



UiO • Matematisk institutt

Det matematisk-naturvitenskapelige fakultet

STK-4051/9051 Computational Statistics Spring 2021
Chaper 4

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



Last time

- Examples IRLS, combinatorial optimization
- EM algorithm $Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)}) = E[\log f_Y(\mathbf{y}|\boldsymbol{\theta}) | \mathbf{x}, \boldsymbol{\theta}^{(t)}]$
 - Missing data (Moths)
 - Hidden structure (mixture Gaussian - Galaxy)
 - Proof of increasing log likelihood

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

Today

- EM in Exponential family
- Bootstrap
- Variance estimate in EM
- EM for hidden Markov model

- Stochastic gradient decent

EM recap

- Notation:
 - $Y = (X, Z)$ are complete data
 - X observed,
 - Z missing
 - Have $f_Y(y|\theta)$
 - Want $\max_{\theta} f_X(x|\theta)$

$$f_X(x|\theta) = \int_z f_Y(x, z|\theta) dz$$

Marginal likelihood

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

Complete likelihood

We maximize:

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

Expected value of the complete log likelihood given the observed data using the current estimate of the parameter

EM in exponential family

- The Exponential family:

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

- Includes
 - binomial, multinomial, Poisson, Gaussian, Gamma,...
- $\mathbf{s}(\mathbf{y})$ is a **sufficient** statistic:

$$\begin{aligned} f_s(\mathbf{s}|\boldsymbol{\theta}) &= \int_{\mathbf{y}:\mathbf{s}(\mathbf{y})=\mathbf{s}} f_y(\mathbf{y}|\boldsymbol{\theta}) d\mathbf{y} \\ &= \int_{\mathbf{y}:\mathbf{s}(\mathbf{y})=\mathbf{s}} c_1(\mathbf{y})c_2(\boldsymbol{\theta}) \exp\{\boldsymbol{\theta}^T \mathbf{s}(\mathbf{y})\} d\mathbf{y} \\ &= c_2(\boldsymbol{\theta}) \exp\{\boldsymbol{\theta}^T \mathbf{s}\} \int_{\mathbf{y}:\mathbf{s}(\mathbf{y})=\mathbf{s}} c_1(\mathbf{y}) d\mathbf{y} \\ &= c_2(\boldsymbol{\theta}) \exp\{\boldsymbol{\theta}^T \mathbf{s}\} g(\mathbf{s}) \end{aligned}$$

$$f(\mathbf{y}|\mathbf{s}; \boldsymbol{\theta}) = \frac{f_y(\mathbf{y}|\boldsymbol{\theta})}{f_s(\mathbf{s}|\boldsymbol{\theta})} = \frac{c_1(\mathbf{y})c_2(\boldsymbol{\theta}) \exp\{\boldsymbol{\theta}^T \mathbf{s}\}}{c_2(\boldsymbol{\theta}) \exp\{\boldsymbol{\theta}^T \mathbf{s}\} g(\mathbf{s})} = \frac{c_1(\mathbf{y})}{g(\mathbf{s})}$$

which do not depend on $\boldsymbol{\theta}$!

Why? We can do computations in advance and just identify terms afterwards

Simplifies a lot of standard problems

The EM algorithms in exponential families E & M

- Log-likelihood

$$l(\theta) = \log c_1(\mathbf{y}) + \log c_2(\theta) + \theta^T \mathbf{s}(\mathbf{y})$$

- E-step:

$$Q(\theta|\theta^{(t)}) = k + \log c_2(\theta) + \int \theta^T \mathbf{s}(\mathbf{y}) f_{z|x}(z|x, \theta^{(t)}) dz$$

- M-step:

What is this?

$$\frac{\partial}{\partial \theta} Q(\theta|\theta^{(t)}) = \frac{c_2'(\theta)}{c_2(\theta)} + \int \mathbf{s}(\mathbf{y}) f_{z|x}(z|x, \theta^{(t)}) dz$$

$$\frac{\partial}{\partial \theta} Q(\theta|\theta^{(t)}) = 0$$

$$\int_{\mathbf{y}} c_1(\mathbf{y}) c_2(\theta) \exp\{\theta^T \mathbf{s}(\mathbf{y})\} d\mathbf{y} = 1 \Leftrightarrow$$

$$\frac{\partial}{\partial \theta} \int_{\mathbf{y}} c_1(\mathbf{y}) c_2(\theta) \exp\{\theta^T \mathbf{s}(\mathbf{y})\} d\mathbf{y} = 0 \Leftrightarrow$$

$$\int_{\mathbf{y}} c_1(\mathbf{y}) [c_2'(\theta) \exp\{\theta^T \mathbf{s}(\mathbf{y})\} + c_2(\theta) \exp\{\theta^T \mathbf{s}(\mathbf{y})\} \mathbf{s}(\mathbf{y})] d\mathbf{y} = 0 \Leftrightarrow$$

$$\frac{c_2'(\theta)}{c_2(\theta)} + E[\mathbf{s}(Y); \theta] = 0$$

$$\int \mathbf{s}(\mathbf{y}) f_{z|x}(z|x, \theta^{(t)}) dz = E[\mathbf{s}(Y); \theta]$$

$$\frac{c_2'(\theta)}{c_2(\theta)} = -E[\mathbf{s}(Y); \theta]$$

1st term: Multiply with $\frac{c_2(\theta)}{c_2(\theta)}$ and put $\frac{c_2'(\theta)}{c_2(\theta)}$ outside integral, What remains inside integrates to 1

Results

- Algorithm

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

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

$E[\mathbf{s}(Y)|\mathbf{x}, \theta]$ is the conditional expectation of the missing data given the observed data.

$E[\mathbf{s}(Y)|\theta]$ is the unconditional expectation of the complete data

- Peppered Moths

- Multinomial distribution part of exponential family with $\theta = (\log p_C, \log p_I, \log p_T)$
- Sufficient statistics:

$$S_1 = 2n_{CC} + n_{CI} + n_{CT}$$

$$S_2 = 2n_{II} + n_{CI} + n_{IT}$$

$$S_3 = 2n_{TT} + n_{CT} + n_{IT}$$

$$E[S_1] = 2nP_C$$

$$E[S_2] = 2nP_I$$

$$E[S_3] = 2nP_T$$

- Gives directly

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

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

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

Compute this integral with your old theta
 A bit sloppy and deceiving to say that this is the E-step the real E-step is: to maximize:

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

In the exponential family it turns out that what you need for computations is the expectation of the sufficient statistics

Variance estimate in EM (4.2.3.4)

- Many approaches
- Bootstrapping (4.2.3.3 for EM)
 - General approach (9.1 & 9.2)
- Approximation using, information matrix
$$\text{var}(\hat{\theta}) \approx J_X(\theta)^{-1}$$
 - $J_X(\theta) = -\ell''(\theta|x)$ (observed information matrix)
 - Louis method (4.2.3.1), Just the part about complete and missing information
 - The SEM algorithm (4.2.3.2)
 - Empirical information (4.2.3.4)
 - Numerical Differentiation (4.2.3.5)

Bootstrapping (9.1-9.2.2)

General for exchangeable observations (x_1, \dots, x_n) ,
e.g. iid from $f(x | \theta)$

$$\hat{F}(x) = \frac{1}{n} \sum_{i=1}^n I(x_i < x)$$

The target parameter is $\theta = T(F)$

We make the **estimate** $\hat{\theta} = T(\hat{F})$ (plug in)

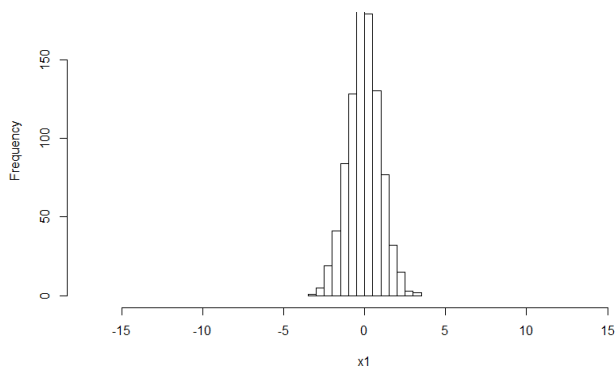
or just $\hat{\theta} = R(x_1, \dots, x_n)$ (some function of data)

- In frequentist inference, the randomness in the estimator comes from the uncertainty in the sampled values. This uncertainty is modelled by the probability density $f(x | \theta)$.
- We could compute the uncertainty by generating many samples from $f(x | \theta)$, and recompute the estimator,
 - but we need many samples from true distribution, we only have one ☹
 - And we do not know the value of θ ☹.
- Two solutions
 - We can get approximate sample from $\hat{F}(x)$ [nonparametric bootstrap]
 - We can sample from the distribution $f(x | \hat{\theta})$ [parametric bootstrap]

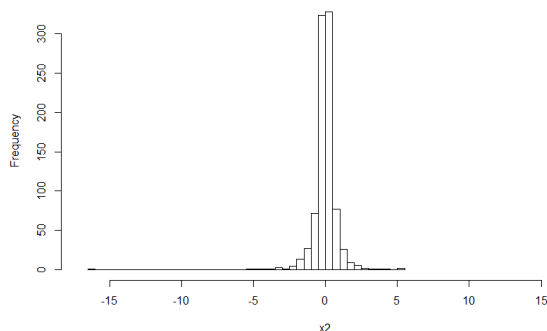
Example: In a symmetric distribution should you estimate the center using the mean or the median?

- If data have a normal distribution the theory says mean.
- But what if the distribution is not known to be normal?

Case 1: normal



Case 2 Student-t

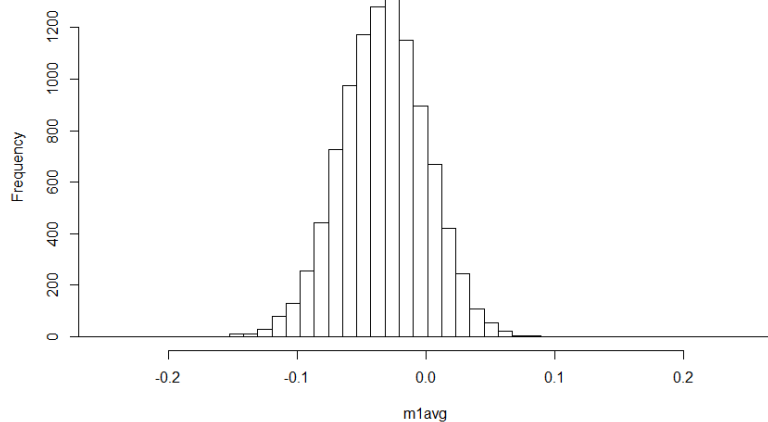


Bootstrap code:

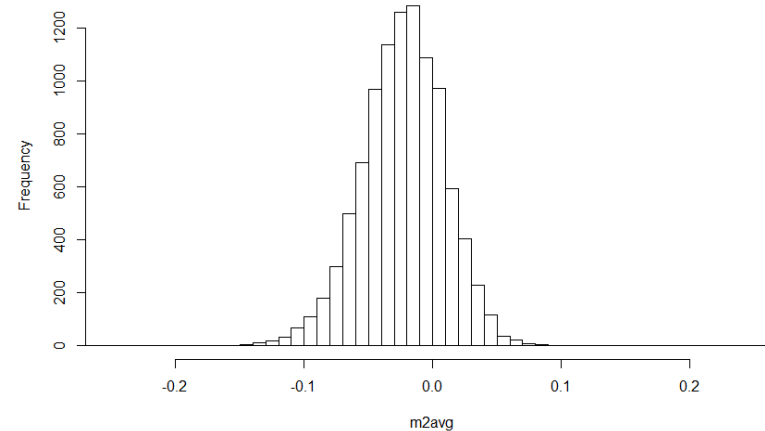
```
m1avg=rep(0,B)
m1med=rep(0,B)

for(b in 1:B)
{
  ind=sample(1:n,n,replace=TRUE)
  m1avg[b]=mean(x1[ind])
  m1med[b]=median(x1[ind])
}
```

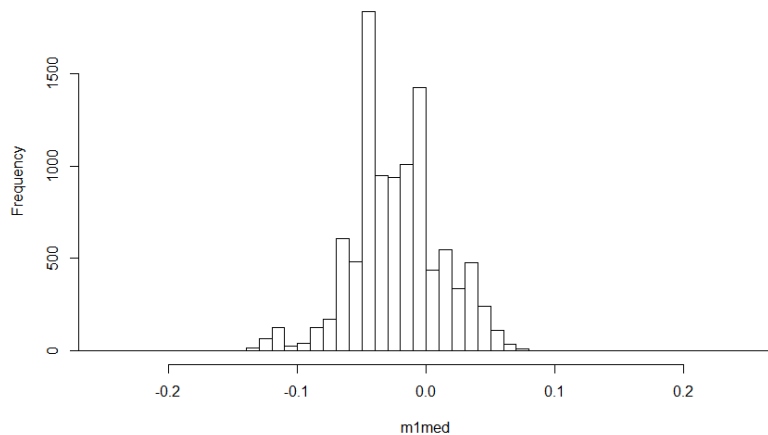
Histogram of m1avg



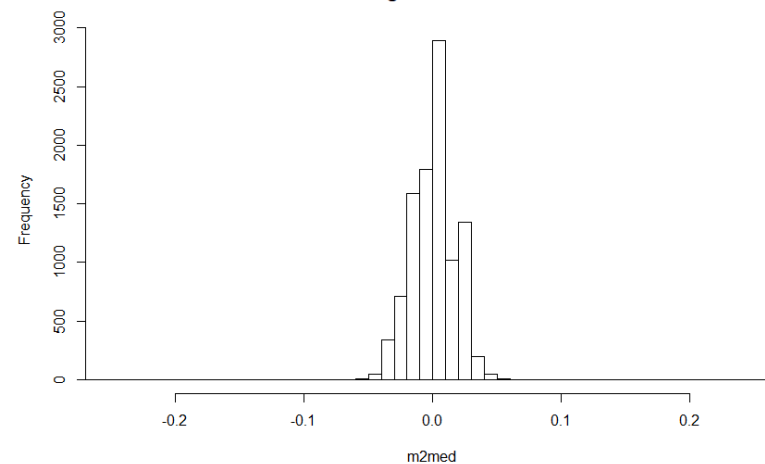
Histogram of m2avg



Histogram of m1med



Histogram of m2med



```
mean(m1avg),mean(m1med),sd(m1avg), sd(m1med))  
-0.03323639 -0.02249270 0.03307427 0.03477931
```

```
:(mean(m2avg),mean(m2med),sd(m2avg), sd(m2med))  
-0.023154073 0.000701744 0.031569047 0.017299653
```

Bootstrapping EM algorithm

- General Bootstrap
- We have $\hat{\theta} = \hat{\theta}(x_1, \dots, x_n)$, i.e. a way to compute an estimate
- Algorithm:
 - For $j=1, \dots, B$
 - Generate sample $\{x_1^*, \dots, x_n^*\}$, from an approximation of $f_X(x|\theta)$ (parametric / nonparametric)
 - Calculate $\hat{\theta}_j = \hat{\theta}(x_1^*, \dots, x_n^*)$
 - $\{\hat{\theta}_j\}_{j=1}^B$ can be seen as a samples from the sampling distribution of $\hat{\theta}$
- Compute variance, quantiles, etc. empirically from $\{\hat{\theta}_j\}_{j=1}^B$
- For the EM algorithm
 - $\hat{\theta}(x_1, \dots, x_n)$ is computed by EM algorithm, $\hat{\theta}_{EM}(x_1, \dots, x_n)$
 - Parametric: Sample $\{y_1^*, \dots, y_n^*\}$ iid $\sim f_Y(y|\theta)$, keep only x , i.e. $\{x_1^*, \dots, x_n^*\}$
 - Nonparametric: Sample $\{x_1^*, \dots, x_n^*\}$ with replacement from $\{x_1, \dots, x_n\}$

Missing information (4.2.3.1)

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

$$-\ell''(\boldsymbol{\theta}|x) = -Q''(\boldsymbol{\theta}|\boldsymbol{\omega})\Big|_{\boldsymbol{\omega}=\boldsymbol{\theta}} + H''(\boldsymbol{\theta}|\boldsymbol{\omega})\Big|_{\boldsymbol{\omega}=\boldsymbol{\theta}}$$

So you do not
differentiate with
respect to
second
argument

$$J_X(\boldsymbol{\theta}) = J_Y(\boldsymbol{\theta}) - J_{Z|X}(\boldsymbol{\theta})$$

Observed information Complete information Missing information

- Nice way of understanding the information loss in missing data
- Sometimes easier to compute $J_Y(\boldsymbol{\theta})$ and $J_{Z|X}(\boldsymbol{\theta})$

Empirical information

- The (expected) Fisher information is defined by

$$\mathbf{I}(\theta) = E\{\ell'(\theta|\mathbf{X})\ell'(\theta|\mathbf{X})\}$$

- Further, since $E[\ell'(\theta|\mathbf{X})] = \mathbf{0}$, we have

$$\mathbf{I}(\theta) = \text{var}[\ell'(\theta|\mathbf{X})]$$

- If we have **IID data**,

$$\ell'(\theta|\mathbf{x}) = \frac{\partial \log f_{\mathbf{X}}(\mathbf{x}|\theta)}{\partial \theta} = \sum_{i=1}^n \frac{\partial \log f_{X_i}(\mathbf{x}_i|\theta)}{\partial \theta} \equiv \sum_{i=1}^n \ell'(\theta|\mathbf{x}_i)$$

- We can **estimate** the information for **one** observation, $\mathbf{I}_1(\theta)$ by

$$\hat{\mathbf{I}}_1(\theta) = \frac{1}{n} \sum_{i=1}^n \ell'(\theta|\mathbf{x}_i)\ell'(\theta|\mathbf{x}_i)^T - \frac{1}{n^2} \ell'(\theta|\mathbf{x})\ell'(\theta|\mathbf{x})^T$$

while information for **all** data can be estimated by

$$\hat{\mathbf{I}}(\theta) = n \cdot \hat{\mathbf{I}}_1(\theta)$$

But how do we compute $\ell'(\theta|\mathbf{x})$?

Computing the score function in EM

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

- giving

$$\begin{aligned} Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)}) - \ell(\boldsymbol{\theta}|\mathbf{x}) &= H(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)}) \\ &\leq H(\boldsymbol{\theta}^{(t)}|\boldsymbol{\theta}^{(t)}) \\ &= Q(\boldsymbol{\theta}^{(t)}|\boldsymbol{\theta}^{(t)}) - \ell(\boldsymbol{\theta}^{(t)}|\mathbf{x}) \end{aligned}$$

so $\boldsymbol{\theta}^{(t)}$ is a max point of $Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)}) - \ell(\boldsymbol{\theta}|\mathbf{x})$.

- Assuming smooth functions,

$$Q'(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)}) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^{(t)}} = \ell'(\boldsymbol{\theta}|\mathbf{x}) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^{(t)}}$$

$$\frac{\partial}{\partial \boldsymbol{\theta}} (Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)}) - \ell(\boldsymbol{\theta}|\mathbf{x})) = \mathbf{0}$$

- $Q'(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)}) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^{(t)}}$ typically calculated in the M-step of the EM-algorithm!

EM-Hidden Markov model

- Assume now $Y_i = (X_i, C_i)$ where i refer to timepoint.
- Model:

$$\Pr(C_1 = k) = \pi_1(k)$$

$$\Pr(C_i = k | C_{i-1} = j) = p(k|j),$$

$$X_i | C_i = k \sim N(\mu_k, \sigma_k) = f_k(x_i)$$

- $\{C_i\}$ is a **Markov chain** and **hidden/missing**
- Complete likelihood:

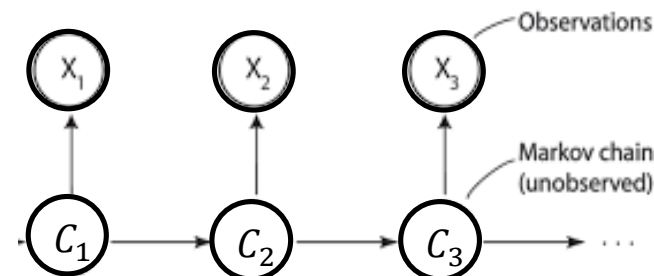
$$L(\theta) = \pi_1(c_1) \phi(x_1; \mu_{c_1}, \sigma_{c_1}) \prod_{i=2}^n p(c_i | c_{i-1}) \phi(x_i; \mu_{c_i}, \sigma_{c_i})$$

$$\ell(\theta) = \log(\pi_1(c_1)) + \log[\phi(x_1; \mu_{c_1}, \sigma_{c_1})] +$$

$$\sum_{i=2}^n [\log[p(c_i | c_{i-1})] + \log[\phi(x_i; \mu_{c_i}, \sigma_{c_i})]]$$

$$= \sum_{k=1}^K I(c_1 = k) [\log(\pi_1(k)) + \log[\phi(x_1; \mu_k, \sigma_k)]] +$$

$$\sum_{i=2}^n \sum_{k=1}^K \sum_{j=1}^K I(c_i = k, c_{i-1} = j) [\log[p(k|j)] + \log[\phi(x_i; \mu_k, \sigma_k)]]$$



EM - Hidden Markov model

- We get

$$Q(\theta|\theta^{(t)}) = \sum_{k=1}^K \Pr(C_1 = k|\mathbf{x}, \theta^{(t)}) [\log(\pi_1(k)) + \log[\phi(x_1; \mu_k, \sigma_k)]] +$$
$$\sum_{i=2}^n \sum_{k=1}^K \sum_{j=1}^K \Pr(C_i = k, C_{i-1} = j|\mathbf{x}, \theta^{(t)}) [\log[p(k|j)] + \log[\phi(x_i; \mu_k, \sigma_k)]]$$

- Main problem now: Calculation of

$$\Pr(C_1 = k|\mathbf{x}, \theta^{(t)})$$

$$\Pr(C_i = k, C_{i-1} = l|\mathbf{x}, \theta^{(t)}), \quad i = 2, \dots, n$$

General idea in hidden Markov Model

Any time \swarrow $P(C_i | \mathbf{x}_{1:n})$ \nwarrow All data

- We compute «forward» (filtering)

$P(C_{i-1} \mathbf{x}_{1:(i-1)})$	$P(C_i \mathbf{x}_{1:(i-1)})$	$P(C_i \mathbf{x}_{1:i})$	$P(C_{i+1} \mathbf{x}_{1:i})$...	$P(C_n \mathbf{x}_{1:n})$
Update	Predict	Update	Predict		Update

- Then we compute backward (smoothing)

$P(C_{i+1} | \mathbf{x}_{1:n})$ $P(C_i | \mathbf{x}_{1:n})$ $P(C_{i-1} | \mathbf{x}_{1:n})$ \cdots $P(C_1 | \mathbf{x}_{1:n})$

- At the end we combine to get

$$P(C_i, C_{i-1} | \mathbf{x}_{1:n})$$

Hidden Markov model

- We have

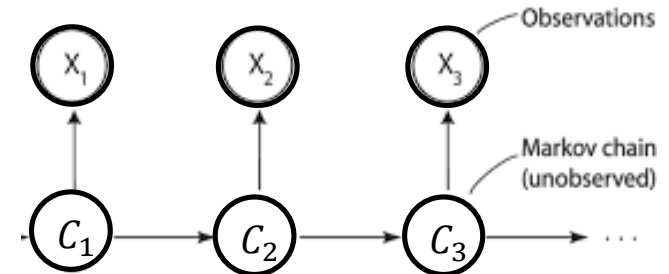
$$L(\theta) = f(\mathbf{x}|\theta)$$

$$= f(x_1|\theta) \prod_{i=2}^n f(x_i|\mathbf{x}_{1:i-1}; \theta)$$

where

$$f(x_1|\theta) = \sum_{k=1}^K \pi_k f_k(x_1; \theta)$$

$$f(x_i|\mathbf{x}_{1:i-1}; \theta) = \sum_{k=1}^K q_{i|i-1}(k) f_k(x_i|\theta)$$



HMM forward equations

- Define $q_{i|j}(k) = \Pr(C_i = k | x_{1:j})$, $x_{1:j} = (x_1, \dots, x_j)$.

- Initialization

$$q_{1|1}(k) = \Pr(C_1 = k | x_1) = \frac{\pi_1(k) f_k(x_1)}{\sum_{j=1}^K \pi_1(j) f_j(x_1)}$$

- Prediction:

$$\begin{aligned} q_{i|i-1}(k) &= \Pr(C_i = k | x_{1:i-1}) \\ &= \sum_{j=1}^K \Pr(C_i = k | C_{i-1} = j, x_{1:i-1}) \Pr(C_{i-1} = j | x_{1:i-1}) \\ &= \sum_{j=1}^K p(k|j) q_{i-1|i-1}(j) \end{aligned}$$

- Updating:

$$\begin{aligned} q_{i|i}(k) &= \Pr(C_i = k | x_{1:i}) = \frac{\Pr(C_i = k | x_{1:i-1}) p(x_i | C_i = k)}{p(x_i | x_{1:i-1})} \\ &\propto q_{i|i-1}(k) f_k(x_i) \end{aligned}$$

Backward equations

- $q_{n|n}(k) = \Pr(C_n = k | x_{1:n})$ obtained from forward equations
- Going backwards:

$$q_{i|n}(k) = \Pr(C_i = k | x_{1:n})$$

Expand

$$= \sum_{\ell=1}^K \Pr(C_i = k, C_{i+1} = \ell | x_{1:n})$$

$$= \sum_{\ell=1}^K \Pr(C_i = k | C_{i+1} = \ell, x_{1:n}) \Pr(C_{i+1} = \ell | x_{1:n})$$

Recognize

$$= \sum_{\ell=1}^K \Pr(C_i = k | C_{i+1} = \ell, x_{1:i}) q_{i+1|n}(\ell)$$

Bayes formula

$$= \sum_{\ell=1}^K \frac{\Pr(C_i = k | x_{1:i}) \Pr(C_{i+1} = \ell | C_i = k, x_{1:i})}{\Pr(C_{i+1} = \ell | x_{1:i})} q_{i+1|n}(\ell)$$

Recognize

$$= \sum_{\ell=1}^K \frac{q_{i|i}(k) p(\ell | k)}{q_{i+1|i}(\ell)} q_{i+1|n}(\ell)$$

Because of the Markov structure

$$= P(C_{i+1} = \ell | C_i = k)$$

Sequence probability

- Needed $\Pr(C_i = k, C_{i-1} = \ell | x_{1:n})$ within EM

$$\begin{aligned} & \Pr(C_i = k, C_{i-1} = \ell | x_{1:n}) \\ &= \Pr(C_i = k | x_{1:n}) \Pr(C_{i-1} = \ell | C_i = k, x_{1:n}) \end{aligned}$$

Recognize

$$= q_{i|n}(k) \Pr(C_{i-1} = \ell | C_i = k, x_{1:i-1})$$

Bayes formula

$$= q_{i|n}(k) \frac{\Pr(C_{i-1} = \ell | x_{1:i-1}) \Pr(C_i = k | C_{i-1} = \ell, x_{1:i-1})}{\Pr(C_i = k | x_{1:i-1})}$$

Recognize

$$= q_{i|n}(k) \frac{q_{i-1|i-1}(\ell) p(k|\ell)}{q_{i|i-1}(k)}$$

HMM - M-step

- Estimation of probabilities

$$\pi_1^{(t+1)}(k) = \Pr(C_1 = k | x_{1:n}, \theta^{(t)})$$

$$p^{(t+1)}(k|j) = \frac{\sum_{i=2}^n \Pr(C_i = k, C_{i-1} = j | x_{1:n}, \theta^{(t)})}{\sum_{i=2}^n \Pr(C_{i-1} = j | x_{1:n}, \theta^{(t)})}$$

- Estimation of parameters in Gaussian distribution

$$\mu_k^{(t+1)} = \frac{\sum_{i=1}^n \Pr(C_i = k | x_{1:n}, \theta^{(t)}) x_i}{\sum_{i=1}^n \Pr(C_i = k | x_{1:n}, \theta^{(t)})}$$

$$(\sigma_k^2)^{(t+1)} = \frac{\sum_{i=1}^n \Pr(C_i = k | x_{1:n}, \theta^{(t)}) (x_i - \mu_k^{(t+1)})^2}{\sum_{i=1}^n \Pr(C_i = k | x_{1:n}, \theta^{(t)})}$$