



UiO : **Matematisk institutt**

Det matematisk-naturvitenskapelige fakultet

STK-4051/9051 Computational Statistics Spring 2024
Ex5

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

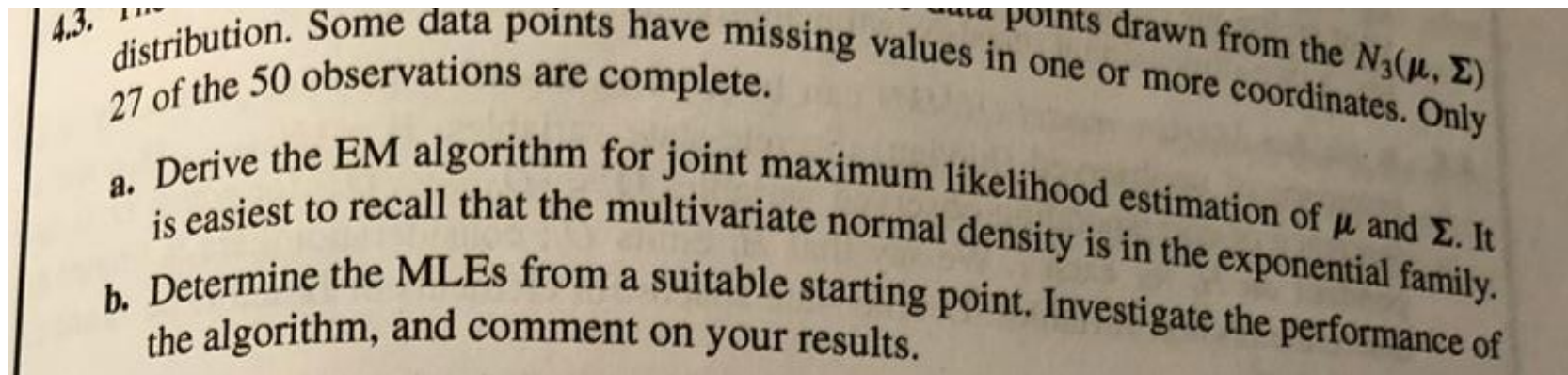


3.8. Thirteen chemical measurements were carried out on each of 178 wines from three regions of Italy [53]. These data are available from the website for this book. Using one or more heuristic search methods from this chapter, partition the wines into three groups for which the total of the within-group sum of squares is minimal. Comment on your work and the results. This is a search problem of size 3^p where $p = 178$. If you have access to standard cluster analysis routines, check your results using a standard method like that of Hartigan and Wong [317].

```
#Calculating objective function
rss = function(x,lab)
{
  rss = 0
  for(j in 1:3)
  {
    z = x[lab==j,]
    mu = colMeans(z)
    res = t(z)-mu
    rss = rss + sum(colSums(res^2))
  }
  rss
}
```

```
#Steepest descent
more = TRUE;r=rep(NA,3)
while(more)
{
  more = FALSE
  for(i in 1:n)
  {
    lab.cur = lab[i]
    for(j in 1:3)
    {
      lab[i] = j
      r[j] = rss(wine[, -1], lab)
    }
    jopt = which.min(r)
    lab[i] = jopt
    if(jopt!=lab.cur)
      more = TRUE
  }
  show(min(r))
}
show(table(wine[, 1], lab))
```

```
#Simulated annealing
lab = sample(1:3,n,replace=T)
rss.cur = rss(wine[, -1], lab)
for(it in c(1:100))
{
  temp = 10/(1+it)
  for(i in 1:n)
  {
    lab.cur = lab[i]
    lab[i] = sample(1:3,1)
    rss.prop = rss(wine[, -1], lab)
    if(runif(1)<exp(-(rss.prop-rss.cur)/temp))
      rss.cur = rss.prop
    else
      lab[i] = lab.cur
  }
  show(rss(wine[, -1], lab))
}
show(table(wine[, 1], lab))
```



Solution to exercise (4.3)

(a). Define \mathbf{y}_i to be the complete data and \mathbf{x}_i the observed data. Note in general that we have

$$\begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} \sim N \left(\begin{pmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{pmatrix}, \begin{pmatrix} \boldsymbol{\Sigma}_{xx} & \boldsymbol{\Sigma}_{xy} \\ \boldsymbol{\Sigma}_{yx} & \boldsymbol{\Sigma}_{yy} \end{pmatrix} \right)$$

then

$$E[\mathbf{y}|\mathbf{x}] = \boldsymbol{\mu}_y + \boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_{xx}^{-1} (\mathbf{x} - \boldsymbol{\mu}_x)$$

$$V[\mathbf{y}|\mathbf{x}] = \boldsymbol{\Sigma}_{yy} - \boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_{xx}^{-1} \boldsymbol{\Sigma}_{xy}$$

$$\begin{aligned} \text{tr}\{\mathbf{A}\mathbf{B}\mathbf{C}\} &= \text{tr}\{\mathbf{B}\mathbf{C}\mathbf{A}\} = \text{tr}\{\mathbf{C}\mathbf{A}\mathbf{B}\} \\ \frac{\partial \log(|\boldsymbol{\Sigma}|)}{\partial \boldsymbol{\Sigma}} &= \boldsymbol{\Sigma}^{-1} \\ \frac{\partial \text{tr}\{\mathbf{A}\boldsymbol{\Sigma}^{-1}\mathbf{B}\}}{\partial \boldsymbol{\Sigma}} &= -(\boldsymbol{\Sigma}^{-1}\mathbf{B}\mathbf{A}\boldsymbol{\Sigma}^{-1})^T \end{aligned}$$

Normal distribution is in the exponential family

- The Exponential family:

$$f_{\mathbf{y}}(\mathbf{y}|\boldsymbol{\theta}) = c_1(\mathbf{y})c_2(\boldsymbol{\theta}) \exp\{\boldsymbol{\theta}^T \mathbf{s}(\mathbf{y})\}$$

$$\frac{1}{\sqrt{2\pi\sigma^2}} \exp -\frac{1}{2} \left(\frac{y - \mu}{\sigma} \right)^2$$

$$\frac{1}{(2\pi)^{\frac{p}{2}} |\Sigma|^{1/2}} \exp -\frac{1}{2} (\mathbf{y} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{y} - \boldsymbol{\mu})$$

$$\frac{|Q|^{\frac{1}{2}}}{(2\pi)^{\frac{p}{2}}} \exp -\frac{1}{2} (\mathbf{y} - \boldsymbol{\mu})^T Q (\mathbf{y} - \boldsymbol{\mu})$$

$$= \frac{|Q|^{\frac{1}{2}}}{(2\pi)^{\frac{p}{2}}} \exp -\frac{1}{2} \boldsymbol{\mu}^T Q \boldsymbol{\mu} \cdot \exp -\frac{1}{2} (\mathbf{y}^T Q \mathbf{y} - 2\boldsymbol{\mu}^T Q \mathbf{y})$$

Alternative:

- The Exponential family:

$$f_{\mathbf{y}}(\mathbf{y}|\boldsymbol{\theta}) = c_1(\mathbf{y})c_2(\boldsymbol{\theta}) \exp\{\boldsymbol{\theta}^T \mathbf{s}(\mathbf{y})\}$$

$$L^{compl}(\boldsymbol{\theta}) = \prod_{i=1}^n \frac{1}{(2\pi)^{3/2} |\boldsymbol{\Sigma}|^{1/2}} \exp\left(-0.5(\mathbf{y}_i - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{y}_i - \boldsymbol{\mu})\right)$$

- Algorithm

E-step $\mathbf{s}^{(t)} = E[\mathbf{s}(\mathbf{Y})|\mathbf{x}; \boldsymbol{\theta}^{(t)}]$

M-step $\boldsymbol{\theta}^{(t+1)}$ solves $E[\mathbf{s}(\mathbf{Y})|\boldsymbol{\theta}] = \mathbf{s}^{(t)}$

Sufficient statistics for normal distribution:

$$\bar{\mathbf{y}} = \frac{1}{n} \sum_{i=1}^n \mathbf{y}_i \quad \overline{\mathbf{y}^2} = \frac{1}{n} \sum_{i=1}^n \mathbf{y}_i \mathbf{y}_i^T$$

$$E(\bar{\mathbf{y}}|\boldsymbol{\theta}^{(t)}, \mathbf{x}) = \frac{1}{n} \sum_{i=1}^n E(\mathbf{y}_i|\mathbf{x}_i, \boldsymbol{\theta}^{(t)})$$

$$E(\bar{\mathbf{y}}|\boldsymbol{\theta}) = \boldsymbol{\mu}$$

$$E[\mathbf{y}|\mathbf{x}] = \boldsymbol{\mu}_y + \boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_{xx}^{-1}(\mathbf{x} - \boldsymbol{\mu}_x)$$

$$V[\mathbf{y}|\mathbf{x}] = \boldsymbol{\Sigma}_{yy} - \boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_{xx}^{-1} \boldsymbol{\Sigma}_{xy}$$

$$\overline{\mathbf{y}^2} = \frac{1}{n} \sum_{i=1}^n \mathbf{y}_i \mathbf{y}_i^T$$

M-step $E(\overline{\mathbf{y}^2} | \boldsymbol{\theta}) = \boldsymbol{\Sigma} + \boldsymbol{\mu} \boldsymbol{\mu}^T$

E-step
$$E(\overline{\mathbf{y}^2} | \boldsymbol{\theta}^{(t)}, \mathbf{x}) = \frac{1}{n} \sum_{i=1}^n E(\mathbf{y}_i \mathbf{y}_i^T | \mathbf{x}_i, \boldsymbol{\theta}^{(t)})$$

$$= \frac{1}{n} \sum_{i=1}^n V(\mathbf{y}_i | \mathbf{x}_i, \boldsymbol{\theta}^{(t)}) + E(\mathbf{y}_i | \mathbf{x}_i, \boldsymbol{\theta}^{(t)}) E(\mathbf{y}_i | \mathbf{x}_i, \boldsymbol{\theta}^{(t)})^T$$

$$E[\mathbf{y} | \mathbf{x}] = \boldsymbol{\mu}_y + \boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_{xx}^{-1} (\mathbf{x} - \boldsymbol{\mu}_x)$$

$$V[\mathbf{y} | \mathbf{x}] = \boldsymbol{\Sigma}_{yy} - \boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_{xx}^{-1} \boldsymbol{\Sigma}_{xy}$$

From scratch...

$$L^{\text{compl}}(\boldsymbol{\theta}) = \prod_{i=1}^n \frac{1}{(2\pi)^{3/2} |\boldsymbol{\Sigma}|^{1/2}} \exp\left(-0.5(\mathbf{y}_i - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{y}_i - \boldsymbol{\mu})\right)$$

$$l^{\text{compl}}(\boldsymbol{\theta}) = -\frac{3n}{2} \log(2\pi) - \frac{n}{2} \log(|\boldsymbol{\Sigma}|) - 0.5 \sum_{i=1}^n (\mathbf{y}_i - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{y}_i - \boldsymbol{\mu})$$

$$\begin{aligned} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^t) &= -\frac{3n}{2} \log(2\pi) - \frac{n}{2} \log(|\boldsymbol{\Sigma}|) - 0.5 \sum_{i=1}^n E[(\mathbf{y}_i - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{y}_i - \boldsymbol{\mu}) | \mathbf{x}_i, \boldsymbol{\theta}^t] \\ &= -\frac{3n}{2} \log(2\pi) - \frac{n}{2} \log(|\boldsymbol{\Sigma}|) - 0.5 \sum_{i=1}^n E[\text{tr}\{(\mathbf{y}_i - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{y}_i - \boldsymbol{\mu})\} | \mathbf{x}_i, \boldsymbol{\theta}^t] \\ &= -\frac{3n}{2} \log(2\pi) - \frac{n}{2} \log(|\boldsymbol{\Sigma}|) - 0.5 \sum_{i=1}^n E[\text{tr}\{\boldsymbol{\Sigma}^{-1}(\mathbf{y}_i - \boldsymbol{\mu})(\mathbf{y}_i - \boldsymbol{\mu})^T\} | \mathbf{x}_i, \boldsymbol{\theta}^t] \\ &= -\frac{3n}{2} \log(2\pi) - \frac{n}{2} \log(|\boldsymbol{\Sigma}|) - 0.5 \sum_{i=1}^n \text{tr}\{\boldsymbol{\Sigma}^{-1} E[(\mathbf{y}_i - \boldsymbol{\mu})(\mathbf{y}_i - \boldsymbol{\mu})^T | \mathbf{x}_i, \boldsymbol{\theta}^t]\} \\ &= -\frac{3n}{2} \log(2\pi) - \frac{n}{2} \log(|\boldsymbol{\Sigma}|) - 0.5 \text{tr}\{\boldsymbol{\Sigma}^{-1} \sum_{i=1}^n E[(\mathbf{y}_i - \boldsymbol{\mu})(\mathbf{y}_i - \boldsymbol{\mu})^T | \mathbf{x}_i, \boldsymbol{\theta}^t]\} \end{aligned}$$

Further,

$$\begin{aligned}
 & E[(\mathbf{y}_i - \boldsymbol{\mu})(\mathbf{y}_i - \boldsymbol{\mu})^T | \mathbf{x}_i, \boldsymbol{\theta}^t] \\
 &= E[(\mathbf{y}_i - E[\mathbf{y}_i | \mathbf{x}_i] + E[\mathbf{y}_i | \mathbf{x}_i, \boldsymbol{\theta}^t] - \boldsymbol{\mu})(\mathbf{y}_i - E[\mathbf{y}_i | \mathbf{x}_i, \boldsymbol{\theta}^t] + E[\mathbf{y}_i | \mathbf{x}_i] - \boldsymbol{\mu})^T | \mathbf{x}_i, \boldsymbol{\theta}^t] \\
 &= E[(\mathbf{y}_i - E[\mathbf{y}_i | \mathbf{x}_i, \boldsymbol{\theta}^t])(\mathbf{y}_i - E[\mathbf{y}_i | \mathbf{x}_i])^T | \mathbf{x}_i, \boldsymbol{\theta}^t] + \\
 &\quad E[(E[\mathbf{y}_i | \mathbf{x}_i, \boldsymbol{\theta}^t] - \boldsymbol{\mu})(E[\mathbf{y}_i | \mathbf{x}_i, \boldsymbol{\theta}^t] - \boldsymbol{\mu})^T | \mathbf{x}_i, \boldsymbol{\theta}^t] \\
 &= V(\mathbf{y}_i | \mathbf{x}_i, \boldsymbol{\theta}^t) + (E[\mathbf{y}_i | \mathbf{x}_i, \boldsymbol{\theta}^t] - \boldsymbol{\mu})(E[\mathbf{y}_i | \mathbf{x}_i, \boldsymbol{\theta}^t] - \boldsymbol{\mu})^T \\
 &= \mathbf{V}_i^t + (\hat{\mathbf{y}}_i^t - \boldsymbol{\mu})(\hat{\mathbf{y}}_i^t - \boldsymbol{\mu})^T
 \end{aligned}$$

Inserting this, we get

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^t) = -\frac{3n}{2} \log(2\pi) - \frac{n}{2} \log(|\boldsymbol{\Sigma}|) - 0.5 \text{tr}\{\boldsymbol{\Sigma}^{-1} \sum_{i=1}^n \mathbf{V}_i^t + (\hat{\mathbf{y}}_i^t - \boldsymbol{\mu})(\hat{\mathbf{y}}_i^t - \boldsymbol{\mu})^T\}$$

$$\begin{aligned}
 & E[(\mathbf{y}_i - E(\mathbf{y}_i | \mathbf{x}, \boldsymbol{\theta}^t)) \cdot (E(\mathbf{y}_i | \mathbf{x}, \boldsymbol{\theta}^t) - \boldsymbol{\mu}) | \mathbf{x}_i, \boldsymbol{\theta}^t] = \\
 & E[(\mathbf{y}_i - E(\mathbf{y}_i | \mathbf{x}, \boldsymbol{\theta}^t)) | \mathbf{x}_i, \boldsymbol{\theta}^t] \cdot (E(\mathbf{y}_i | \mathbf{x}, \boldsymbol{\theta}^t) - \boldsymbol{\mu}) = \\
 & (E[\mathbf{y}_i | \mathbf{x}_i, \boldsymbol{\theta}^t] - E(\mathbf{y}_i | \mathbf{x}, \boldsymbol{\theta}^t)) \cdot (E(\mathbf{y}_i | \mathbf{x}, \boldsymbol{\theta}^t) - \boldsymbol{\mu}) = 0
 \end{aligned}$$

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^t) = -\frac{3n}{2} \log(2\pi) - \frac{n}{2} \log(|\boldsymbol{\Sigma}|) - 0.5 \text{tr}\left\{\boldsymbol{\Sigma}^{-1} \sum_{i=1}^n \mathbf{V}_i^t + (\hat{\mathbf{y}}_i^t - \boldsymbol{\mu})(\hat{\mathbf{y}}_i^t - \boldsymbol{\mu})^T\right\}$$

Differentiating wrt $\boldsymbol{\mu}$, we get

$$\boldsymbol{\Sigma}^{-1} \sum_{i=1}^n (\hat{\mathbf{y}}_i^t - \boldsymbol{\mu}^{t+1}) = \mathbf{0}$$

$$\sum_{i=1}^n (\hat{\mathbf{y}}_i^t - \boldsymbol{\mu}^{t+1}) = \mathbf{0}$$

$$\boldsymbol{\mu}^{t+1} = \frac{1}{n} \sum_{i=1}^n \hat{\mathbf{y}}_i^t$$

$$\begin{aligned} \text{tr}\{\mathbf{A}\mathbf{B}\mathbf{C}\} &= \text{tr}\{\mathbf{B}\mathbf{C}\mathbf{A}\} = \text{tr}\{\mathbf{C}\mathbf{A}\mathbf{B}\} \\ \frac{\partial \log(|\boldsymbol{\Sigma}|)}{\partial \boldsymbol{\Sigma}} &= \boldsymbol{\Sigma}^{-1} \\ \frac{\partial \text{tr}\{\mathbf{A}\boldsymbol{\Sigma}^{-1}\mathbf{B}\}}{\partial \boldsymbol{\Sigma}} &= -(\boldsymbol{\Sigma}^{-1}\mathbf{B}\mathbf{A}\boldsymbol{\Sigma}^{-1})^T \end{aligned}$$

Can similarly obtain

$$\boldsymbol{\Sigma}^{t+1} = \frac{1}{n} \sum_{i=1}^n [\mathbf{V}_i^t + (\hat{\mathbf{y}}_i^t - \boldsymbol{\mu}^{t+1})(\hat{\mathbf{y}}_i^t - \boldsymbol{\mu}^{t+1})^T]$$

```

7 EM.trinormal = function(x)
8 {
9   n = nrow(x)
10  #Make table of missing values
11  num.mis = apply(is.na(x),1,sum)
12  #Initial values based on complete data
13  x.compl = x[num.mis==0,]
14  mu = colMeans(x.compl)
15  sigma = var(x.compl)
16
17  #Make a matrix with estimates of all observations (equal to observed when given)
18  y.hat = x
19  #Array for conditional covariance matrix for all observations
20  v = array(0,c(n,3,3))
21  more = TRUE; eps=0.001
22  while(more)
23  {
24    mu.new = rep(0,3)
25    sigma.new = matrix(0,nrow=3,ncol=3)
26    for(i in 1:n)
27    {
28      #First find y.hat and v
29      if(num.mis[i]>0)
30      {
31        ind = c(1:3)[!is.na(x[i,])]
32        y.hat[i,] = mu[1:3]+sigma[1:3,ind]%%solve(sigma[ind,ind])%%matrix(x[i,ind]-mu[ind],ncol=1)
33        v[i,,] = sigma-sigma[1:3,ind]%%solve(sigma[ind,ind])%%sigma[ind,1:3]
34      }
35    }
36    #Then estimate
37    mu.new = colMeans(y.hat)
38    sigma.new = matrix(0,3,3)
39    for(i in 1:n)
40      sigma.new = sigma.new + v[i,,]+(y.hat[i,]-mu)%o%(y.hat[i,]-mu)
41    sigma.new = sigma.new/n
42    more = (max(abs(mu-mu.new))>eps) | (max(abs(sigma-sigma.new))>eps)
43    mu = mu.new
44    sigma = sigma.new
45  }
46  list(mu=mu,sigma=sigma)
47 }

```

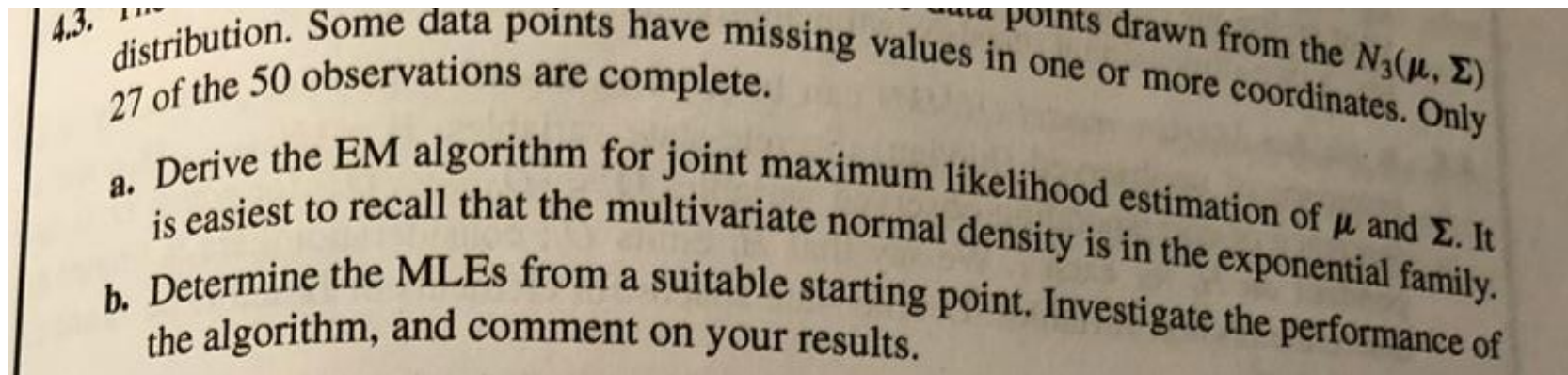
Initialize mu and sigma

Set data for all

Only update data which have missing

$$E[\mathbf{y}|\mathbf{x}] = \boldsymbol{\mu}_y + \boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_{xx}^{-1} (\mathbf{x} - \boldsymbol{\mu}_x)$$

$$V[\mathbf{y}|\mathbf{x}] = \boldsymbol{\Sigma}_{yy} - \boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_{xx}^{-1} \boldsymbol{\Sigma}_{xy}$$



$$\mu^{t+1} = \frac{1}{n} \sum_{i=1}^n \hat{y}_i^t$$

$$E[\mathbf{y}|\mathbf{x}] = \mu_y + \Sigma_{yx} \Sigma_{xx}^{-1} (\mathbf{x} - \mu_x)$$

$$V[\mathbf{y}|\mathbf{x}] = \Sigma_{yy} - \Sigma_{yx} \Sigma_{xx}^{-1} \Sigma_{xy}$$

Can similarly obtain

$$\Sigma^{t+1} = \sum_{i=1}^n [V_i^t + (\hat{y}_i^t - \mu^{t+1})(\hat{y}_i^t - \mu^{t+1})^T]$$

(a). Using the R-script for Exercise 4.3, write a function that has as input the data and some initial estimate of the parameters and returns the ML estimates.

- (b). Write a script that generates $B = 1000$ bootstrap samples by sampling with replacement from the $n = 48$ observations. Run the script and obtain
- (d). Consider now a parametric bootstrap approach, where first you generate $\mathbf{y}_i \sim N(\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}})$ for $i = 1, \dots, n$ and then put those observations to missing that are missing in the original data set. Run this procedure and obtain

Bootstrap

```
#Non-parametric Bootstrapping
B = 1000
n = nrow(x)
mu.boot = matrix(nrow=B,ncol=3)
for(b in 1:B)
{
  ind = sample(1:n,n,replace=T)
  x.new = x[ind,]
  res.new = EM.trinormal(x.new)
  mu.boot[b,] = res.new$mu
}
```

```
#Parametric Bootstrapping
B = 1000
n = nrow(x)
mu.boot = matrix(nrow=B,ncol=3)
for(b in 1:B)
{
  x.new = rmvnorm(n,mu.hat,sigma.hat)
  #Insert missing values
  x.new[is.na(x)] <- NA
  res.new = EM.trinormal(x.new)
  mu.boot[b,] = res.new$mu
}
```

library(mvtnorm)

```
mu[1] mu[2] mu[3]
Estimate 0.879 2.850 9.026
Bias -0.001 0.000 0.006
Variance 0.033 0.018 0.059
SE 0.182 0.135 0.243
2.5% 0.520 2.578 8.571
97.5% 1.237 3.105 9.504
```

```
bias.mu = colMeans(mu.boot)-mu.hat
var.mu = apply(mu.boot,2,var)
ci.mu = apply(mu.boot,2,quantile,probs=c(0.025,0.975))
tab = rbind(mu.hat,bias.mu,var.mu,sqrt(var.mu),ci.mu)
colnames(tab) = c("mu[1]","mu[2]","mu[3]")
rownames(tab)[1:4] = c("Estimate","Bias","Variance","SE")
show("Non-parametric bootstrapping")
show(round(tab,3))
```

```
mu[1] mu[2] mu[3]
Estimate 0.879 2.850 9.026
Bias -0.003 -0.004 0.004
Variance 0.031 0.017 0.056
SE 0.176 0.130 0.237
2.5% 0.546 2.610 8.562
97.5% 1.251 3.130 9.487
```

(c). For the bootstrap procedure above, do it take into account

- wrong model specifications?
- the actual amount of missing data?

(e). For the parametric bootstrap procedure, do it take into account

- wrong model specifications?
- the actual amount of missing data?

Comment on possible differences from the non-parametric version.

Exercise 11 (Stochastic gradient and neural network) (a). Take the *Stoc_grad_NN.R* script from the course web page. Modify the script to $n = 10\,000$. Run the script and inspect how the algorithm perform. Repeat the algorithm several times and look at the variability of the results.

Hint: You probably have to modify the script a bit, perhaps divide it into two. For the first part, the simulation of data, you want to keep that fixed (either by storing the data or by be sure that you use the same seed every time. However, for the run of the SG algorithm you should vary the seed!

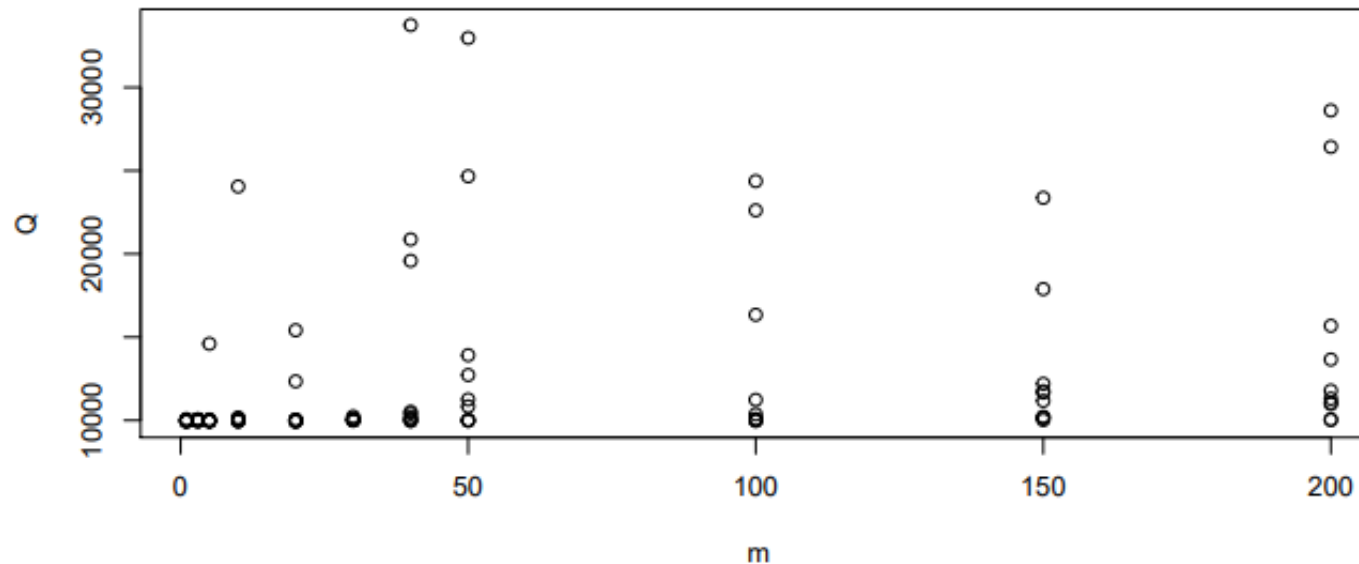
(a). 10 repetitions with random initial points on the α 's and linear regression estimates on the β 's given the α 's of the algorithm with $n = 10\,000$ gave Q - values ranging from 10025.54 to 10368.43:

```
10217.33  10144.09  10105.26  10368.43  10097.87  10188.39  
10120.00  10025.54  10079.28  10200.39
```

This indicates that the algorithms is quite sensitive to starting values.

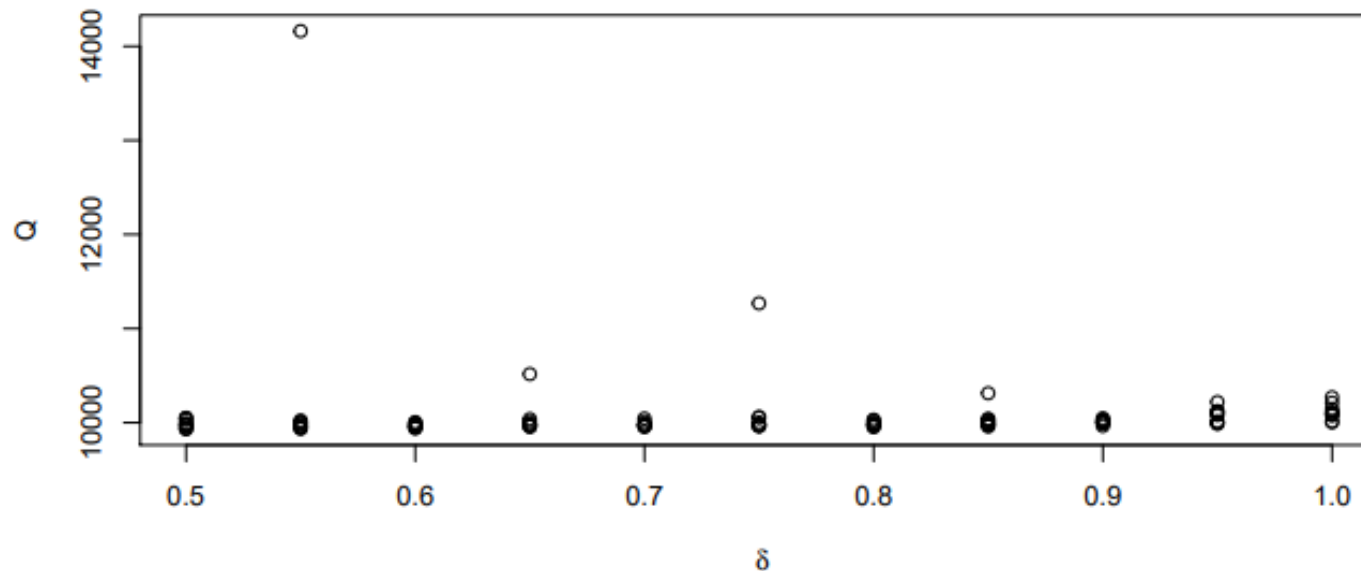
(b). Try to modify the size of the subsample and look at the performance.

(b). Using $m = 1, 3, 5, 10, 20, 30, 40, 50, 100, 150, 200$, repeating the procedure 10 times for each value of m gave the following results:



The minimum value 9925.46 was obtained for $m = 20$ while for $m = 1$ we obtained 9925.61. It seems like the procedure is working best for small m .

- (c). Try to modify the learning rate $\gamma_t = c/(1+t)^\delta$ with $\delta \in [0.5, 1]$. How do the algorithm perform for different choices?
- (c). Using $\delta = 0.5, 0.55, 0.6, 0.65, 0.7, 0.75, 0.8, 0.85, 0.9, 0.95, 1.0$ and $m = 1$, repeating the procedure 10 times for each value of δ gave the following results:



The smallest value obtained was 9931.51 for $\delta = 0.5$. It is not very clear what is preferable here, but top large values should be avoided.

(d). Try to change the initial values of the β 's to be completely random as well. How do the algorithm perform then?

(d). Running now with $m = 1$ and $\delta = 0.5$ 22 times gave the following values of Q :

10041.637	9955.682	10016.338	9950.949	9948.319	9960.379
9941.765	9938.026	9956.522	9986.015	9938.413	9972.349
10025.531	10078.068	9926.557	10004.114	9947.645	9941.453
9964.286	9967.233				

with the smallest value equal to 9926.557. Note that this is smaller than the values obtained by using LS estimates on β !

- Exercise 12 (SG and Geostatistics) (a). Consider the *Geostat_SG.R* script from the course web page. Try out different values of m and evaluate the performance of the algorithm.
- (b). Include a call to the *geoR* package (last commands in the script) which should give the optimal solution. Does the results from the SG algorithm look reasonable?
- (c). Try to increase n by doubling it and run the routine from the *geoR* package. How many times are you able to double it before the routines take too much time to produce results?
- (d). Now try to increase n to 1000 and 10 000. Try out different values of m and learning rates to obtain as good performance as possible.

SGD for dependent data

- Consider spatial data:

$$\mathbf{Y} = \begin{pmatrix} Y(\mathbf{s}_1) \\ Y(\mathbf{s}_2) \\ \vdots \\ Y(\mathbf{s}_n) \end{pmatrix} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

where $\boldsymbol{\mu} = \boldsymbol{\mu}\mathbf{I}$ and $\boldsymbol{\Sigma} = \sigma^2\mathbf{R} + \tau^2\mathbf{I}$, that is

$$\text{cov}[Y(\mathbf{s}_i), Y(\mathbf{s}_j)] = \begin{cases} \sigma^2 r(\|\mathbf{s}_i - \mathbf{s}_j\|; \boldsymbol{\phi}) & \mathbf{s}_i \neq \mathbf{s}_j \\ \sigma^2 + \tau^2 & \mathbf{s}_i = \mathbf{s}_j \end{cases}$$

- Realisation of a process defined continuously in a space \mathcal{S}
- Log-likelihood with $\boldsymbol{\theta} = (\boldsymbol{\mu}, \sigma^2, \tau^2, \boldsymbol{\phi})$

$$l(\boldsymbol{\theta}) = -\frac{n}{2} \log 2\pi - \frac{1}{2} \log(|\boldsymbol{\Sigma}|) - \frac{1}{2} (\mathbf{y} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \boldsymbol{\mu})$$

- In general, computational burden is $O(n^3)$, problematic for large n

KL and Geostatistics

- Liang et al. (2013): Approximate KL by

$$\widehat{KL}(f_{\theta}, g) = C - \frac{1}{\binom{n}{m}} \sum_{k=1}^{\binom{n}{m}} \log(f_{\theta}(\mathbf{y}_k | \mathbf{s}_k))$$

where $(\mathbf{y}_k, \mathbf{s}_k)$ is a subset of (\mathbf{y}, \mathbf{s}) of size m .

- Find θ as the solution of

$$\frac{\partial}{\partial \theta} \widehat{KL}(f_{\theta}, g) = - \frac{1}{\binom{n}{m}} \sum_{k=1}^{\binom{n}{m}} H(\theta, \mathbf{y}_k, \mathbf{s}_k)$$

$$H(\theta, \mathbf{y}_k, \mathbf{s}_k) = \frac{\partial}{\partial \theta} \log(f_{\theta}(\mathbf{y}_k | \mathbf{s}_k))$$

by the **stochastic gradient algorithm!**

Example

$$\log f(\mathbf{y}_k | \mathbf{s}_k) = -\frac{m}{2} \log 2\pi - \frac{1}{2} \log |\Sigma_k| - \frac{1}{2} (\mathbf{y}_k - \mu \mathbf{1}_m)^T \Sigma_k^{-1} (\mathbf{y}_k - \mu \mathbf{1}_m)$$

$$(\Sigma_k)_{i,j} = \text{cov}(Y(\mathbf{s}_{k,i}) - Y(\mathbf{s}_{k,j})) = \tau^2 I(j = i) + \sigma^2 \exp(-(\|\mathbf{s}_{k,i} - \mathbf{s}_{k,j}\|/\phi))$$

$$(\mathbf{R}_k)_{i,j} = \exp(-(\|\mathbf{s}_{k,i} - \mathbf{s}_{k,j}\|/\phi))$$

$$H_\mu(\boldsymbol{\theta}, \mathbf{y}_k, \mathbf{s}_k) = \mathbf{1}_m^T \Sigma_k^{-1} (\mathbf{y}_k - \mu \mathbf{1}_m)$$

$$H_{\sigma^2}(\boldsymbol{\theta}, \mathbf{y}_k, \mathbf{s}_k) = -\frac{1}{2} \text{tr}(\Sigma_k^{-1} \mathbf{R}_k) + \frac{1}{2} (\mathbf{y}_k - \mu \mathbf{1}_m)^T \Sigma_k^{-1} \mathbf{R}_k \Sigma_k^{-1} (\mathbf{y}_k - \mu \mathbf{1}_m)$$

$$H_{\tau^2}(\boldsymbol{\theta}, \mathbf{y}_k, \mathbf{s}_k) = -\frac{1}{2} \text{tr}(\Sigma_k^{-1}) + \frac{1}{2} (\mathbf{y}_k - \mu \mathbf{1}_m)^T \Sigma_k^{-2} (\mathbf{y}_k - \mu \mathbf{1}_m)$$

$$H_\phi(\boldsymbol{\theta}, \mathbf{y}_k, \mathbf{s}_k) = -\frac{\sigma^2}{2} \text{tr}(\Sigma_k^{-1} \frac{d\mathbf{R}_k}{d\phi}) + \frac{\sigma^2}{2} (\mathbf{y}_k - \mu \mathbf{1}_m)^T \Sigma_k^{-1} \frac{d\mathbf{R}_k}{d\phi} \Sigma_k^{-1} (\mathbf{y}_k - \mu \mathbf{1}_m)$$

$$\frac{d(\mathbf{R}_k)_{i,j}}{d\phi} = \|\mathbf{s}_{k,i} - \mathbf{s}_{k,j}\|/\phi^2 \cdot \exp(-(\|\mathbf{s}_{k,i} - \mathbf{s}_{k,j}\|/\phi))$$

```
#Geostatistical model
rm(list=ls())
library(mvtnorm)
set.seed(345)
mu=1
sigma = 1
tau = 1
phi=0.25
n = 500*2
s = cbind(runif(n),runif(n)) # sampling spatial locations
d = as.matrix(dist(s,diag=TRUE,upper=TRUE))
sigma = sigma^2*exp(-d/phi)+tau^2*diag(1,n)

sigma.chol = chol(sigma)
y = mu + rnorm(n)%*%sigma.chol # sampling y"values"

logl = function(y,d,mu,tau2,sig2,phi)
{
  n = length(y)
  sigma = sig2*exp(-d/phi)+tau2*diag(1,n)
  # res = matrix(y-mu,ncol=1)
  # show(c(mu,sig2,tau2,phi))
  # show(class(sigma))
  # l = -0.5*n*log(2*pi)-0.5*determinant(sigma,log=TRUE)$mod-0.5*sum(res*solve(sigma,res))
  d = dmvnorm(y,rep(mu,n),sigma,log=TRUE)
  # show(c(d,l))
  d
}
tr = function(M)
{
  sum(diag(M))
}
```

```

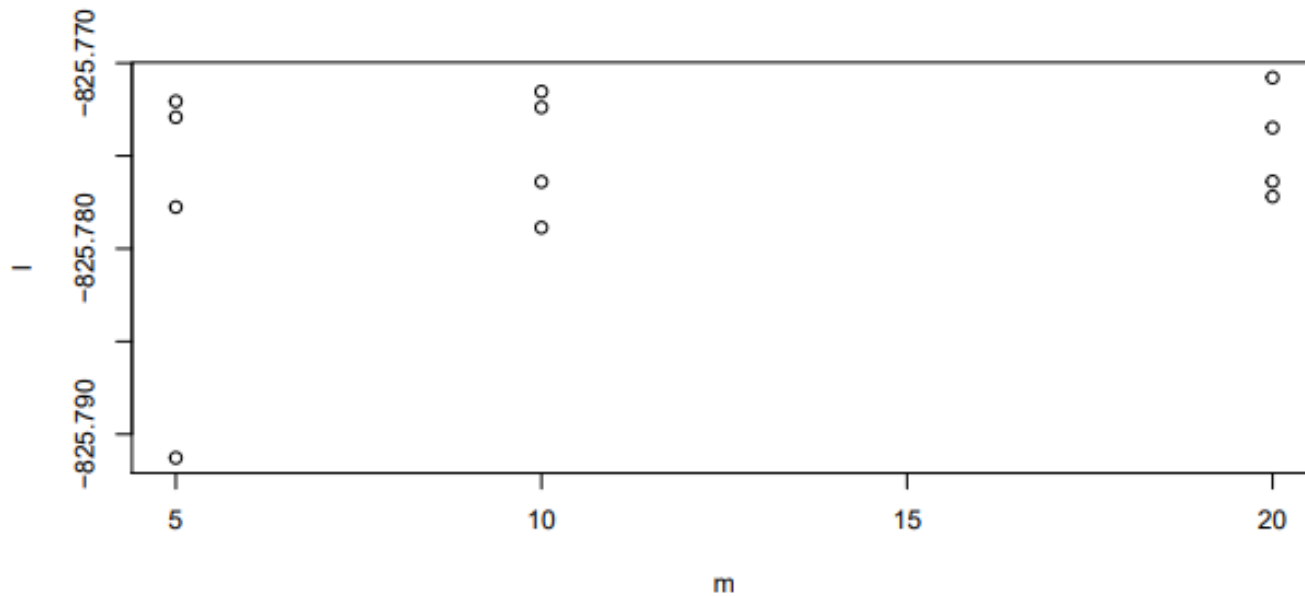
for(it in 1:100)
{
  # sample & set up
  ind = sample(1:n,m,replace=FALSE)
  y.k = y[ind]
  s.k = s[ind,]
  d.k = as.matrix(dist(s.k,diag=TRUE,upper=TRUE))
  R.k = exp(-d.k/phi.hat)
  sigma.k = sig2.hat*R.k+tau2.hat*diag(1,m)
  dR.k = d.k*R.k/phi.hat^2
  res = y.k-mu.hat
  res2 = solve(sigma.k,res)
  res3 = R.k%%res2
  res4 = dR.k%%res2
  H.mu = sum(res2)
  H.sig2 = -0.5*tr(solve(sigma.k,R.k))+0.5*sum(res2*res3)
  H.tau2 = -0.5*tr(solve(sigma.k))+0.5*sum(res2*res2)
  H.phi = -0.5*tr(solve(sigma.k,sig2.hat*dR.k))+0.5*sig2.hat*sum(res2*res4)

  ## stepsize
  alpha = 10/(m*(1+it)^delta)

  ## update
  mu.hat = mu.hat + alpha*H.mu
  a = log(sig2.hat) + alpha*H.sig2/sig2.hat
  sig2.hat = exp(a)
  a = log(tau2.hat) + alpha*H.tau2/tau2.hat
  tau2.hat = exp(a)
  a = log(phi.hat) + alpha*H.phi/phi.hat
  phi.hat = exp(a)
  if(it%%100==0)
  {
    #show(c(H.mu,H.sig2,H.tau2,H.phi))
    l = logl(y,d,mu.hat,tau2.hat,sig2.hat,phi.hat)
    lvec = c(lvec,l)
    itvec = c(itvec,it)
    show(c(it,l,mu.hat,tau2.hat,sig2.hat,phi.hat))
  }
}

```

(a). $m = 5, 10, 20$), 10 000 iterations, 4 repetitions gave



(b). *geor* gave log-likelihood of -791.7.