$$\pi_k^{(t+1)} = \frac{\sum_{i=1}^{n}\Pr(C_i = k|\mathbf{x}, \theta^{(t)})}{\lambda}$$

$$= \frac{1}{n}\sum_{i=1}^{n}\Pr(C_i = k|\mathbf{x}, \theta^{(t)})$$

$$\sum_{k=1}^{K}\sum_{i=1}^{n}\frac{\Pr(C_i = k|x, \theta^{(t)})}{\lambda} = 1$$

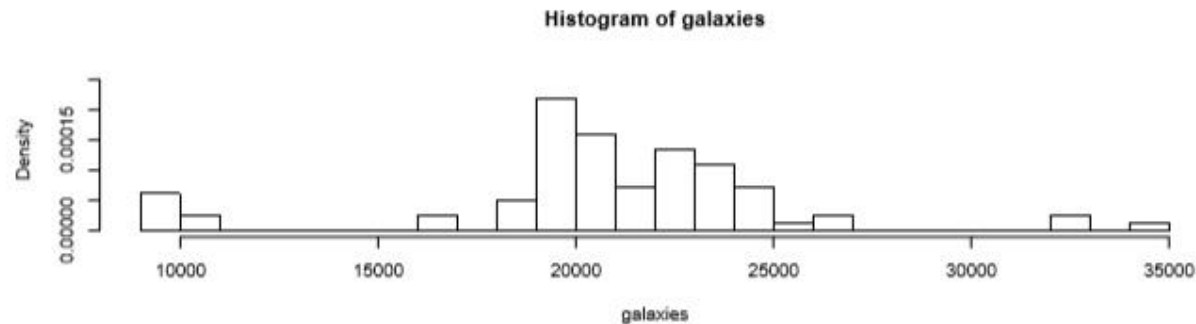$$\frac{1}{\lambda}\sum_{i=1}^{n}\underbrace{\sum_{k=1}^{K}\Pr(C_i = k|x, \theta^{(t)})}_{= 1} = 1$$

Similarly

$$\mu_k^{(t+1)} = \frac{1}{n\pi_k^{(t+1)}}\sum_{i=1}^{n}\Pr(C_i = k|\mathbf{x}, \theta^{(t)})x_i$$
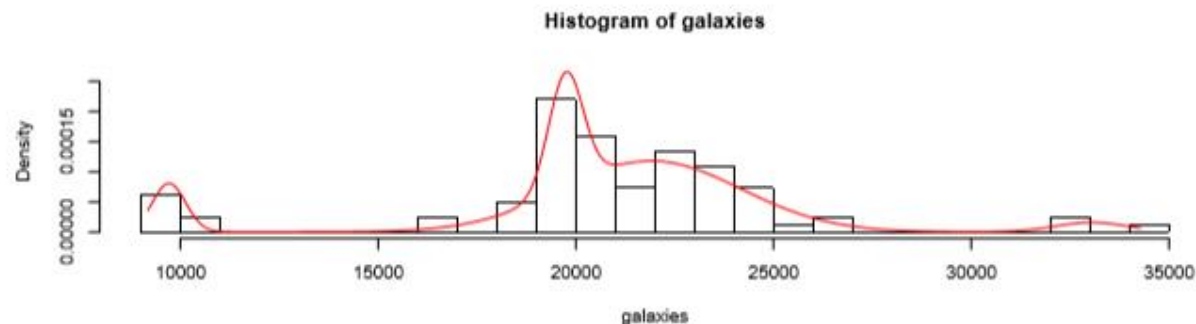
$$(\sigma_k^2)^{(t+1)} = \frac{1}{n\pi_k^{(t+1)}}\sum_{i=1}^{n}\Pr(C_i = k|\mathbf{x}, \theta^{(t)})(x_i - \mu_k^{(t+1)})^2$$
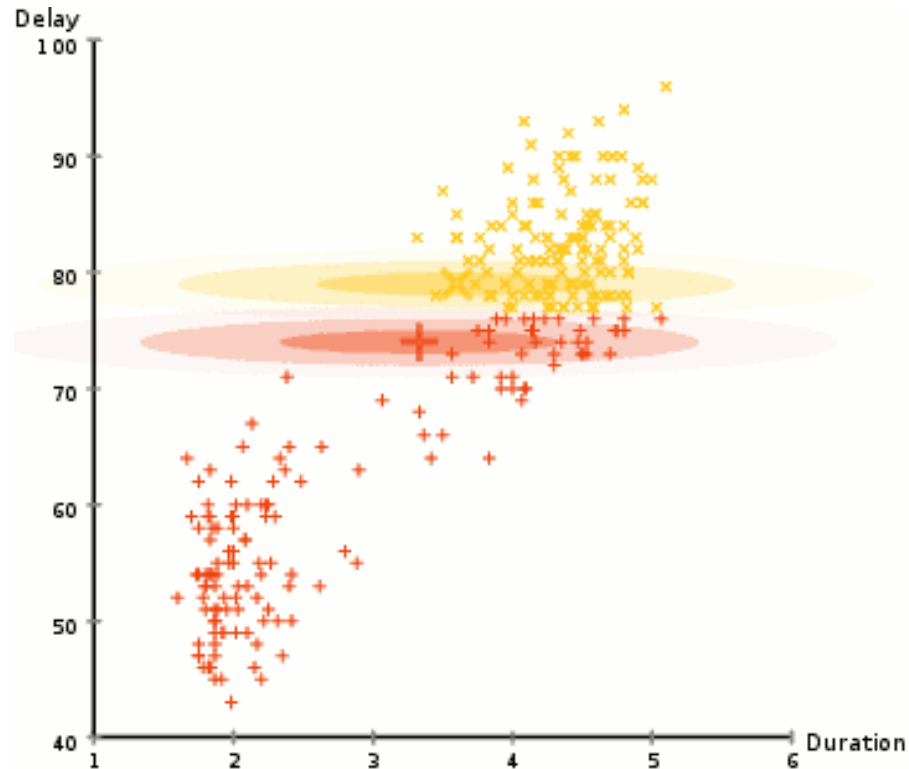
# Examples galaxy

- A numeric vector of velocities in km/sec of 82 galaxies from 6 well-separated conic sections of an unfilled survey of the Corona Borealis region. Multimodality in such surveys is evidence for voids and superclusters in the far universe.



Histogram of galaxies

- `galaxies_EM.R`



Histogram of galaxies

# EM clustering of Old Faithful eruption data.



By Chire - Own work, CC BY-SA 3.0,
https://commons.wikimedia.org/w/index.php?curid=20494862