

Theorem: The Bell numbers satisfy the following recurrence relation:

$$B_{n+1} = \sum_{k=0}^n \binom{n}{k} B_k$$

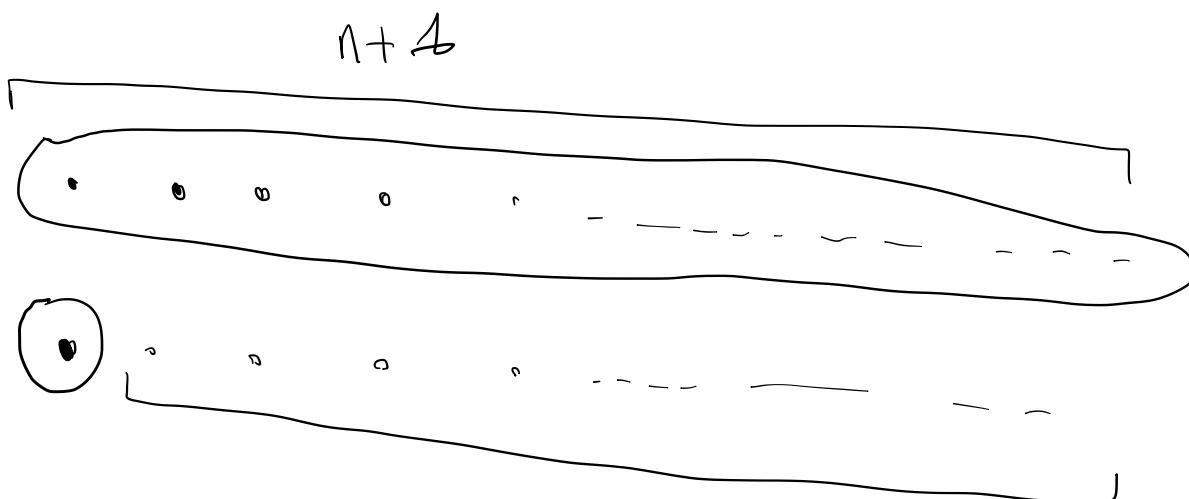
Proof: Proof of the recurrence relation:

First, consider a partition  $\mathcal{P}$  of  $(n+1)$  observations in  $(C_k)_k$  clusters. The cluster containing the first observation is removed. The number of observations left varies from 0 (the cluster removed contained all the observations) to  $n$  (it contained only the first observation).

$n_1$  being always removed, the  $k$  observations left have to be chosen among  $n$  elements.

$\Rightarrow \binom{n}{k}$ . Moreover,  $B_k$  is the number of ways to partition the  $k$  elements.

So:  $B_{n+1} = \sum_{k=0}^n \binom{n}{k} B_k$



Remark:  $S = \frac{1}{n} \sum_{i=1}^n \|x_i - \hat{x}\|^2 = \sum_{j=1}^p \hat{\sigma}_j^2$  (with  $x_i \in \mathbb{R}^p$ )

where  $\hat{\sigma}_j^2$  is the empirical (biased) variance of variable  $j$ .

Proof:

$$S = \frac{1}{n} \sum_{i=1}^n \|x_i - \hat{x}\|^2 = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^p (x_{ij} - \hat{x}_j)^2 = \sum_{j=1}^p \left\{ \underbrace{\frac{1}{n} \sum_{i=1}^n (x_{ij} - \hat{x}_j)^2}_{\hat{\sigma}_j^2} \right\}$$

$$\Rightarrow S = \sum_{j=1}^p \hat{\sigma}_j^2$$

Theorem (Huygens) :

$$S = W + B$$

Proof:

$$\begin{aligned} \sum_{i=1}^n \|x_i - \hat{x}\|^2 &= \sum_{k=1}^K \sum_{x_i \in C_k} \|(x_i - m_k) + (m_k - \hat{x})\|^2 \\ &= \sum_{k=1}^K \sum_{x_i \in C_k} \left\{ \|x_i - m_k\|^2 + 2(x_i - m_k)^T (m_k - \hat{x}) + \|m_k - \hat{x}\|^2 \right\} \\ &= \sum_{k=1}^K \sum_{x_i \in C_k} \|x_i - m_k\|^2 + 2 \sum_{k=1}^K \underbrace{\left( \underbrace{\sum_{x_i \in C_k} x_i}_{n_k m_k} - m_k \underbrace{\sum_{x_i \in C_k} 1}_{n_k} \right)^T (m_k - \hat{x})}_{=0} + \sum_{k=1}^K n_k \|m_k - \hat{x}\|^2 \\ &= n_1 W + n B \end{aligned}$$

Properties: Consider the random sample  $(x_1, \dots, x_n)$  with  $x_i \sim \text{ind}(m, \Sigma)$ .

The MLE of  $\mu$  and  $\Sigma$  are:

$$\rightarrow \hat{\mu} = \bar{x} \quad (\text{empirical mean})$$

$$\rightarrow \hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^\top$$

Proof:

$$\begin{aligned} L(\mu, \Sigma) &= \log p(x_1, \dots, x_n | \mu, \Sigma) \\ &= \sum_{i=1}^n \log p(x_i | \mu, \Sigma) \\ &= -\frac{n\theta}{2} \log(2\pi) - \frac{n}{2} \log|\Sigma| - \frac{1}{2} \sum_{i=1}^n (x_i - \mu)^\top \Sigma^{-1} (x_i - \mu) \end{aligned}$$

$$\begin{aligned} \text{for } \mu: \quad L(\mu, \Sigma) &= -\frac{1}{2} \sum_{i=1}^n \left\{ x_i^\top \Sigma^{-1} x_i - 2x_i^\top \Sigma^{-1} \mu + \mu^\top \Sigma^{-1} \mu \right\} + \text{cte} \quad (\text{cte does not depend on } \mu) \\ &= -\frac{1}{2} \sum_{i=1}^n \left\{ -2x_i^\top \Sigma^{-1} \mu + \mu^\top \Sigma^{-1} \mu \right\} + \text{cte} \end{aligned}$$

Recall: consider  $f: \mathbb{R}^d \rightarrow \mathbb{R}$ ,  $f(x) = y^\top x = \sum_{j=1}^d y_j x_j \Rightarrow \nabla_x f(x) = y$

$$\begin{aligned} \text{Moreover: consider } g: \mathbb{R}^d \rightarrow \mathbb{R}, \quad g(x) &= x^\top A x = \sum_{i=1}^d \sum_{j=1}^d x_i A_{ij} x_j = \sum_{i=1}^d x_i^2 A_{ii} + \sum_{i \neq j} x_i A_{ij} x_j \\ &\Rightarrow \frac{\partial g(x)}{\partial x_k} = 2x_k A_{kk} + \underbrace{\sum_{j \neq k} x_j A_{kj} x_j + \sum_{i \neq k} x_i A_{ik} x_i}_{\sum_{j \neq k} (A_{kj} + A_{jk}) x_j} \\ &= 2 \sum_{j \neq k} A_{kj} x_j \quad \text{if } A \text{ is symmetric} \end{aligned}$$

$$\Rightarrow \frac{\partial g(x)}{\partial x_a} = 2x_a A_{aa} + \sum_{j \neq a} x_j A_{aj} = 2 \sum_{j=1}^d x_j A_{aj} = 2 A_{aa} x_a$$

$$\Rightarrow \nabla_x g(x) = 2Ax$$

Going back to our problem, we have:

$$\nabla_\mu L(\mu, \Sigma) = \sum_{i=1}^n \Sigma^{-1} x_i - n \Sigma^{-1} \mu$$

$$\begin{aligned} \text{We solve: } \nabla_\mu L(\mu, \Sigma) &= 0 \iff \sum_{i=1}^n \Sigma^{-1} x_i = n \Sigma^{-1} \mu \\ \Rightarrow \hat{\mu} &= \frac{\sum_{i=1}^n x_i}{n} \quad (\text{only if } \Sigma \text{ can be inverted}) \end{aligned}$$

$$\text{For } \Sigma: L(\mu, \Sigma) = -\frac{n}{2} \log |\Sigma| - \frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) + \text{const}$$

Recall:  $\mathbf{x}^\top A \mathbf{x} = \sum_{q,p} x_q A_{qp} x_p = \sum_{q,p} A_{qp} B_{pq}$  where  $B_{pq} = x_p x_q \Leftrightarrow B = \mathbf{x} \mathbf{x}^\top$

Moreover:  $(AB)_{ij} = \sum_k A_{ik} B_{kj} \Rightarrow \text{Tr}(AB) = \sum_i (AB)_{ii} = \sum_i \sum_k A_{ik} B_{ki}$

$$\Rightarrow \text{so: } \mathbf{x}^\top A \mathbf{x} = \text{Tr}(A \mathbf{x} \mathbf{x}^\top)$$

$$\text{Therefore: } L(\mu, \Sigma) = -\frac{n}{2} \log |\Sigma| - \frac{1}{2} \sum_{i=1}^n \text{Tr} \left( \Sigma^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) (\mathbf{x}_i - \boldsymbol{\mu})^\top \right) + \text{const}$$

$$\text{Recall that: } \nabla_X \log |X| = (X^{-1})^\top = (X^\top)^{-1}$$

$$\text{Moreover: } \nabla_X \mathbf{a}^\top X \mathbf{b} = \mathbf{a} \mathbf{b}^\top \text{ and } \nabla_X \text{Tr}(A X^{-1} B) = -X^{-\top} A B^\top X^{-\top}$$

$$\text{So, we find: } \nabla_\Sigma L(\mu, \Sigma) = -\frac{n}{2} \Sigma^{-1} + \frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu}) (\mathbf{x}_i - \boldsymbol{\mu})^\top \Sigma^{-1}$$

$$\text{We solve: } \nabla_\Sigma L(\mu, \Sigma) = \mathbf{0}_{d \times d} \Leftrightarrow -\frac{n}{2} \Sigma + \frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu}) (\mathbf{x}_i - \boldsymbol{\mu})^\top = \mathbf{0}_{d \times d}$$

$$\Rightarrow \hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu}) (\mathbf{x}_i - \boldsymbol{\mu})^\top$$

Property: Given the observations (and the parameters), all the  $Z_i$  are independent.

Proof:

$$P((z_i)_{i=1}^n | (x_i)_{i=1}^n, \pi, \theta) = \frac{P((x_i, z_i)_{i=1}^n | \pi, \theta)}{P(x_i)_{i=1}^n | \pi, \theta)} = \frac{\prod_{i=1}^n P(x_i, z_i | \pi, \theta)}{\prod_{i=1}^n P(x_i | \pi, \theta)} = \prod_{i=1}^n \frac{P(x_i, z_i | \pi, \theta)}{P(x_i | \pi, \theta)}$$

$$= \prod_{i=1}^n P(z_i | x_i, \pi, \theta)$$

Bayes rule

Property:  $P(z_i | x_i, \pi, \theta) = \mathcal{M}(z_i | \mu_i, \Sigma_i)$  where  $\mu_i^\top = (\mu_{i1}, \dots, \mu_{ik})$  and  $\Sigma_{ik} = \frac{\pi_k \mathcal{N}(x_i; \mu_k, \Sigma_k)}{\sum_{l=1}^k \pi_l \mathcal{N}(x_i; \mu_l, \Sigma_l)}$

Proof:

$$P(z_i | x_i, \pi, \theta) = \frac{P(x_i | z_i, \theta) P(z_i | \pi)}{P(x_i | \pi, \theta)}$$

$$= \frac{\prod_{k=1}^K \left( \pi_k \mathcal{N}(x_i; \mu_k, \Sigma_k) \right)^{z_{ik}}}{\sum_{l=1}^K \pi_l \mathcal{N}(x_i; \mu_l, \Sigma_l)}$$

$$= \prod_{k=1}^K \left( \frac{\pi_k \mathcal{N}(x_i; \mu_k, \Sigma_k)}{\sum_{l=1}^K \pi_l \mathcal{N}(x_i; \mu_l, \Sigma_l)} \right)^{z_{ik}}$$

Property: If  $Y \sim \mathcal{M}(\mu, \Sigma)$  where  $\mu^\top = (\mu_1, \dots, \mu_K)$

then:  $E[Y_k] = \mu_k$

Proof:  $E[Y] = \sum_y y P(Y=y) = \sum_y y \prod_{k=1}^K \frac{\mu_k^{y_k}}{\binom{m}{y_k}} = \begin{pmatrix} m \\ \vdots \\ m_k \end{pmatrix}$

$\Rightarrow$  in particular:  $E[Y_k] = \mu_k$ .

Property: The EM algorithm does increase the log-likelihood  $L_{(x_i)_i}(\pi, \theta)$  at each iteration.

Proof (general): for any law  $\rho((z_i)_i)$ , the log-likelihood can always be decomposed in two terms:

$$\textcircled{1} \quad L_{(x_i)_i}(\mu, \theta) = \log p((x_i)_i | \mu, \theta)$$

$$= \mathcal{L}(\rho((z_i)_i; \mu, \theta)) + KL(\rho((z_i)_i) \| p((z_i)_i | (x_i)_i; \pi, \theta))$$

$$\text{where: } \mathcal{L}(\rho((z_i)_i; \mu, \theta)) = \sum_{(z_i)_i} \rho((z_i)_i) \log \frac{p((x_i, z_i)_i | \pi, \theta)}{\rho((z_i)_i)}$$

$$\text{and } KL(\rho((z_i)_i) \| p((z_i)_i | (x_i)_i; \pi, \theta)) = - \sum_{(z_i)_i} \rho((z_i)_i) \log \frac{p((z_i)_i | (x_i)_i; \pi, \theta)}{\rho((z_i)_i)}$$

$$\text{Proof of } \textcircled{1}: \mathcal{L}(\rho((z_i)_i; \mu, \theta)) + KL(\rho((z_i)_i) \| p((z_i)_i | (x_i)_i; \pi, \theta))$$

$$= \sum_{(z_i)_i} \rho((z_i)_i) \log \frac{p((x_i, z_i)_i | \pi, \theta)}{p((z_i)_i | (x_i)_i; \pi, \theta)}$$

$$= \underbrace{\sum_{(z_i)_i} \rho((z_i)_i)}_{\mathcal{L}(\rho((z_i)_i; \mu, \theta))} \log p((x_i)_i | \pi, \theta)$$

$$= \log p((x_i)_i | \pi, \theta)$$

$$= L_{(x_i)_i}(\pi, \theta)$$

2) Therefore, since  $L_{(x_i)}(\pi, \theta)$  does not depend on  $R((z_i)_i)$ , with  $\mathcal{L} \uparrow$  then  $KL \downarrow$

In fact,  $\mathcal{L}$  acts as a lower bound of the log-likelihood.

$\Rightarrow$  we focus the optimisation on  $\mathcal{L}$  with respect to  $R((z_i)_i), \pi, \theta$ .

EII: E-step:  $\pi, \theta$  fixed,  $\mathcal{L}$  maximized with respect to  $R((z_i)_i)$ . In the case of Gaussian mixture models,  $R((z_i)_i | (x_i)_i, \pi, \theta)$  has an analytical form, so  $KL$  can reach 0 by setting  $R((z_i)_i) = \sum_i p(z_i | x_i, \pi, \theta)$  [Optimal]

$$= \prod_{i=1}^n \mathcal{M}(z_i; 1, \tau_i)$$

$\Leftrightarrow$  characterized by the  $\tau_i$ .

M-step:  $\tau_i$  fixed,  $\mathcal{L}$  maximized with respect to  $\pi$  and  $\theta$ :

$$\begin{aligned} (\hat{\pi}, \hat{\theta}) &= \underset{\pi, \theta}{\operatorname{argmax}} \mathcal{L}(R((z_i)_i), \pi, \theta) = \underset{\pi, \theta}{\operatorname{argmax}} \sum_{(z_i)_i} R((z_i)_i) \log p(z_i | x_i, z_i, \pi, \theta) \\ &= \underset{\pi, \theta}{\operatorname{argmax}} E_{(z_i)_i} \left[ L_{(x_i, z_i)_i}(\pi, \theta) \right] \end{aligned}$$

③ Recall that after the E-step,  $KL = 0$  so  $L_{(x_i)}(\pi, \theta) = \mathcal{L}(R((z_i)_i), \pi, \theta)$

Thus, Init:  $\tau_i^{(0)}$   
after M-step:  $\pi^*, \theta^* \rightarrow \mathcal{L}((\tau_i^{(0)})_i; \pi^*, \theta^*)$

E-step:  $\tau_i^1 \rightarrow \mathcal{L}((\tau_i^1)_i; \pi^*, \theta^*) = L_{(x_i)}(\pi^*, \theta^*)$

M-step:  $\pi^1, \theta^1 \rightarrow \mathcal{L}((\tau_i^1)_i; \pi^1, \theta^1) \geq \mathcal{L}((\tau_i^1)_i; \pi^*, \theta^*)$

E-step:  $\tau_i^2 \rightarrow \mathcal{L}((\tau_i^2)_i; \pi^1, \theta^1) = L_{(x_i)}(\pi^1, \theta^1) \geq \mathcal{L}((\tau_i^1)_i; \pi^1, \theta^1)$

Yoo, in summary:  $L_{(x_i)}(\pi^1, \theta^1) = \mathcal{L}((\tau_i^2)_i; \pi^1, \theta^1) \geq \mathcal{L}((\tau_i^1)_i; \pi^1, \theta^1) \geq \mathcal{L}((\tau_i^1)_i; \pi^*, \theta^*) = L_{(x_i)}(\pi^*, \theta^*)$

that is

$$L_{(x_i)}(\pi^1, \theta^1) \geq L_{(x_i)}(\pi^*, \theta^*)$$