

# The Expectation-Maximization (EM) Algorithm

**Algoritmul EM, fundamentare teoretică:  
pasul E [și pasul M]**

CMU, 2008 fall, Eric Xing, HW4, pr. 1.1-3

Algoritmul EM (Expectation-Maximization) permite crearea unor modele probabiliste care pe de o parte depind de un set de parametri  $\theta$  iar pe de altă parte includ pe lângă variabilele obișnuite (“observabile” sau “vizibile”)  $x$  și variabile necunoscute (“neobservabile”, “ascunse” sau “latente”)  $z$ .

În general, în astfel de situații/modele, nu se poate face în mod direct o estimare a parametrilor modelului ( $\theta$ ), în așa fel încât să se garanteze atingerea maximului verosimilității datelor observabile  $x$ .

În schimb, algoritmul EM procedează în manieră iterativă, constituind astfel o modalitate foarte convenabilă de estimare a parametrilor  $\theta$ .

Definim log-verosimilitatea datelor observabile ( $x$ ) ca fiind  $\log P(x \mid \theta)$ , iar log-verosimilitatea datelor *complete* (observabile,  $x$ , și neobservabile,  $z$ ) ca fiind  $\log P(x, z \mid \theta)$ .

*Observație:* Pe tot parcursul acestui exercițiu se va considera funcția log ca având baza supraunitară, fixată.

a. Log-verosimilitatea datelor *observabile* ( $x$ ) se poate exprima în funcție de datele neobservabile ( $z$ ), astfel:<sup>a</sup>

$$\ell(\theta) \stackrel{not.}{=} \log P(x \mid \theta) = \log \left( \sum_z P(x, z \mid \theta) \right)$$

În continuare vom nota cu  $q$  o funcție / distribuție de probabilitate definită peste variabilele ascunse/neobservabile  $z$ .

Folosiți *inegalitatea lui Jensen* pentru a demonstra că are loc următoarea inegalitate:

$$\log P(x \mid \theta) \geq \sum_z q(z) \log \left( \frac{P(x, z \mid \theta)}{q(z)} \right) \quad (1)$$

pentru orice  $x$  (fixat), pentru orice valoare a parametrului  $\theta$  și pentru orice distribuție probabilistă  $q$  definită peste variabilele neobservabile  $z$ .

---

<sup>a</sup>Rețineți că  $x$ , vectorul de date observabile, este fixat (dat), în vreme ce  $z$ , vectorul de date neobservabile, este liber (variabil).

## Observație (1)

Semnificația inegalității (1) este următoarea:

Funcția (de fapt, orice funcție de forma)

$$F(q, \theta) \stackrel{\text{def.}}{=} \sum_z q(z) \log \left( \frac{P(x, z | \theta)}{q(z)} \right)$$

constituie o margină inferioară pentru funcția de log-verosimilitate a datelor incomplete / observabile,  $\ell(\theta) \stackrel{\text{not.}}{=} \log P(x | \theta)$ .

Remarcați faptul că  $F$  este o funcție de două variabile, iar prima variabilă nu este de tip numeric (cum este  $\theta$ ), ci este de tip funcțional.

Mai mult, se observă că expresia funcției  $F$  este de fapt o medie,

$$E_{q(z)} \left[ \log \frac{P(x, z | \theta)}{q(z)} \right],$$

atunci când  $x$ ,  $q$  și parametrul  $\theta$  se consideră fixați, iar  $z$  este lăsat să varieze.

## Soluție

**Inegalitatea lui Jensen, în contextul teoriei probabilităților:**

**considerând  $X$  este o variabilă aleatoare (unară),**

**dacă  $\varphi$  este o funcție (reală) convexă, atunci  $\varphi(E[X]) \leq E[\varphi(X)]$ ;**

**dacă  $\varphi$  este funcție concavă, atunci  $\varphi(E[X]) \geq E[\varphi(X)]$ .**

**Aici vom folosi funcția  $\log$  cu bază supraunitară (funcție concavă), deci aplicând inegalitatea lui Jensen vom obține:  $\log(E[X]) \geq E[\log(X)]$ .**

**Log-verosimilitatea datelor observabile este:**

$$\begin{aligned} \ell(\theta) \stackrel{not.}{=} \log P(x | \theta) &= \log \left( \sum_z P(x, z | \theta) \right) = \log \left( \sum_z q(z) \frac{P(x, z | \theta)}{q(z)} \right) \\ &\stackrel{def.}{=} \log \left( E_{q(z)} \left[ \frac{P(x, z | \theta)}{q(z)} \right] \right) \end{aligned}$$

**Conform inegalității lui Jensen (înlocuind  $X$  de mai sus cu  $\frac{P(x, z | \theta)}{q(z)}$ ), rezultă:**

$$\ell(\theta) \stackrel{not.}{=} \log P(x | \theta) \geq E_{q(z)} \left[ \log \frac{P(x, z | \theta)}{q(z)} \right] \stackrel{def.}{=} \sum_z q(z) \log \frac{P(x, z | \theta)}{q(z)},$$

## Notăție

În continuare, pentru a vă aduce mereu aminte că distribuția  $q$  se referă la datele neobservabile  $z$ , vom folosi notația  $q(z)$  în loc de  $q$ .

În consecință, în cele ce urmează, în funcție de context,  $q(z)$  va desemna fie distribuția  $q$ , fie valoarea acestei distribuții pentru o valoare oarecare [a variabilei neobservabile]  $z$ .

(Este adevărat că această lejeră ambiguitate poate induce în eroare cititorul neexperimentat.)

b. Vă reamintim definiția entropiei relative (numită și divergența Kullback-Leibler):

$$KL(q(z) \parallel P(z \mid x, \theta)) = - \sum_z q(z) \log \left( \frac{P(z \mid x, \theta)}{q(z)} \right)$$

Arătați că

$$\log P(x \mid \theta) = F(q(z), \theta) + KL(q(z) \parallel P(z \mid x, \theta)).$$

Observație (2):

Semnificația egalității care trebuie demonstrată la acest punct este foarte interesantă: diferența dintre funcția obiectiv  $\ell(\theta) \stackrel{not.}{=} \log P(x \mid \theta)$  și marginea sa inferioară  $F(q(z), \theta)$  — a se vedea punctul a — este  $KL(q(z) \parallel P(z \mid x, \theta))$ . Tocmai pe această chestiune se va “construi” punctul final, și cel mai important, al problemei noastre.



## Observație (3)

Ideile de bază ale algoritmului EM sunt două:

1. În loc să calculeze maximul funcției de log-verosimilitate  $\log P(x | \theta)$  în raport cu  $\theta$ , algoritmul EM va maximiza marginea sa inferioară,  $F(q(z), \theta)$ , în raport cu ambele argumente,  $q(z)$  și  $\theta$ .
2. Pentru a căuta maximul (de fapt, un maxim local al) marginii inferioare  $F(q(z), \theta)$ , algoritmul EM aplică metoda *creșterii pe coordonate* (engl., coordinate ascent): după ce inițial se fixează  $\theta^{(0)}$  eventual aleatoriu, se maximizează *iterativ* funcția  $F(q(z), \theta)$ , în mod *alternativ*: mai întâi în raport cu distribuția  $q(z)$  și apoi în raport cu parametrul  $\theta$ .

$$\text{Pasul E:} \quad q^{(t)}(z) = \operatorname{argmax}_{q(z)} F(q(z), \theta^{(t)})$$

$$\text{Pasul M:} \quad \theta^{(t+1)} = \operatorname{argmax}_{\theta} F(q^{(t)}(z), \theta)$$

## Soluție

$$\begin{aligned}
 F(q(z), \theta) &\stackrel{\text{def.}}{=} \sum_z q(z) \log \left( \frac{P(x, z \mid \theta)}{q(z)} \right) \\
 &= \sum_z q(z) \log \left( \frac{P(z \mid x, \theta) \cdot P(x \mid \theta)}{q(z)} \right) \\
 &= \sum_z q(z) \left[ \log \frac{P(z \mid x, \theta)}{q(z)} + \log P(x \mid \theta) \right] \\
 &= \sum_z q(z) \log \left( \frac{P(z \mid x, \theta)}{q(z)} \right) + \sum_z q(z) \log P(x \mid \theta) \\
 &= -KL(q(z) \parallel P(z \mid x, \theta)) + \log P(x \mid \theta) \cdot \underbrace{\sum_z q(z)}_{=1}
 \end{aligned}$$

$$\Rightarrow \log P(x \mid \theta) = F(q(z), \theta) + KL(q(z) \parallel P(z \mid x, \theta)).$$

### Observație (4)

Conform proprietății  $KL(p \parallel q) \geq 0$  pentru  $\forall p, q$ , rezultă  $KL(q(z) \parallel P(z \mid x, \theta)) \geq 0$ .

Așadar, din egalitatea care tocmai a fost demonstrată la punctul  $b$  obținem (din nou!, după rezultatul de la punctul  $a$ ) că  $F(q(z), \theta)$  este o margine inferioară pentru log-verosimilitatea datelor observabile,  $\ell(\theta) \stackrel{not.}{=} \log P(x \mid \theta)$ .

c. Fie  $\theta^{(t)}$  valoarea obținută pentru parametrul / parametrii  $\theta$  la iterația  $t$  a algoritmului EM. Considerând această valoare fixată, arătați că maximul lui  $F$  în raport cu argumentul / distribuția  $q(z)$  este atins pentru distribuția  $P(z \mid x, \theta^{(t)})$ , iar valoarea maximului este:

$$\max_{q(z)} F(q(z), \theta^{(t)}) = E_{P(z|x, \theta^{(t)})} [\log P(x, z \mid \theta^{(t)})] + H(P(z \mid x, \theta^{(t)}))$$

## Soluție

Trebuie să maximizăm  $F(q(z), \theta^{(t)})$  — marginea inferioară a log-verosimilității datelor observabile  $x$  — în raport cu distribuția  $q(z)$ .

Pe de o parte, rezultatul de la punctul  $a$  ne spune că  $F(q(z), \theta) \leq \log P(x | \theta)$ , pentru orice valoare a lui  $\theta$ ; în particular, pentru  $\theta^{(t)}$  avem

$$\log P(x | \theta^{(t)}) \geq F(q(z), \theta^{(t)})$$

Pe de altă parte, dacă în egalitatea demonstrată la punctul  $b$  se înlocuiește  $\theta$  cu  $\theta^{(t)}$ , rezultă:

$$\log P(x | \theta^{(t)}) = F(q(z), \theta^{(t)}) + KL(q(z) || P(z | x, \theta^{(t)}))$$

În fine, dacă alegem  $q(z) = P(z | x, \theta^{(t)})$ , atunci termenul  $KL(q(z) || P(z | x, \theta^{(t)}))$  din dreapta egalității de mai sus devine zero (vezi exercițiul [PS-31] de la capitolul de Probabilități și statistică).

Așadar, valoarea  $\max_{q(z)} F(q(z), \theta^{(t)})$  se obține pentru distribuția  $q(z) = P(z | x, \theta^{(t)})$ .

Acum vom calcula această valoare maximă:

$$\begin{aligned}
 \max_{q(z)} F(q(z), \theta^{(t)}) &= \log P(x | \theta^{(t)}) \stackrel{\text{def.}}{=} \sum_z P(z | x, \theta^{(t)}) \log \left( \frac{P(x, z | \theta^{(t)})}{P(z | x, \theta^{(t)})} \right) \\
 &= E_{P(z|x, \theta^{(t)})} \left[ \log \frac{P(x, z | \theta^{(t)})}{P(z | x, \theta^{(t)})} \right] \\
 &= E_{P(z|x, \theta^{(t)})} [\log P(x, z | \theta^{(t)}) - \log P(z | x, \theta^{(t)})] \\
 &= E_{P(z|x, \theta^{(t)})} [\log P(x, z | \theta^{(t)})] - E_{P(z|x, \theta^{(t)})} [\log P(z | x, \theta^{(t)})] \\
 &= E_{P(z|x, \theta^{(t)})} [\log P(x, z | \theta^{(t)})] + H[P(z | x, \theta^{(t)})] \\
 &= Q(\theta^{(t)} | \theta^{(t)}) + H[P(z | x, \theta^{(t)})]
 \end{aligned}$$

unde  $Q(\theta | \theta^{(t)}) \stackrel{\text{not.}}{=} E_{P(z|x, \theta^{(t)})} [\log P(x, z | \theta)]$ .

## Observație (5)

Notând  $G_t(\theta) \stackrel{\text{def.}}{=} F(P(z \mid x, \theta^{(t)}), \theta)$ , din calculul de mai sus rezultă că  $G_t(\theta^{(t)}) = \log P(x \mid \theta^{(t)}) = Q(\theta^{(t)} \mid \theta^{(t)}) + H(P(z \mid x, \theta^{(t)}))$ . Se poate demonstra ușor — procedând similar cu calculul de mai sus — egalitatea

$$G_t(\theta) = Q(\theta \mid \theta^{(t)}) + H[P(z \mid x, \theta^{(t)})]$$

Observând că termenul  $H[P(z \mid x, \theta^{(t)})]$  din această ultimă egalitate nu depinde de  $\theta$ , rezultă imediat că

$$\operatorname{argmax}_{\theta} G_t(\theta) = \operatorname{argmax}_{\theta} Q(\theta \mid \theta^{(t)})$$

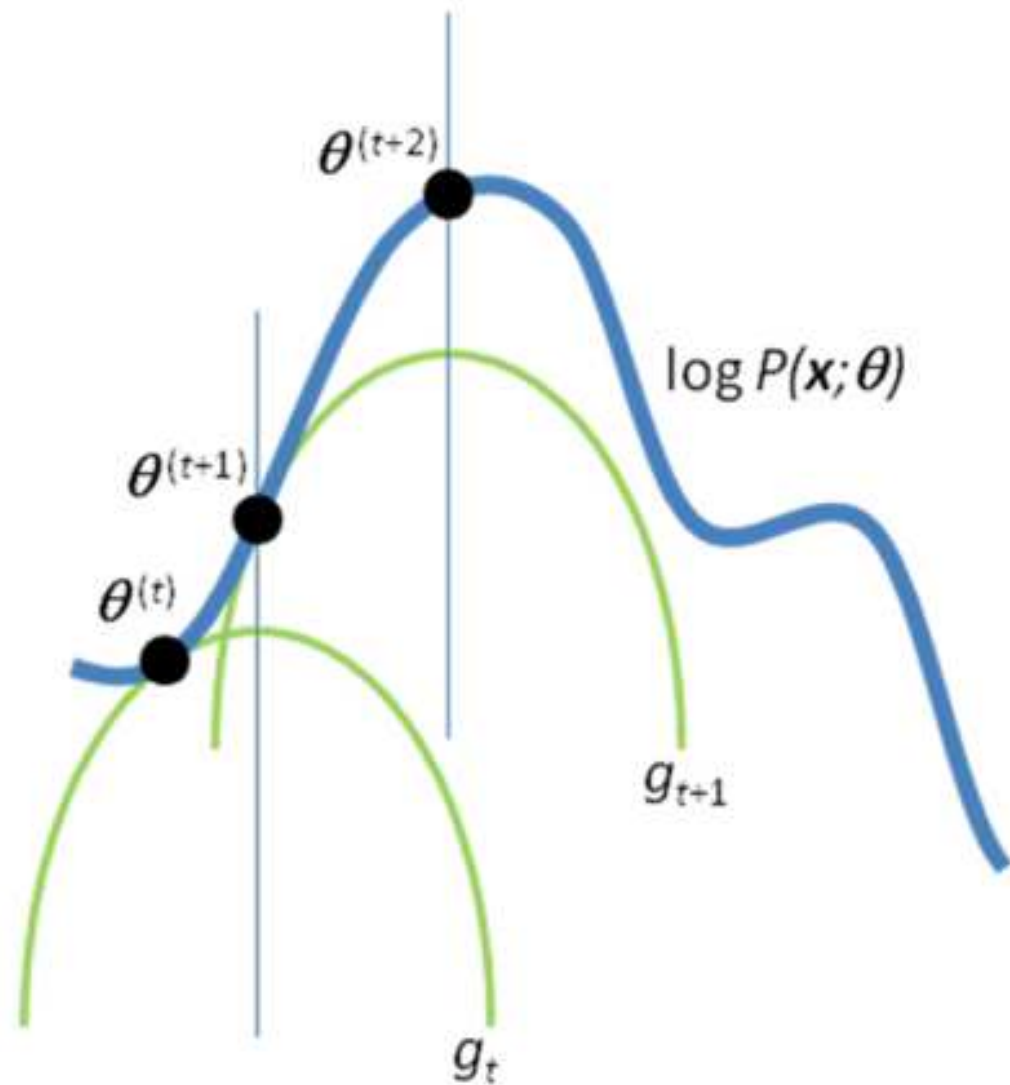
În consecință,

$$\theta^{(t+1)} \stackrel{\text{def.}}{=} \operatorname{argmax}_{\theta} F(P(z \mid x, \theta^{(t)}), \theta) = \operatorname{argmax}_{\theta} G_t(\theta) = \operatorname{argmax}_{\theta} Q(\theta \mid \theta^{(t)})$$

From:

*What is the expectation maximization algorithm?*

Chuong B. Do, Serafim Batzoglou,  
Nature Biotechnology,  
vol. 26, no. 8, 2008, pp. 897-899





Egalitatea precedentă este responsabilă pentru următoarea reformulare (cea uzuală!) a algoritmului EM:

Pasul E': calculează  $Q(\theta \mid \theta^{(t)}) = E_{P(z \mid x, \theta^{(t)})}[\log P(x, z \mid \theta)]$

Pasul M': calculează  $\theta^{(t+1)} = \operatorname{argmax}_{\theta} Q(\theta \mid \theta^{(t)})$

# Algoritmul EM: corectitudine / convergență

prelucrare de Liviu Ciortuz, după  
[en.wikipedia.org/wiki/Expectation-maximization](https://en.wikipedia.org/wiki/Expectation-maximization)

Pentru a maximiza funcția de log-verosimilitate a datelor observabile, i.e.  $\ell(\theta) \stackrel{\text{def.}}{=} \log P(x \mid \theta)$ , unde baza logaritmului (nespecificată) este considerată supraunitară, algoritmul EM procedează în mod iterativ, optimizând la pasul  $M$  al fiecărei iterații ( $t$ ) o funcție “auxiliară”

$$Q(\theta \mid \theta^{(t)}) \stackrel{\text{def.}}{=} E_{P(z|x, \theta^{(t)})}[\log P(x, z \mid \theta)],$$

reprezentând media log-verosimilității datelor complete (observabile și neobservabile) în raport cu distribuția condițională  $P(z \mid x, \theta^{(t)})$ .

Vom considera iterațiile  $t = 0, 1, \dots$  și  $\theta^{(t+1)} = \operatorname{argmax}_{\theta} Q(\theta \mid \theta^{(t)})$ , cu  $\theta^{(0)}$  ales în mod arbitrar.

**Demonstrați** că pentru orice  $t$  fixat (arbitrar) și pentru orice  $\theta$  astfel încât  $Q(\theta \mid \theta^{(t)}) \geq Q(\theta^{(t)} \mid \theta^{(t)})$  are loc inegalitatea:

$$\log P(x \mid \theta) - \log P(x \mid \theta^{(t)}) \geq Q(\theta \mid \theta^{(t)}) - Q(\theta^{(t)} \mid \theta^{(t)}) \quad (2)$$

## Observații

1. Semnificația imediată a relației (2):

Orice îmbunătățire a valorii funcției auxiliare  $Q(\theta \mid \theta^{(t)})$  conduce la o îmbunătățire cel puțin la fel de mare a valorii funcției obiectiv,  $\ell(\theta)$ .

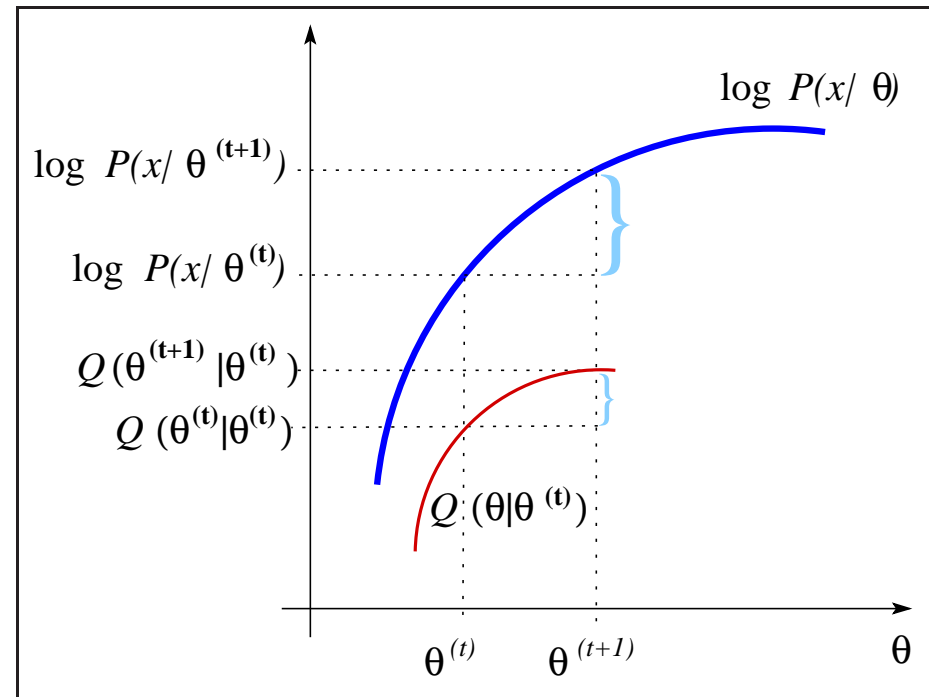
2. Dacă în inegalitatea (2) se înlocuiește  $\theta$  cu  $\theta^{(t+1)} = \operatorname{argmax}_{\theta} Q(\theta \mid \theta^{(t)})$ , va rezulta

$$\log P(x \mid \theta^{(t+1)}) \geq \log P(x \mid \theta^{(t)}).$$

În final, vom avea

$$\ell(\theta^{(0)}) \leq \ell(\theta^{(t)}) \leq \ell(\theta^{(t+1)}) \leq \dots$$

Șirul acesta (monoton) este mărginit superior de 0 (vezi definiția lui  $\ell$ ), deci converge la o anumită valoare  $\ell^*$ . În anumite cazuri / condiții, această valoare este un maxim (în general, local) al funcției de log-verosimilitate.



## Observații (cont.)

3. Conform aceleiași inegalități (2), la pasul M de la iterația  $t$  a algoritmului EM, este suficient ca în loc să se ia  $\theta^{(t+1)} = \operatorname{argmax}_{\theta} Q(\theta \mid \theta^{(t)})$ , să se aleagă  $\theta^{(t+1)}$  astfel încât  $Q(\theta^{(t+1)} \mid \theta^{(t)}) > Q(\theta^{(t)} \mid \theta^{(t)})$ . Aceasta constituie *versiunea “generalizată”* a algoritmului EM.

## Demonstrația relației (2)

$$P(x, z \mid \theta) = P(z \mid x, \theta) \cdot P(x \mid \theta) \Rightarrow \log P(x \mid \theta) = \log P(x, z \mid \theta) - \log P(z \mid x, \theta) \Rightarrow$$

$$\underbrace{\sum_z P(z \mid x, \theta^{(t)}) \cdot \log P(x \mid \theta)}_1 =$$

$$\sum_z P(z \mid x, \theta^{(t)}) \cdot \log P(x, z \mid \theta) - \sum_z P(z \mid x, \theta^{(t)}) \cdot \log P(z \mid x, \theta) \Rightarrow$$

$$\log P(x \mid \theta) = Q(\theta \mid \theta^{(t)}) - \sum_z P(z \mid x, \theta^{(t)}) \cdot \log P(z \mid x, \theta)$$

Ultimul termen din egalitatea aceasta reprezintă o cross-entropie, pe care o vom nota cu  $CH(\theta \mid \theta^{(t)})$ . Așadar,

$$\log P(x \mid \theta) = Q(\theta \mid \theta^{(t)}) + CH(\theta \mid \theta^{(t)})$$

Această egalitate este valabilă pentru toate valorile posibile ale parametrului  $\theta$ . În particular pentru  $\theta = \theta^{(t)}$ , vom avea:

$$\log P(x \mid \theta^{(t)}) = Q(\theta^{(t)} \mid \theta^{(t)}) + CH(\theta^{(t)} \mid \theta^{(t)})$$

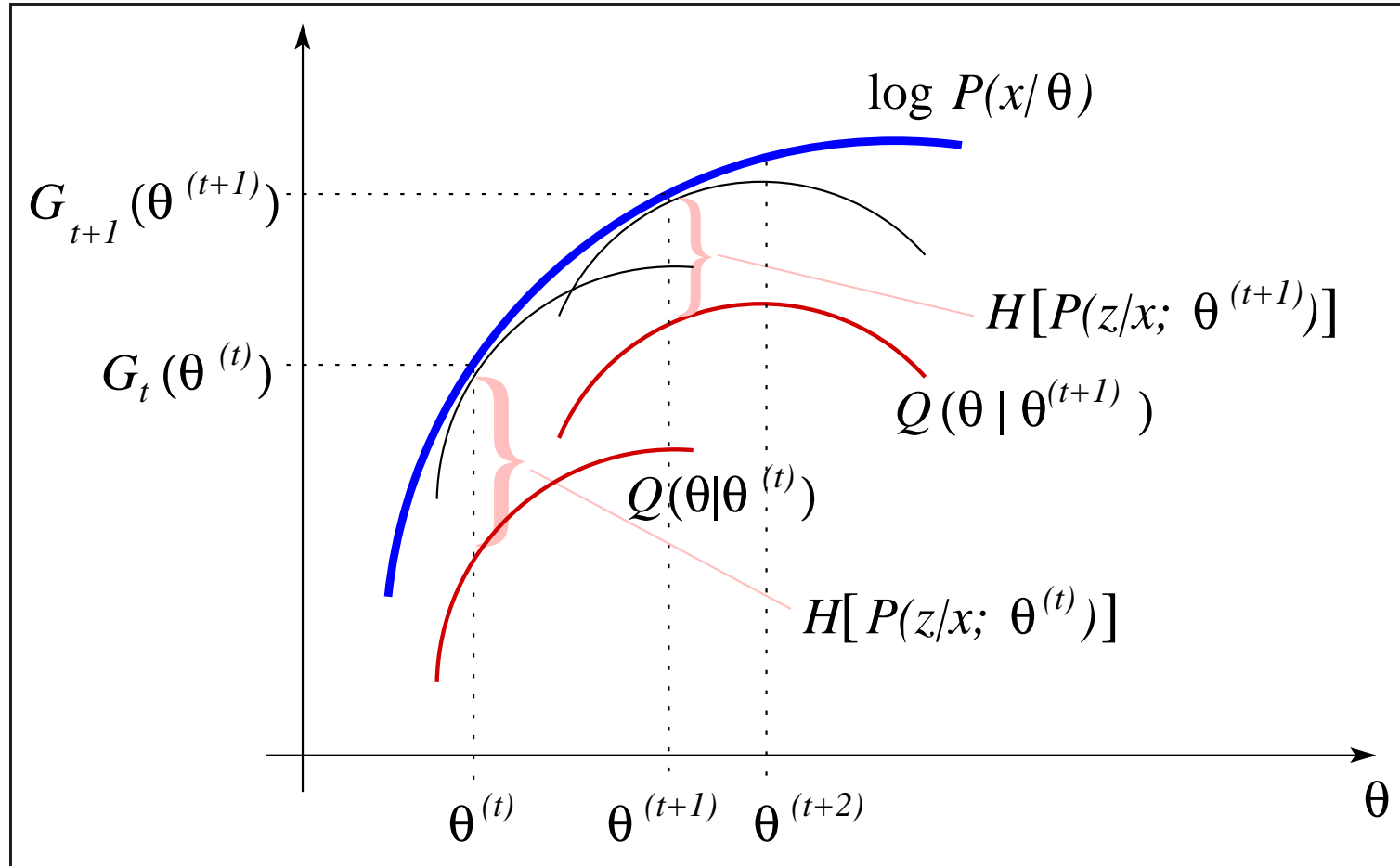
## Demonstrația relației (2), cont.

Scăzând membru cu membru ultimele două egalități, obținem:

$$\log P(x \mid \theta) - \log P(x \mid \theta^{(t)}) = Q(\theta \mid \theta^{(t)}) - Q(\theta^{(t)} \mid \theta^{(t)}) + CH(\theta \mid \theta^{(t)}) - CH(\theta^{(t)} \mid \theta^{(t)})$$

Conform inegalității lui Gibbs, avem  $CH(\theta \mid \theta^{(t)}) \geq CH(\theta^{(t)} \mid \theta^{(t)})$ , deci în final rezultă:

$$\log P(x \mid \theta) - \log P(x \mid \theta^{(t)}) \geq Q(\theta \mid \theta^{(t)}) - Q(\theta^{(t)} \mid \theta^{(t)})$$





Using the EM algorithm for  
*learning a categorical distribution*  
[and implicitly a multinomial distribution]

Application to the ABO blood model

Liviu Ciortuz, following

Anirban DasGupta, *Probability for Statistics and Machine Learning*, Springer, 2011, Ex. 20.10

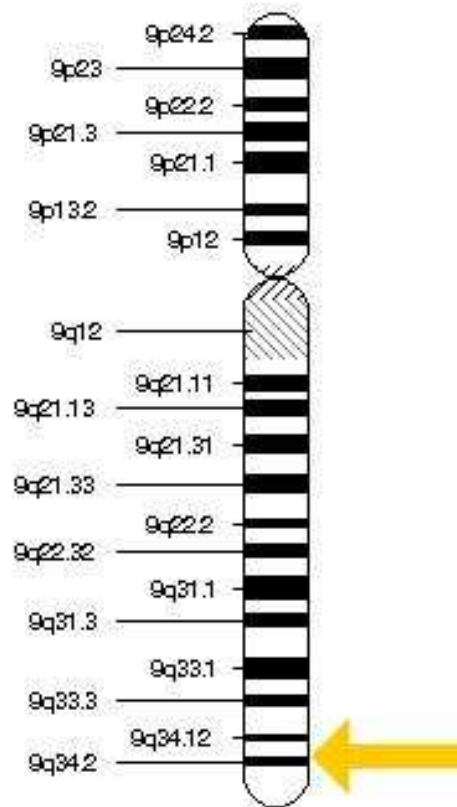
Grupele sangvine ale oamenilor sunt determinate de variantele („alelele“) unei gene situate pe cromozomul 9 (mai exact, la poziția 9q34.2), numită gena *ABO*. Se știe că fiecare dintre noi dispunem de câte o pereche de astfel de cromozomi, deci de câte două copii ale genei *ABO*, și anume una moștenită de la tată și una moștenită de la mamă.

Se notează cu *A*, *B* și *O* cele trei tipuri de alele ale genei *ABO*. Alelele *A* și *B* sunt *dominante* în raport cu alela *O*. (Altfel spus, alela *O* este recesivă în raport cu fiecare dintre alelele *A* și *B*.) Alelele *A* și *B* sunt *codominante*. Prin urmare, grupele sangvine pot fi specificate conform tabelului alăturat.

grupe sangvine	alele moștenite
<i>A</i>	<i>AA</i> <i>AO</i>
<i>B</i>	<i>BB</i> <i>BO</i>
<i>AB</i>	<i>AB</i>
<i>O</i>	<i>OO</i>

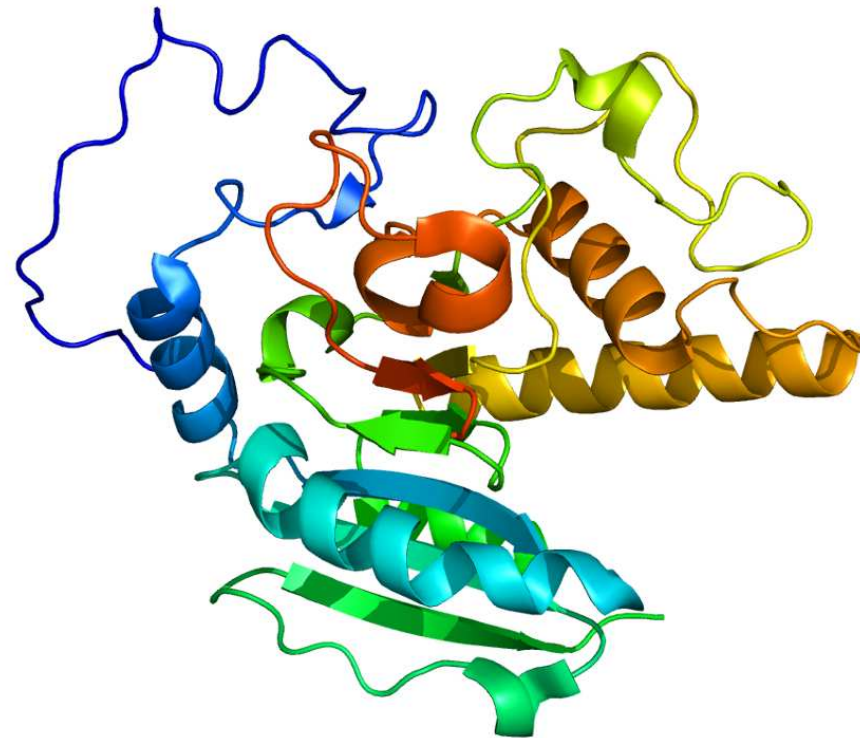
Presupunem că, într-o populație oarecare, fiecare dintre alelele *A*, *B* și *O* are la bărbați aceeași frecvență ca și la femei. În cele ce urmează, aceste frecvențe (notate respectiv cu  $p_A$ ,  $p_B$  și  $p_O$ ) vor fi considerate a priori necunoscute, dar urmează să le determinăm.

## Location of the ABO gene on chromosome 9

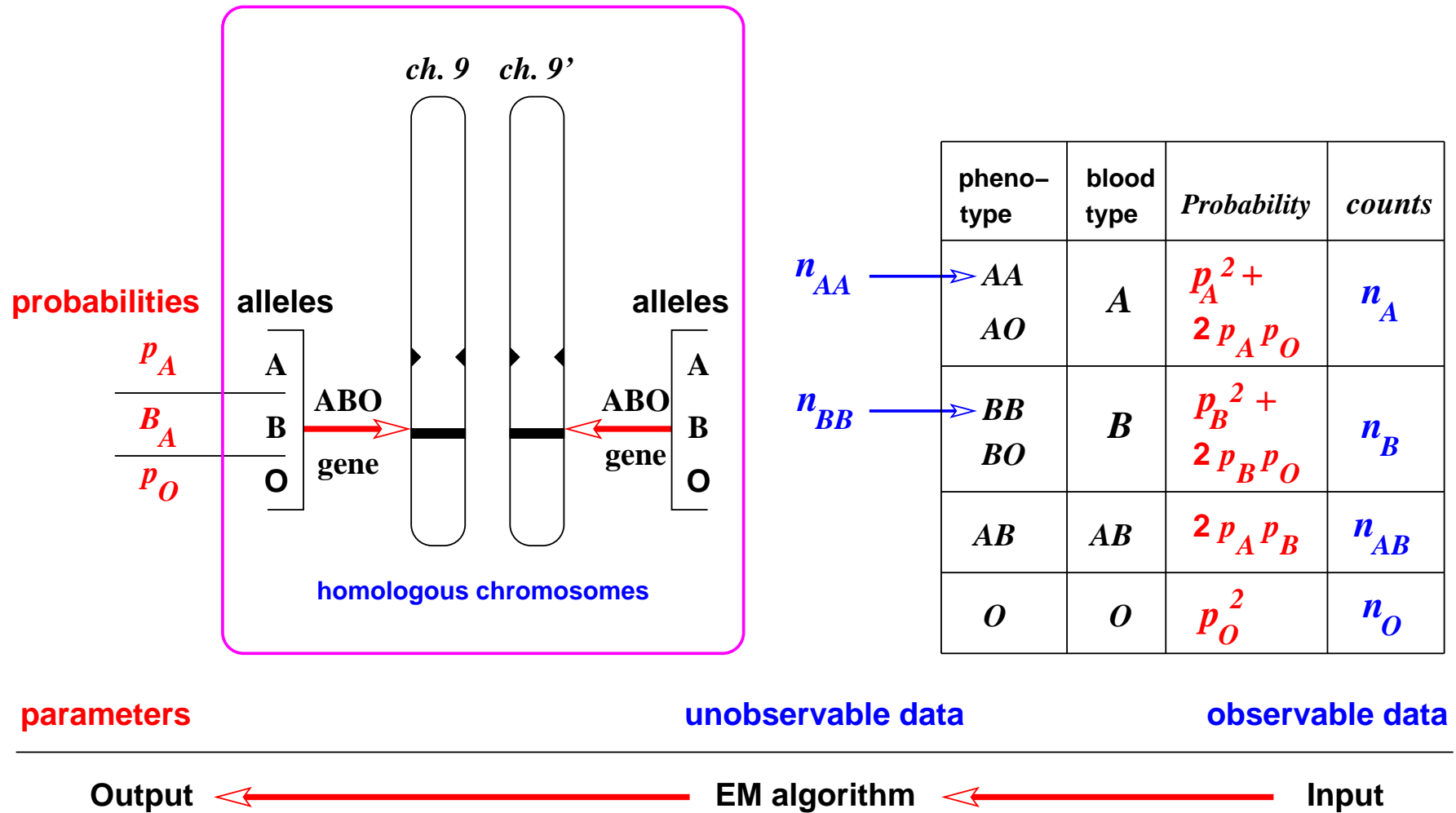


Source: wiki

## The ABO protein



Source: wiki



a. Considerăm că într-un eșantion de populație format din  $n$  persoane sunt  $n_A$  persoane cu grupa sanguină  $A$ ,  $n_B$  persoane cu grupa sanguină  $B$ ,  $n_{AB}$  persoane cu grupa sanguină  $AB$  și  $n_O$  persoane cu grupa sanguină  $O$ . (Evident,  $n = n_A + n_B + n_{AB} + n_O$ .)

Pornind de la aceste date „observabile“, să se deriveze *algoritmul EM* pentru determinarea probabilităților  $p_A$ ,  $p_B$  și  $p_O$ . (Evident, întrucât suma lor este 1, va fi suficient să se determine doar două dintre ele.)

## Indicații

1. Presupunând că procesul de asociere a alelelor moștenite de către un individ de la părinții lui respectă proprietățile specifice evenimentelor aleatoare independente, rezultă că probabilitățile de realizare a combinațiilor (în termeni genetici: „fenotipurile“)  $AA$ ,  $AO$ ,  $BB$ ,  $BO$ ,  $AB$  și  $OO$  într-o populație oarecare sunt  $p_A^2$ ,  $2p_Ap_O$ ,  $p_B^2$ ,  $2p_Bp_O$ ,  $2p_Ap_B$ , și respectiv  $p_O^2$ .
2. Considerând  $n_A = n_{AA} + n_{AO}$  și  $n_B = n_{BB} + n_{BO}$ , unde semnificațiile numerelor  $n_{AA}$ ,  $n_{AO}$ ,  $n_{BB}$  și  $n_{BO}$  sunt similare cu semnificațiile numerelor  $n_A$ ,  $n_B$ ,  $n_{AB}$ , și  $n_O$  care au fost precizate mai sus, este natural ca în formularea algoritmului EM datele  $n_{AA}$  și  $n_{BB}$  să fie considerate „neobservabile“. Ca *parametri* ai modelului, se vor considera probabilitățile  $p_A$ ,  $p_B$  și  $p_O$ .

3. La pasul E al algoritmului EM veți calcula  $\hat{n}_{AA}$  și  $\hat{n}_{BB}$ , care reprezintă respectiv numărul „așteptat“ de apariții ale combinației de alele  $AA$  și numărul „așteptat“ de apariții ale combinației de alele  $BB$  în populația dată.

Veți scrie apoi funcția de log-verosimilitate a datelor complete („observabile“ și „neobservabile“), exprimată cu ajutorul distribuției  $Multinomial(n; p_A^2, 2p_Ap_O, p_B^2, 2p_Bp_O, 2p_Ap_B, p_O^2)$ .

4. La pasul M al algoritmului EM, pornind de la media funcției de log-verosimilitate care a fost calculată la pasul E, veți stabili regulile de actualizare pentru probabilitățile  $p_A$ ,  $p_B$  și  $p_O$ .

Atenție: suma acestor probabilități fiind 1, problema de optimizare pe care va trebui să o rezolvați la pasul M este una cu restricții. În acest sens, metoda multiplicatorilor lui Lagrange vă poate fi de folos.

## Solution

### Notations:

- **parameters:**  $p = \{p_A, p_B, p_O\}$   
(in fact, it will be enough to estimate  $p_A$  and  $p_B$ , since  $p_A + p_B + p_O = 1$ )
- **observable data:**  $n_{obs} = \{n_A, n_B, n_{AB}, n_O\}$
- **unobservable data:**  $n_{unobs} = \{n_{AA}, n_{AO}, n_{BB}, n_{BO}\}$   
(in fact, we will restrict  $n_{unobs}$  to  $\{n_{AA}, n_{BB}\}$ )
- **complete data:**  $n_{compl} = n_{obs} \cup n_{unobs}$
- $n = n_A + n_B + n_{AB} + n_O, \quad n_A = n_{AA} + n_{AO}, \quad n_B = n_{BB} + n_{BO}.$

### The likelihood of complete data:

$$L(p) \stackrel{not.}{=} P(n_{compl}|p) = \frac{n!}{n_{AA}! n_{AO}! n_{BB}! n_{BO}! n_{AB}! n_O!} \cdot (p_A^2)^{n_{AA}} \cdot (2 p_A p_O)^{n_{AO}} \cdot (p_B^2)^{n_{BB}} \cdot (2 p_B p_O)^{n_{BO}} \cdot (2 p_A p_B)^{n_{AB}} \cdot (p_O^2)^{n_O}$$



**The likelihood function:**

$$\begin{aligned}
\ell(p) &\stackrel{\text{def.}}{=} \ln L(p) \\
&= \ln c + n_{AA} \ln(p_A^2) + n_{AO} \ln(2 p_A p_O) + n_{BB} \ln(p_B^2) + n_{BO} \ln(2 p_B p_O) + n_{AB} \ln(2 p_A p_B) + n_O \ln(p_O^2) \\
&= \ln c' + 2 n_{AA} \ln p_A + n_{AO} (\ln p_A + \ln p_O) + \\
&\quad 2 n_{BB} \ln p_B + n_{BO} (\ln p_B + \ln p_O) + n_{AB} (\ln p_A + \ln p_B) + 2 n_O \ln p_O \\
&= \ln c' + 2 n_{AA} \ln p_A + (n_A - n_{AA}) (\ln p_A + \ln p_O) + \\
&\quad 2 n_{BB} \ln p_B + (n_B - n_{BB}) (\ln p_B + \ln p_O) + n_{AB} (\ln p_A + \ln p_B) + 2 n_O \ln p_O,
\end{aligned}$$

where  $c$  and  $c'$  are constants which do not depend on the  $p$  parameter.

**The “auxiliary” function:**

$$\begin{aligned}
Q(p|p^{(t)}) &\stackrel{\text{def.}}{=} E[\ell(p)|n_{obs}; p^{(t)}] \\
&= \ln c + 2 \hat{n}_{AA} \ln p_A + (n_A - \hat{n}_{AA}) (\ln p_A + \ln p_O) + \\
&\quad 2 \hat{n}_{BB} \ln p_B + (n_B - \hat{n}_{BB}) (\ln p_B + \ln p_O) + n_{AB} (\ln p_A + \ln p_B) + 2 n_O \ln p_O,
\end{aligned}$$

where

$$\begin{aligned}
\hat{n}_{AA} &\stackrel{\text{not.}}{=} E[n_{AA}|n_{obs}; p^{(t)}] = E[n_{AA}|n_A, n_B, n_{AB}, n_O; p_A^{(t)}, p_B^{(t)}, p_O^{(t)}] \\
\hat{n}_{BB} &\stackrel{\text{not.}}{=} E[n_{BB}|n_{obs}; p^{(t)}] = E[n_{BB}|n_A, n_B, n_{AB}, n_O; p_A^{(t)}, p_B^{(t)}, p_O^{(t)}]
\end{aligned}$$

**E step:** As indicated by the expression of  $Q$ , here we have to compute  $\hat{n}_{AA}$  and  $\hat{n}_{BB}$  the *expected number* of individuals having the *phenotype* (aka, the allele pair)  $AA$ , and respectively  $BB$ .

Firstly, taking into account that  $n_{AA}$  (or, equivalently, the phenotype  $AA$ ), when seen as a random variable, follows a *binomial distribution* of parameters  $n_A$  and

$\frac{(p_A^{(t)})^2}{(p_A^{(t)})^2 + 2 p_A^{(t)} p_O^{(t)}}$ , its expectation will be:

$$\begin{aligned}\hat{n}_{AA} &\stackrel{not.}{=} E[n_{AA} | n_A, n_B, n_{AB}, n_O; p_A^{(t)}, p_B^{(t)}, p_O^{(t)}] \\ &= \frac{(p_A^{(t)})^2}{(p_A^{(t)})^2 + 2 p_A^{(t)} p_O^{(t)}} \cdot n_A\end{aligned}$$

Similarly,

$$\begin{aligned}\hat{n}_{BB} &\stackrel{not.}{=} E[n_{BB} | n_A, n_B, n_{AB}, n_O; p_A^{(t)}, p_B^{(t)}, p_O^{(t)}] \\ &= \frac{(p_B^{(t)})^2}{(p_B^{(t)})^2 + 2 p_B^{(t)} p_O^{(t)}} \cdot n_B\end{aligned}$$

### M step:

Given the constraint  $p_A + p_B + p_O = 1$ , we will introduce the Lagrange variable/“multiplier”  $\lambda$  and solve for the optimisation problem

$$\begin{aligned}
 p^{(t+1)} &\stackrel{not.}{=} (p_A^{(t+1)}, p_B^{(t+1)}, p_O^{(t+1)}) \\
 &= \operatorname{argmax}_{p_A, p_B, p_O} [Q(p_A, p_B, p_O | p_A^{(t)}, p_B^{(t)}, p_O^{(t)}) + \lambda(1 - (p_A + p_B + p_O))] \\
 &= \operatorname{argmax}_{p_A, p_B, p_O} [\ln c' + 2 \hat{n}_{AA} \ln p_A + (n_A - \hat{n}_{AA})(\ln p_A + \ln p_O) + 2 \hat{n}_{BB} \ln p_B + \\
 &\quad (n_B - \hat{n}_{BB})(\ln p_B + \ln p_O) + n_{AB}(\ln p_A + \ln p_B) + 2 n_O \ln p_O + \\
 &\quad \lambda(1 - (p_A + p_B + p_O))]
 \end{aligned}$$

Taking the partial derivatives of the objective function w.r.t.  $p_A$ ,  $p_B$  and  $p_O$  and solving them leads to:

$$\frac{1}{p_A} (2 \hat{n}_{AA} + n_A - \hat{n}_{AA} + n_{AB}) - \lambda = 0 \Rightarrow \hat{p}_A = \frac{1}{\lambda} (\hat{n}_{AA} + n_A + n_{AB})$$

$$\frac{1}{p_B} (2 \hat{n}_{BB} + n_B - \hat{n}_{BB} + n_{AB}) - \lambda = 0 \Rightarrow \hat{p}_B = \frac{1}{\lambda} (\hat{n}_{BB} + n_B + n_{AB})$$

$$\frac{1}{p_O} (n_A - \hat{n}_{AA} + n_B - \hat{n}_{BB} + 2 n_O) - \lambda = 0 \Rightarrow \hat{p}_O = \frac{1}{\lambda} (n_A - \hat{n}_{AA} + n_B - \hat{n}_{BB} + 2 n_O)$$

Now, enforcing the constraint  $\hat{p}_A + \hat{p}_B + \hat{p}_O = 1$  gives

$$\begin{aligned} \frac{1}{\lambda}(\hat{n}_{AA} + n_A + n_{AB} + \hat{n}_{BB} + n_B + n_{AB} + n_A - \hat{n}_{AA} + n_B - \hat{n}_{BB} + 2n_O) &= 1 \Leftrightarrow \\ \frac{2}{\lambda}(n_A + n_B + n_{AB} + n_O) &= 1 \end{aligned}$$

Since  $n = n_A + n_B + n_{AB} + n_O$ , it follows immediately that  $\frac{1}{\lambda}2n = 1$ .  
Therefore  $\lambda = 2n$ .

Together with the results obtained on the previous slide, this leads to the following *updating relations*:

$$\begin{aligned} \hat{p}_A^{(t+1)} &= \frac{1}{2n}(\hat{n}_{AA} + n_A + n_{AB}) \\ \hat{p}_B^{(t+1)} &= \frac{1}{2n}(\hat{n}_{BB} + n_B + n_{AB}) \\ \hat{p}_O^{(t+1)} &= \frac{1}{2n}(n_A - \hat{n}_{AA} + n_B - \hat{n}_{BB} + 2n_O) = \frac{1}{2n}(n_{AO} + n_{BO} + 2n_O) \end{aligned}$$

b. Implementați algoritmul EM pe care l-ați conceput la punctul a și rulați-l pentru input-ul  $n_A = 186$ ,  $n_B = 38$ ,  $n_{AB} = 13$ ,  $n_O = 284$  ( $n = 521$ ).

Ca *valori inițiale* pentru parametri, veți lucra mai întâi cu  $p_A = p_B = p_O = \frac{1}{3}$ , iar apoi cu  $p_A = p_O = 0.1$  și  $p_B = 0.98$ .

Pentru *oprire*, veți cere ca  $p_A^{(t)}$ ,  $p_B^{(t)}$  și  $p_O^{(t)}$  să nu difere (fiecare în parte) cu mai mult de  $10^{-4}$  față de valorile calculate la iterația precedentă. Comparați rezultatele obținute pentru fiecare din cele două inițializări.

**Starting values:**

$$p_A = p_B = p_O = 1/3.$$

**Iterations:**

$t$	$p_A$	$p_B$	$p_C$
1	0.2505	0.0611	0.6884
2	0.2185	0.0505	0.7311
3	0.2142	0.0502	0.7357
4	0.2137	0.0501	0.7362
5	0.2136	0.0501	0.7363
6	0.2136	0.0501	0.7363

**Starting values:**

$$p_A = p_O = 0.1, p_B = 0.98.$$

**Iterations:**

$t$	$p_A$	$p_B$	$p_C$
1	0.2505	0.0847	0.6648
2	0.2193	0.0511	0.7296
3	0.2143	0.0502	0.7355
4	0.2137	0.0501	0.7362
5	0.2136	0.0501	0.7363
6	0.2136	0.0501	0.7363

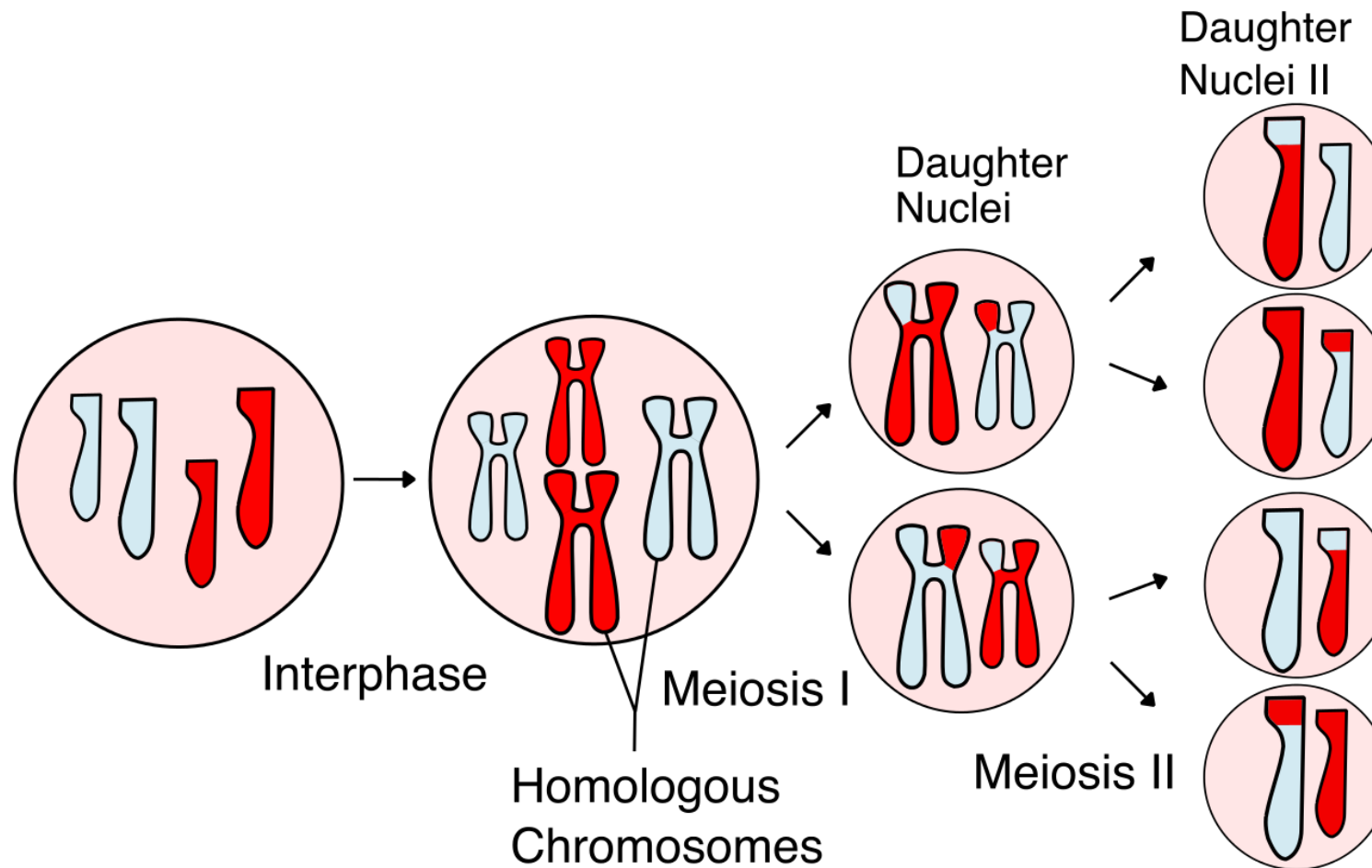
**Note:** The final results are the same.

Using the EM algorithm for  
*learning a multinomial distribution*  
which depends on a single parameter

Application to a bioinformatics task:  
computing probabilities for two linked bi-allelic loci  
in haploid cells

Liviu Ciortuz, following  
Brani Vidakovic, Georgia Institute of Technology,  
Bayesian Statistics course (ISyE 8843A), 2004,  
Handout 12, sec. 1.2.1

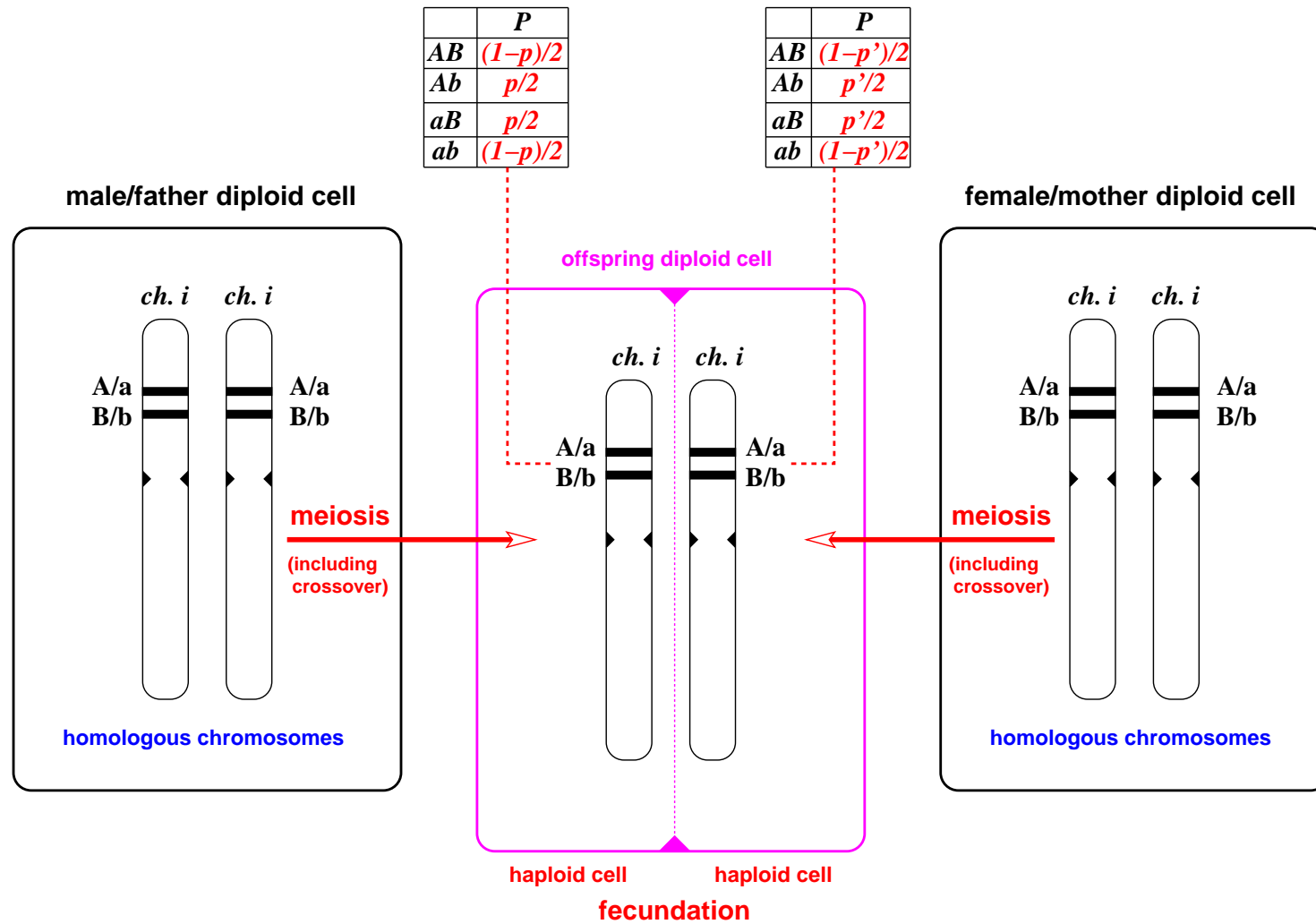
## Biological background: Meiosis



Source: wiki



# Biological background (cont'd)



## Biological background (cont'd)

Note: Alleles A and B are dominant w.r.t. a and respectively b.

Genotype:

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
AA BB	Aa BB	aA BB	AA Bb	Aa Bb	aA Bb	AA bB	Aa bB	aA bB	AA bb	Aa bb	aA bb	aa BB	aa Bb	aa bB	aa bb
<b>AB</b>									<b>Ab</b>			<b>aB</b>		<b>ab</b>	
$n_{AB}$									$n_{Ab}$			$n_{aB}$		$n_{ab}$	
$\frac{2 + \psi}{4}$									$\frac{1 - \psi}{4}$			$\frac{1 - \psi}{4}$		$\frac{\psi}{4}$	

with  $\psi \stackrel{not.}{=} (1 - p)(1 - p')$ .

We have designated in **magenta** colour the **offspring phenotype**, in **blue** the **counts**, i.e. number of individuals of having a certain phenotype in a given population, and in **red** the corresponding **probabilities**.

Note that

$$P(ab) = \frac{(1 - p)(1 - p')}{4} \text{ and}$$

$$P(Ab) = P(aB) = \frac{p}{2} \cdot \frac{p'}{2} + \frac{p}{2} \cdot \frac{1 - p'}{2} + \frac{1 - p'}{2} \cdot \frac{p}{2} = \frac{1 - (1 - p)(1 - p')}{4}.$$

Fie o variabilă aleatoare  $X$  care ia valori în mulțimea  $\{v_1, v_2, v_3, v_4\}$  și urmează distribuția  $Multinomial\left(n; \frac{2+\psi}{4}, \frac{1-\psi}{4}, \frac{1-\psi}{4}, \frac{\psi}{4}\right)$ , unde  $\psi$  este un parametru cu valori în intervalul  $(0, 1)$ . Cele patru probabilități listate în definiția acestei distribuții multinomiale sunt în corespondență directă cu valorile  $v_1, v_2, v_3$  și  $v_4$ .

a. Presupunem că  $n_1, n_2, n_3$  și  $n_4$  sunt numărul de „realizări” ale valorilor  $v_1, v_2, v_3$  și  $v_4$  în totalul celor  $n$  „observații”. Pentru fixarea ideilor, vom considera  $n_1 = 125, n_2 = 18, n_3 = 20$  și  $n_4 = 34$ .

Calculați estimarea de verosimilitate maximă (MLE) a parametru-lui  $\psi$ .

## Solution

By denoting  $n = n_1 + n_2 + n_3 + n_4$ ,  
 and knowing that  $n \sim \text{Multinomial}\left(n; \frac{2+\psi}{4}, \frac{1-\psi}{4}, \frac{1-\psi}{4}, \frac{\psi}{4}\right)$ ,  
 it follows that the verosimilarity function will be

$$L(\psi) \stackrel{\text{def.}}{=} P(D|\psi) = \frac{n!}{n_1! n_2! n_3! n_4!} \cdot \left(\frac{2+\psi}{4}\right)^{n_1} \cdot \left(\frac{1-\psi}{4}\right)^{n_2} \cdot \left(\frac{1-\psi}{4}\right)^{n_3} \cdot \left(\frac{\psi}{4}\right)^{n_4},$$

while

$$\ell(\psi) \stackrel{\text{def.}}{=} \ln L(\psi) = \ln c + n_1 \ln(2+\psi) + (n_2 + n_3) \ln(1-\psi) + n_4 \ln(\psi),$$

where  $c$  is constant w.r.t.  $\psi$ .

$$\begin{aligned}
\frac{\partial}{\partial \psi} \ell(\psi) &= \frac{n_1}{2+\psi} - \frac{n_2+n_3}{1-\psi} + \frac{n_4}{\psi} = \frac{n_1 \cdot \psi(1-\psi) - (n_2+n_3) \cdot \psi(2+\psi) + n_4 \cdot (2+\psi)(1-\psi)}{(2+\psi)(1-\psi)\psi} \\
&= \frac{n_1\psi - n_1\psi^2 - 2n_2\psi - n_2\psi^2 - 2n_3\psi - n_3\psi^2 + 2n_4 - 2n_4\psi + n_4\psi - n_4\psi^2}{(2+\psi)(1-\psi)\psi} \\
&= \frac{-\psi^2(n_1+n_2+n_3+n_4) + \psi(n_1-2n_2-2n_3-n_4) + 2n_4}{(2+\psi)(1-\psi)\psi} \\
&= \frac{-n\psi^2 + \psi(n_1-2n_2-2n_3-n_4) + 2n_4}{(2+\psi)(1-\psi)\psi}
\end{aligned}$$

**Note that the denominator  $(2+\psi)(1-\psi)\psi$  is always positive.**

**By substituting  $n_1, n_2, n_3, n_4$  for 125, 18, 20 and respectively 34, the nominator becomes  $-197\psi^2 + 15\psi + 68$ .**

**By analysing the sign of this expression, it is easy to see that for our data the function  $\ell(\psi)$  reaches its maximum in the interval  $(0, 1)$  for**

$$\hat{\psi} = \frac{-15 + \sqrt{15^2 + 4 \cdot 197 \cdot 68}}{-2 \cdot 197} = \frac{-15 + \sqrt{225 + 53809}}{-394} = \frac{-15 + 231.967}{-394} = 0.626769036.$$

b. Fie acum variabila aleatoare  $X' \sim Multinomial\left(n; \frac{1}{2}, \frac{\psi}{4}, \frac{1-\psi}{4}, \frac{1-\psi}{4}, \frac{\psi}{4}\right)$ . Valorile luate de variabila  $X'$ , corespunzător acestor cinci probabilități, sunt (în ordine)  $v_1, v_1, v_2, v_3$  și  $v_4$ .

Observați că, în raport cu distribuția multinomială de la punctul a, am „descompus” probabilitatea  $\frac{2+\psi}{4}$  în două probabilități,  $\frac{1}{2}$  și  $\frac{\psi}{4}$ , ca și cum ele ar corespunde unor evenimente disjuncte.

Vom considera  $n_{11}, n_{12}, n_2, n_3$  și respectiv  $n_4$  numărul de „realizări” ale acestor valori, însă de data aceasta  $n_{11}$  și  $n_{12}$  vor fi neobservabile. În schimb, vom furniza suma  $n_1 = n_{11} + n_{12}$  ca dată „observabilă”, alături de celelalte date observabile,  $n_2, n_3$  și  $n_4$ .

Să se estimeze parametrul  $\psi$  folosind algoritmul EM. Cum este această estimare față de valoarea obținută la punctul a?

*Indicație:* Veți elabora formulele corespunzătoare pasului E și pasului M. Apoi veți face o implementare și veți rula algoritmul EM (pornind, de exemplu, cu valoarea inițială 0.5 pentru  $\psi$ ) până când valorile acestui parametru până la cea de-a șasea zecimală nu se mai modifică.

## Solution: E-step

The verosimilarity function is now

$$\begin{aligned}
 P(D|\psi) &= \frac{n!}{n_{11}! n_{12}! n_2! n_3! n_4!} \cdot \left(\frac{1}{2}\right)^{n_{11}} \cdot \left(\frac{\psi}{4}\right)^{n_{12}} \cdot \left(\frac{1-\psi}{4}\right)^{n_2} \cdot \left(\frac{1-\psi}{4}\right)^{n_3} \cdot \left(\frac{\psi}{4}\right)^{n_4} \\
 &= \frac{n!}{n_{11}! n_{12}! n_2! n_3! n_4!} \cdot \left(\frac{1}{2}\right)^{n_{11}} \cdot \left(\frac{1-\psi}{4}\right)^{n_2+n_3} \cdot \left(\frac{\psi}{4}\right)^{n_{12}+n_4},
 \end{aligned}$$

while the log-verosimilarity function can be written as

$$\ln P(D|\psi) = \ln d + (n_2 + n_3) \ln(1 - \psi) + (n_{12} + n_4) \ln(\psi),$$

where  $d$  is constant w.r.t.  $\psi$ .

Therefore, the „auxiliary“ function will be

$$\begin{aligned}
 Q(\psi|\psi^{(t+1)}) &= E_{n_{12}|n_1, n_2, n_3, n_4, \psi^{(t)}} [\ln P(D|\psi)] \\
 &= \ln d + (n_2 + n_3) \ln(1 - \psi) + (\hat{n}_{12} + n_4) \ln(\psi),
 \end{aligned}$$

where

$$\hat{n}_{12} \stackrel{not.}{=} E[n_{12}|n_1, n_2, n_3, n_4, \psi^{(t)}] = \frac{\frac{1}{2}}{\frac{1}{2} + \frac{\psi}{4}} = \frac{\psi^{(t)}}{2 + \psi^{(t)}} \cdot n_1.$$

## M-step

$$\frac{\partial}{\partial \psi} E[\ln P(D|\psi)] = -(n_2 + n_3) \frac{1}{1 - \psi} + (\hat{n}_{12} + n_4)$$

$$\Rightarrow \frac{\partial^2}{\partial \psi^2} E[\ln P(D|\psi)] = (n_1 + n_3) \frac{-1}{n_1 + n_3} - (\hat{n}_{12} + n_4) \frac{1}{\psi^2} \leq 0$$

$\Rightarrow$  the auxiliary function is concave.

$$\frac{\partial}{\partial \psi} E[\ln P(D|\psi)] = 0 \Leftrightarrow$$

$$(\hat{n}_{12} + n_4)(1 - \psi) = \psi(n_2 + n_3) \Leftrightarrow$$

$$\psi(\hat{n}_{12} + n_2 + n_3 + n_4) = \hat{n}_{12} + n_4 \Leftrightarrow$$

$$\psi^{(t+1)} = \frac{\hat{n}_{12} + n_4}{n - \hat{n}_{11}}$$

where

$$\hat{n}_{11} \stackrel{not.}{=} E[n_{11}|n_1, n_2, n_3, n_4, \psi^{(t)}] = n_1 - \hat{n}_{12} = \frac{2}{2 + \psi^{(t)}} \cdot n_1.$$

A Matlab implementation of this EM algorithm produces  $\psi = 0.62682139$ .



**The EM algorithm for solving  
a Bernoulli mixture model**

CMU, 2008 fall, Eric Xing, HW4, pr. 1.4-7

Suppose I have two unfair coins. The first lands on heads with probability  $p$ , and the second lands on heads with probability  $q$ .

Imagine  $n$  tosses, where for each toss I choose to use the first coin with probability  $\pi$  and choose to use the second with probability  $1 - \pi$ . The outcome of each toss  $i$  is  $x_i \in \{0, 1\}$ .

Suppose I tell you the outcomes of the  $n$  tosses,  $x \stackrel{not.}{=} \{x_1, x_2, \dots, x_n\}$ , but I don't tell you which coins I used on which toss.

Given only the outcomes,  $x$ , your job is to compute estimates for  $\theta$  which is the set of all parameters,  $\theta = \{p, q, \pi\}$  using the EM algorithm.

To compute these estimates, we will create a latent variable  $Z$ , where  $z_i \in \{0, 1\}$  indicates the coin used for the  $n^{th}$  toss. For example  $z_2 = 1$  indicates the first coin was used on the second toss.

We define the “incomplete” data log-likelihood as  $\log P(x|\theta)$  and the “complete” data log-likelihood as  $\log P(x, z|\theta)$ .

- a. Show that  $E[z_i | x_i, \theta] = P(z_i = 1 | x_i, \theta)$ .
- b. Use Bayes rule to compute  $P(z_i = 1 | x_i, \theta^{(t)})$ , where  $\theta^{(t)}$  denotes the parameters at iteration  $t$ .
- c. Write down the complete log-likelihood,  $\log P(x, z | \theta)$ .
- d. ***E-Step:*** Show that the expected log-likelihood of the complete data  $Q(\theta | \theta^{(t)}) \stackrel{not.}{=} E_{P(z|x, \theta^{(t)})}[\log P(x, z | \theta)]$  is given by

$$\begin{aligned}
 Q(\theta | \theta^{(t)}) &= \sum_{i=1}^n E[z_i | x_i, \theta^{(t)}] \cdot (\log \pi + x_i \log p + (1 - x_i) \log(1 - p)) + \\
 &\quad + (1 - E[z_i | x_i, \theta^{(t)}]) \cdot (\log(1 - \pi) + x_i \log q + (1 - x_i) \log(1 - q))
 \end{aligned}$$

- e. ***M-Step:*** Describe the process you would use to obtain the update equations for  $p^{(t+1)}$ ,  $q^{(t+1)}$  si  $\pi^{(t+1)}$ .

## Solution

a.

$$\begin{aligned}
 E[z_i \mid x_i, \theta] &= \sum_{z \in \{0,1\}} z_i P(z_i \mid x_i, \theta) = 0 \cdot P(z_i = 0 \mid x_i, \theta) + 1 \cdot P(z_i = 1 \mid x_i, \theta) \\
 &\Rightarrow E[z_i \mid x_i, \theta] = P(z_i = 1 \mid x_i, \theta)
 \end{aligned}$$

b.

$$\begin{aligned}
 &P(z_i = 1 \mid x_i, \theta) \\
 &= \frac{P(x_i \mid z_i = 1, \theta) P(z_i = 1 \mid \theta)}{P(x_i \mid z_i = 1, \theta) P(z_i = 1 \mid \theta) + P(x_i \mid z_i = 0, \theta) P(z_i = 0 \mid \theta)} \\
 &= \frac{p^{x_i} \cdot (1 - p)^{1-x_i} \cdot \pi}{p^{x_i} \cdot (1 - p)^{1-x_i} \cdot \pi + q^{x_i} \cdot (1 - q)^{1-x_i} \cdot (1 - \pi)}
 \end{aligned}$$

c.

**Note (in Romanian):** Din punct de vedere metodologic, pentru a calcula  $P(x, z|\theta)$  ne putem inspira de la punctul precedent: observăm că la numitorul fracției care ne dă valoarea lui  $P(z_i|x_i, \theta)$  avem  $P(x_i, z_i = 1|\theta)$  și  $P(x_i, z_i = 0|\theta)$ . Pornind de la aceste două expresii, ne vom propune să le exprimăm în mod unitar, adică sub forma unei singure expresii.

$$\begin{aligned}
 \log P(x, z \mid \theta) &\stackrel{i.i.d.}{=} \log \prod_{i=1}^n P(x_i, z_i \mid \theta) = \log \prod_{i=1}^n P(x_i \mid z_i, \theta) \cdot P(z_i \mid \theta) \\
 &= \log \prod_{i=1}^n \left( p^{x_i} (1-p)^{1-x_i} \pi \right)^{z_i} \left( q^{x_i} (1-q)^{1-x_i} (1-\pi) \right)^{1-z_i} \\
 &= \sum_{i=1}^n \log \left( \left( p^{x_i} (1-p)^{1-x_i} \pi \right)^{z_i} \left( q^{x_i} (1-q)^{1-x_i} (1-\pi) \right)^{1-z_i} \right) \\
 &= \sum_{i=1}^n \left[ z_i \log \left( p^{x_i} (1-p)^{1-x_i} \pi \right) + (1-z_i) \log \left( q^{x_i} (1-q)^{1-x_i} (1-\pi) \right) \right]
 \end{aligned}$$

d.

$$\begin{aligned}
Q(\theta \mid \theta^{(t)}) &\stackrel{not.}{=} E_{P(z|x, \theta^{(t)})} [\log P(x, z \mid \theta)] \\
&= E_{P(z|x, \theta^{(t)})} \left[ \sum_{i=1}^n \left[ z_i \log (p^{x_i} (1-p)^{1-x_i} \pi) + \right. \right. \\
&\quad \left. \left. (1-z_i) \log (q^{x_i} (1-q)^{1-x_i} (1-\pi)) \right] \right] \\
&= \sum_{i=1}^n \left[ E[z_i \mid x_i, \theta^{(t)}] \cdot \log p^{x_i} (1-p)^{1-x_i} \pi + \right. \\
&\quad \left. + (1 - E[z_i \mid x_i, \theta^{(t)}]) \cdot \log q^{x_i} (1-q)^{1-x_i} (1-\pi) \right] \\
&= \sum_{i=1}^n \left[ E[z_i \mid x_i, \theta^{(t)}] \cdot (\log \pi + x_i \log p + (1-x_i) \log(1-p)) + \right. \\
&\quad \left. + (1 - E[z_i \mid x_i, \theta^{(t)}]) \cdot (\log(1-\pi) + x_i \log q + (1-x_i) \log(1-q)) \right]
\end{aligned}$$

e.

$$\mu_i^{(t)} \stackrel{\text{not.}}{=} E[z_i \mid x_i, \theta^{(t)}] = \frac{(p^{(t)})^{x_i} \cdot (1 - p^{(t)})^{1-x_i} \cdot \pi^{(t)}}{(p^{(t)})^{x_i} \cdot (1 - p^{(t)})^{1-x_i} \cdot \pi^{(t)} + (q^{(t)})^{x_i} \cdot (1 - q^{(t)})^{1-x_i} \cdot (1 - \pi^{(t)})}$$

$$\Rightarrow Q(\theta \mid \theta^{(t)}) = \sum_{i=1}^n \left[ \mu_i^{(t)} (\log \pi + x_i \log p + (1 - x_i) \log(1 - p)) + \right. \\ \left. + (1 - \mu_i^{(t)}) \cdot (\log(1 - \pi) + x_i \log q + (1 - x_i) \log(1 - q)) \right]$$

$$\begin{aligned} \frac{\partial Q(\theta \mid \theta^{(t)})}{\partial p} = 0 &\Leftrightarrow \sum_{i=1}^n \mu_i^{(t)} \left( \frac{x_i}{p} - \frac{1 - x_i}{1 - p} \right) = 0 \Leftrightarrow \frac{1}{p} \sum_{i=1}^n \mu_i^{(t)} x_i = \frac{1}{1 - p} \sum_{i=1}^n \mu_i^{(t)} (1 - x_i) \\ &\Leftrightarrow (1 - p) \sum_{i=1}^n \mu_i^{(t)} x_i = p \sum_{i=1}^n \mu_i^{(t)} (1 - x_i) \Leftrightarrow \sum_{i=1}^n \mu_i^{(t)} x_i = p \left( \sum_{i=1}^n \mu_i^{(t)} (1 - x_i) + \sum_{i=1}^n \mu_i^{(t)} x_i \right) \\ &\Leftrightarrow \sum_{i=1}^n \mu_i^{(t)} x_i = p \sum_{i=1}^n \mu_i^{(t)} \Rightarrow p^{(t+1)} = \frac{\sum_{i=1}^n \mu_i^{(t)} x_i}{\sum_{i=1}^n \mu_i^{(t)}} \in [0, 1] \end{aligned}$$

$$\frac{\partial Q(\theta \mid \theta^{(t)})}{\partial q} = 0 \Leftrightarrow \sum_{i=1}^n (1 - \mu_i^{(t)}) \left( \frac{x_i}{q} - \frac{1 - x_i}{1 - q} \right) = 0 \Rightarrow q^{(t+1)} = \frac{\sum_{i=1}^n (1 - \mu_i^{(t)}) x_i}{\sum_{i=1}^n (1 - \mu_i^{(t)})} \in [0, 1]$$

$$\frac{\partial Q(\theta \mid \theta^{(t)})}{\partial \pi} = 0 \Leftrightarrow \sum_{i=1}^n \left( \frac{\mu_i^{(t)}}{\pi} - \frac{1 - \mu_i^{(t)}}{1 - \pi} \right) = 0 \Leftrightarrow \frac{1}{\pi} \sum_{i=1}^n \mu_i^{(t)} = \frac{1}{1 - \pi} \sum_{i=1}^n (1 - \mu_i^{(t)})$$

$$\Leftrightarrow (1 - \pi) \sum_{i=1}^n \mu_i^{(t)} = \pi \sum_{i=1}^n (1 - \mu_i^{(t)}) \Leftrightarrow \sum_{i=1}^n \mu_i^{(t)} = \pi \left( \sum_{i=1}^n (1 - \mu_i^{(t)}) + \sum_{i=1}^n \mu_i^{(t)} \right)$$

$$\Leftrightarrow \sum_{i=1}^n \mu_i^{(t)} = \pi \sum_{i=1}^n 1 \Leftrightarrow \sum_{i=1}^n \mu_i^{(t)} = n\pi \Rightarrow \pi^{(t+1)} = \frac{1}{n} \sum_{i=1}^n \mu_i^{(t)} \in [0, 1]$$

**Note:** It can be easily shown that *the second derivatives*  $\left( \frac{\partial^2}{\partial p^2}, \frac{\partial^2}{\partial q^2} \text{ and } \frac{\partial^2}{\partial \pi^2} \right)$  are all *negative* on the respective domains of definition. Therefore, the above solutions  $(p^{(t+1)}, q^{(t+1)} \text{ and } \pi^{(t+1)})$  maximize the  $Q$  functional.



## Estimarea parametrilor unei mixturi de distribuții [de tip] Bernoulli

- A. când toate variabilele sunt observabile: MLE;
- B. când unele variabile sunt neobservabile: algoritmul EM

prelucrare de Liviu Ciortuz, după

*“What is the expectation maximization algorithm?”*,  
Chuong B. Do, Serafim Batzoglou,  
Nature Biotechnology, vol. 26, no. 8, 2008, pag. 897-899

Fie următorul *experiment probabilist*:

Disponem de două monede,  $A$  și  $B$ .

Efectuăm 5 serii de operațiuni de tipul următor:

- Alegem în mod aleatoriu una dintre monedele  $A$  și  $B$ , cu probabilitate egală ( $1/2$ );
- Aruncăm de 10 ori moneda care tocmai a fost aleasă ( $Z$ ) și notăm rezultatul, sumarizat ca număr ( $X$ ) de fețe ‘head’ (ro. ‘stemă’) obținute în urma aruncării.

A. La acest punct vom considera că s-a obținut următorul rezultat pentru experimentul nostru:

$i$	$Z_i$	$X_i$
1	$B$	$5H \ (5T)$
2	$A$	$9H \ (1T)$
3	$A$	$8H \ (2T)$
4	$B$	$4H \ (6T)$
5	$A$	$7H \ (3T)$

Semnificația variabilelor aleatoare  $Z_i$  și  $X_i$  pentru  $i = 1, \dots, 5$  din tabelul de mai sus este imediată.

i. Calculați  $\hat{\theta}_A$  și  $\hat{\theta}_B$ , probabilitățile de apariție a feței ‘head’/stemă pentru cele două monede, folosind *definiția clasică pentru probabilitatea evenimentelor aleatoare*, și anume raportul dintre numărul de cazuri favorabile și numărul de cazuri posibile, relativ la întregul experiment.

Răspuns:

Analizând datele din tabelul din enunț, rezultă imediat

$$\hat{\theta}_A = \frac{24}{24 + 6} = 0.8 \text{ și } \hat{\theta}_B = \frac{9}{9 + 11} = 0.45.$$

De observat că termenii 6 și respectiv 11 de la numitorii acestor fracții reprezintă numărul de fețe ‘tail’/ban care au fost obținute la aruncarea monedei  $A$  și respectiv  $B$ :  $6T = 1T + 2T + 3T$ ,  $11T = 5T + 6T$ .

## Maximum likelihood



Coin A	Coin B
	5 H, 5 T
9 H, 1 T	
8 H, 2 T	
	4 H, 6 T
7 H, 3 T	
24 H, 6 T	9 H, 11 T

$$\hat{\theta}_A = \frac{24}{24 + 6} = 0.80$$

$$\hat{\theta}_B = \frac{9}{9 + 11} = 0.45$$

## Observație

Dacă în locul variabilelor binare  $Z_i \in \{A, B\}$  pentru  $i = 1, \dots, 5$  introducem în mod natural variabilele-indicator  $Z_{i,A} \in \{0, 1\}$  și  $Z_{i,B} \in \{0, 1\}$  tot pentru  $i = 1, \dots, 5$ , definite prin  $Z_{i,A} = 1$  iff  $Z_i = A$ , și  $Z_{i,A} = 0$  iff  $Z_{i,B} = 0$ , atunci procesările necesare pentru calculul probabilităților/parametrilor  $\hat{\theta}_A$  și  $\hat{\theta}_B$  pot fi prezentate în mod sintetizat ca în tabelul de mai jos.<sup>a</sup>

$i$	$Z_{i,A}$	$Z_{i,B}$	$X_i$
1	0	1	5H
2	1	0	9H
3	1	0	8H
4	0	1	4H
5	1	0	7H

$\Rightarrow$

---

<sup>a</sup>Prezentăm acest “artificiu” ca pregătire pentru rezolvarea (ulterioară a) punctului B al prezentei probleme.

$$\Rightarrow$$

$X_i \cdot Z_{i,A}$	$X_i \cdot Z_{i,B}$
$0H \ (0T)$	$5H \ (5T)$
$9H \ (1T)$	$0H \ (0T)$
$8H \ (2T)$	$0H \ (0T)$
$0H \ (0T)$	$4H \ (6T)$
$7H \ (3T)$	$0H \ (0T)$
$\sum_{i=1}^5 X_i \cdot Z_{i,A} = 24H$	$\sum_{i=1}^5 X_i \cdot Z_{i,B} = 9H$
$\sum_{i=1}^5 (10 - X_i) \cdot Z_{i,A} = 6T$	$\sum_{i=1}^5 (10 - X_i) \cdot Z_{i,B} = 11T$

$$\Rightarrow$$

$$\Rightarrow \begin{cases} \hat{\theta}_A = \frac{24}{24 + 6} = 0.8 \\ \hat{\theta}_B = \frac{9}{9 + 11} = 0.45 \end{cases}$$

*ii.* Calculați  $L_1(\theta_A, \theta_B) \stackrel{not.}{=} P(X, Z \mid \theta_A, \theta_B)$ , funcția de verosimilitate a datelor  $X \stackrel{not.}{=} \langle X_1, \dots, X_5 \rangle$  și  $Z \stackrel{not.}{=} \langle Z_1, \dots, Z_5 \rangle$ , în raport cu parametrii  $\theta_A$  și respectiv  $\theta_B$  ai distribuțiilor *Bernoulli* care modelează aruncarea celor două monede.

*Observație:* Facem *presupunerea* că am reținut succesiunea [tuturor] rezultatelor obținute la aruncările celor două monede. Precizarea *de facto* a acestei succesiuni nu este esențială. Dacă nu s-ar face această *presupunere*, ar trebui să lucrăm cu *distribuția binomială*.



Răspuns:

Calculul verosimilității datelor:

$$\begin{aligned}
 L_1(\theta_A, \theta_B) &\stackrel{\text{def.}}{=} P(X, Z_A, Z_B \mid \theta_A, \theta_B) \stackrel{\text{indep. cdt.}}{=} \prod_{i=1}^5 P(X_i, Z_{i,A}, Z_{i,B} \mid \theta_A, \theta_B) \\
 &= \prod_{i=1}^5 P(X_i \mid Z_{i,A}, Z_{i,B}; \theta_A, \theta_B) \cdot P(Z_{i,A}, Z_{i,B} \mid \theta_A, \theta_B) \\
 &= P(X_1 \mid Z_{B,1} = 1, \theta_B) \cdot 1/2 \cdot \\
 &\quad P(X_2 \mid Z_{A,2} = 1, \theta_A) \cdot 1/2 \cdot \\
 &\quad P(X_3 \mid Z_{A,3} = 1, \theta_A) \cdot 1/2 \cdot \\
 &\quad P(X_4 \mid Z_{B,4} = 1, \theta_B) \cdot 1/2 \cdot \\
 &\quad P(X_5 \mid Z_{A,5} = 1, \theta_A) \cdot 1/2 \\
 &= \theta_B^5 (1 - \theta_B)^5 \cdot \theta_A^9 (1 - \theta_A) \cdot \theta_A^8 (1 - \theta_A)^2 \cdot \theta_B^4 (1 - \theta_B)^6 \cdot \theta_A^7 (1 - \theta_A)^3 \cdot \frac{1}{2^5} \\
 &= \frac{1}{2^5} \theta_A^{24} (1 - \theta_A)^6 \theta_B^9 (1 - \theta_B)^{11}
 \end{aligned}$$

*iii.* Calculați  $\hat{\theta}_A \stackrel{not.}{=} \arg \max_{\theta_A} \log L_1(\theta_A, \theta_B)$  și  $\hat{\theta}_B \stackrel{not.}{=} \arg \max_{\theta_B} \log L_1(\theta_A, \theta_B)$  folosind derivatele parțiale de ordinul întâi.

*Observații:*

1. Baza logaritmului, fixată dar lăsată mai sus nespecificată, se va considera supraunitară (de exemplu 2,  $e$  sau 10).
2. Lucrând corect, veți obține același rezultat ca la punctul  $i$ .

Răspuns:

Funcția de log-verosimilitate a datelor complete se exprimă astfel:

$$\log L_1(\theta_A, \theta_B) = -5 \log 2 + 24 \log \theta_A + 6 \log(1 - \theta_A) + 9 \log \theta_B + 11 \log(1 - \theta_B)$$

Prin urmare, maximul acestei funcții în raport cu parametrul  $\theta_A$  se calculează astfel:

$$\frac{\partial \log L_1(\theta_A, \theta_B)}{\partial \theta_A} = 0 \Leftrightarrow \frac{24}{\theta_A} - \frac{6}{1 - \theta_A} = 0 \Leftrightarrow \frac{4}{\theta_A} = \frac{1}{1 - \theta_A} \Leftrightarrow 4 - 4\theta_A = \theta_A \Leftrightarrow \hat{\theta}_A = 0.8$$

Similar, se face calculul și pentru  $\frac{\partial \log L_1(\theta_A, \theta_B)}{\partial \theta_B}$  și se obține  $\hat{\theta}_B = 0.45$ .<sup>a</sup>

Cele două valori obținute,  $\hat{\theta}_A$  și  $\hat{\theta}_B$ , reprezintă estimarea de verosimilitate maximă (MLE) a probabilităților de apariție a feței ‘head’ (‘stema’) pentru moneda  $A$  și respectiv moneda  $B$ .

---

<sup>a</sup>Se verifică ușor faptul că într-adevăr rădăcinile derivatelor parțiale de ordinul întâi pentru funcția de log-verosimilitate reprezintă puncte de maxim. Pentru aceasta se studiază semnele acestor derivate.

## Observații

1. Am arătat pe acest caz particular că metoda de calculare a probabilităților ( $\hat{\theta}_A$  și  $\hat{\theta}_B$ ) direct din datele observate (așa cum o știm din liceu) corespunde de fapt metodei de estimare în sensul verosimității maxime (MLE).
2. La punctul B vom arăta cum anume se poate face estimarea aceluiași parametri  $\theta_A$  și  $\theta_B$  în cazul în care o parte din date, și anume variabilele  $Z_i$  (pentru  $i = 1, \dots, 5$ ) sunt neobservabile.

B. La acest punct se va relua experimentul de la punctul A, însă de data aceasta vom considera că valorile variabilelor  $Z_i$  nu sunt cunoscute.

$i$	$Z_i$	$X_i$
1	?	$5H \ (5T)$
2	?	$9H \ (1T)$
3	?	$8H \ (2T)$
4	?	$4H \ (6T)$
5	?	$7H \ (3T)$

*iv.* Pentru conveniență, în locul variabilelor „neobservabile“  $Z_i$  pentru  $i = 1, \dots, 5$  vom considera variabilele-indicator (de asemenea neobservabile)  $Z_{i,A}, Z_{i,B} \in \{0, 1\}$ , cu  $Z_{i,A} = 1$  iff  $Z_{i,B} = 0$  și  $Z_{i,B} = 1$  iff  $Z_{i,A} = 0$ . Evident, întrucât variabilele  $Z_i$  sunt aleatoare, rezultă că și variabilele  $Z_{i,A}$  și  $Z_{i,B}$  sunt aleatoare.

Folosind teorema lui Bayes, calculați mediile variabilelor neobservabile  $Z_{i,A}$  și  $Z_{i,B}$  condiționate de variabilele observabile  $X_i$ . Veți considera că parametrii acestor distribuții Bernoulli care modelează aruncarea monedelor  $A$  și  $B$  au valorile  $\theta_A^{(0)} = 0.6$  și respectiv  $\theta_B^{(0)} = 0.5$ .

Așadar, se cer:  $E[Z_{i,A} \mid X_i, \theta_A^{(0)}, \theta_B^{(0)}]$  și  $E[Z_{i,B} \mid X_i, \theta_A^{(0)}, \theta_B^{(0)}]$  pentru  $i = 1, \dots, 5$ . Ca și mai înainte, probabilitățile a priori  $P(Z_{i,A} = 1)$  și  $P(Z_{i,B} = 1)$  se vor considera egale cu  $1/2$ .

## Notă

Algoritmul EM ne permite să facem în mod iterativ estimarea parametrilor  $\theta_A$  și  $\theta_B$  în funcție de valorile variabilelor observabile,  $X_i$ , și de valorile inițiale atribuite parametrilor (în cazul nostru,  $\theta_A^{(0)} = 0.6$  și  $\theta_B^{(0)} = 0.5$ ).

Vom face o sinteză a calculelor de la prima iterație a algoritmului EM — detaliate la punctele *iv*, *v* și *vi* de mai jos — sub forma următoare, care seamănă într-o anumită măsură cu tabelele de la punctul A:

$i$	$Z_{i,A}$	$Z_{i,B}$	$X_i$		$E[Z_{i,A}]$	$E[Z_{i,B}]$	
1	—	—	$5H$		$0.45H$	$0.55H$	
2	—	—	$9H$	$\xRightarrow{E}$	$0.80H$	$0.20H$	$\xRightarrow{M}$
3	—	—	$8H$		$0.73H$	$0.27H$	
4	—	—	$4H$		$0.35H$	$0.65H$	
5	—	—	$7H$		$0.65H$	$0.35H$	

## Notă (cont.)

$$\begin{array}{c}
 \xRightarrow{M}
 \end{array}
 \begin{array}{c|c}
 X_i \cdot E[Z_{i,A}] & X_i \cdot E[Z_{i,B}] \\
 \hline
 \begin{array}{c}
 2.2H \ (2.2T) \\
 7.2H \ (0.8T) \\
 5.9H \ (1.5T) \\
 1.4H \ (2.1T) \\
 4.5H \ (1.9T)
 \end{array}
 &
 \begin{array}{c}
 2.8H \ (2.8T) \\
 1.8H \ (0.2T) \\
 2.1H \ (0.5T) \\
 2.6H \ (3.9T) \\
 2.5H \ (1.1T)
 \end{array}
 \\
 \hline
 \begin{array}{c}
 \sum_{i=1}^5 X_i \cdot E[Z_{i,A}] = 21.3H \\
 \sum_{i=1}^5 (10 - X_i) \cdot E[Z_{i,A}] = 8.7T
 \end{array}
 &
 \begin{array}{c}
 \sum_{i=1}^5 X_i \cdot E[Z_{i,B}] = 11.7H \\
 \sum_{i=1}^5 (10 - X_i) \cdot E[Z_{i,B}] = 8.3T
 \end{array}
 \end{array}
 \Rightarrow
 \begin{array}{c}
 \Rightarrow \left\{ \begin{array}{l} \hat{\theta}_A^{(1)} = \frac{21.3}{21.3 + 8.7} \approx 0.71 \\ \hat{\theta}_B^{(1)} = \frac{11.7}{11.7 + 8.3} \approx 0.58 \end{array} \right.
 \end{array}$$



## Notă (cont.)

Vom arăta că:

- într-adevăr, este posibilă calcularea mediilor variabilelor neobservabile  $Z_{i,A}$  și  $Z_{i,B}$ , condiționate de variabilele observabile  $X_i$  și în funcție de valorile asignate inițial (0.6 și 0.5 în enunț, dar în general ele se pot asigna în mod aleatoriu) pentru parametrii  $\theta_A$  și  $\theta_B$ ;
- față de tabloul de sinteză de la punctul precedent, când toate variabilele erau observabile și se calculau produsele  $X_i \cdot Z_{i,A}$  și  $X_i \cdot Z_{i,B}$ , aici se înlocuiesc variabilele  $Z_{i,A}$  și  $Z_{i,B}$  cu mediile  $E[Z_{i,A}]$  și  $E[Z_{i,B}]$  în produsele respective. De fapt, în loc să se calculeze  $\sum_i X_i \cdot Z_{i,A}$  se calculează media  $E[\sum_i X_i \cdot Z_{i,A}]$ , și similar pentru  $B$ .

**Notăție:** Pentru simplitate, în cele de mai sus (inclusiv în tabelele precedente), prin  $E[Z_{i,A}]$  am notat  $E[Z_{i,A} \mid X_i, \theta^{(0)}]$ , iar prin  $E[Z_{i,B}]$  am notat  $E[Z_{i,B} \mid X_i, \theta^{(0)}]$ , unde  $\theta^{(0)} \stackrel{not.}{=} (\theta_A^{(0)}, \theta_B^{(0)})$ .

## Răspuns (iv)

Întrucât variabilele  $Z_{i,A}$  au valori boolene (0 sau 1), rezultă că

$$\begin{aligned} E[Z_{i,A} | X_i, \theta^{(0)}] &= 0 \cdot P(Z_{i,A} = 0 | X_i, \theta^{(0)}) + 1 \cdot P(Z_{i,A} = 1 | X_i, \theta^{(0)}) \\ &= P(Z_{i,A} = 1 | X_i, \theta^{(0)}) \end{aligned}$$

Probabilitățile  $P(Z_{i,A} = 1 | X, \theta^{(0)}) = P(Z_{i,A} = 1 | X_i, \theta^{(0)})$ , pentru  $i = 1, \dots, 5$ , se pot calcula folosind teorema lui Bayes:

$$\begin{aligned} P(Z_{i,A} = 1 | X_i, \theta^{(0)}) &= \frac{P(X_i | Z_{i,A} = 1, \theta^{(0)}) \cdot P(Z_{i,A} = 1 | \theta^{(0)})}{\sum_{j \in \{0,1\}} P(X_i | Z_{i,A} = j, \theta^{(0)}) \cdot P(Z_{i,A} = j | \theta^{(0)})} \\ &= \frac{P(X_i | Z_{i,A} = 1, \theta_A^{(0)})}{P(X_i | Z_{i,A} = 1, \theta_A^{(0)}) + P(X_i | Z_{i,B} = 1, \theta_B^{(0)})} \end{aligned}$$

S-a ținut cont că  $P(Z_{i,A} = 1 | \theta^{(0)}) = P(Z_{i,B} = 1 | \theta^{(0)}) = 1/2$  (a se vedea enunțul).

De exemplu, pentru  $i = 1$  vom avea:

$$E[Z_{A,1} \mid X_1, \theta^{(0)}] = \frac{0.6^5(1 - 0.6)^5}{0.6^5(1 - 0.6)^5 + 0.5^5(1 - 0.5)^5} = \frac{1}{1 + \left(\frac{0.25}{0.24}\right)^5} \approx 0.45$$

Similar cu  $E[Z_{A,1} \mid X_1, \theta^{(0)}]$  se calculează și celelalte medii  $E[Z_{i,A} \mid X_i, \theta^{(0)}]$  pentru  $i = 2, \dots, 5$  și  $E[Z_{i,B} \mid X_i, \theta^{(0)}]$  pentru  $i = 1, \dots, 5$ .

Am înregistrat aceste valori/medii în cel de-al doilea tabel din *Nota* de mai sus.

**Observație:** Se poate ține cont că, de îndată ce s-a calculat  $E[Z_{i,A} \mid X_i, \theta^{(0)}]$ , se poate obține imediat și  $E[Z_{i,B} \mid X_i, \theta^{(0)}] = 1 - E[Z_{i,A} \mid X_i, \theta^{(0)}]$ , fiindcă  $Z_{i,A} + Z_{i,B} = 1$ .

v. Calculați media funcției de log-verosimilitate a datelor complete,  $X$  (observabile) și  $Z$  (neobservabile):

$$L_2(\theta_A, \theta_B) \stackrel{def.}{=} E_{P(Z|X, \theta^{(0)})}[\log P(X, Z | \theta)],$$

unde  $\theta \stackrel{not.}{=} (\theta_A, \theta_B)$  și  $\theta^{(0)} \stackrel{not.}{=} (\theta_A^{(0)}, \theta_B^{(0)})$ .

*Semnificația* notației de mai sus este următoarea:

funcția  $L_2(\theta_A, \theta_B)$  este o medie a variabilei aleatoare reprezentată de log-verosimilitatea datelor complete (observabile și, respectiv, neobservabile), iar această medie se calculează în raport cu distribuția probabilistă condițională a datelor neobservabile,  $P(Z | X, \theta^{(0)})$ .

*Observație:* La elaborarea calculului, veți folosi mai întâi proprietatea de liniaritate a mediilor variabilelor aleatoare, și apoi rezultatele de la punctul *iv*.

Răspuns:

Media funcției de log-verosimilitate a datelor complete,  $L_2(\theta_A, \theta_B)$ , se calculează astfel:

$$\begin{aligned}
 L_2(\theta_A, \theta_B) &\stackrel{def.}{=} E_{P(Z|X, \theta^{(0)})} [\log P(X, Z | \theta)] \\
 &\stackrel{indep. \ cdt.}{=} E_{P(Z|X, \theta^{(0)})} \left[ \log \prod_{i=1}^5 P(X_i, Z_{i,A}, Z_{i,B} | \theta_A, \theta_B) \right] \\
 &\stackrel{reg. \ de \ mult.}{=} E_{P(Z|X, \theta^{(0)})} \left[ \log \prod_{i=1}^5 P(X_i | Z_{i,A}, Z_{i,B}; \theta_A, \theta_B) \cdot P(Z_{i,A}, Z_{i,B} | \theta_A, \theta_B) \right]
 \end{aligned}$$

În continuare, omițând din nou distribuția probabilistă în raport cu care se calculează media aceasta întrucât ea poate fi subînțeleasă, vom scrie:

$$L_2(\theta_A, \theta_B)$$

$$= E\left[\log \prod_{i=1}^5 \cdot (\theta_A^{Z_{i,A}})^{X_i} \cdot [(1 - \theta_A)^{Z_{i,A}}]^{10-X_i} \cdot (\theta_B^{Z_{i,B}})^{X_i} \cdot [(1 - \theta_B)^{Z_{i,B}}]^{10-X_i} \cdot \frac{1}{2} \right]$$

$$= E\left[\sum_{i=1}^5 [X_i \cdot Z_{i,A} \cdot \log \theta_A + (10 - X_i) \cdot Z_{i,A} \cdot \log(1 - \theta_A) + \right. \\ \left. X_i \cdot Z_{i,B} \cdot \log \theta_B + (10 - X_i) \cdot Z_{i,B} \cdot \log(1 - \theta_B) - \log 2] \right]$$

$$\stackrel{lin. med.}{=} \sum_{i=1}^5 [X_i \cdot E[Z_{i,A}] \cdot \log \theta_A + (10 - X_i) \cdot E[Z_{i,A}] \cdot \log(1 - \theta_A) + \\ X_i \cdot E[Z_{i,B}] \cdot \log \theta_B + (10 - X_i) \cdot E[Z_{i,B}] \cdot \log(1 - \theta_B) - \log 2]$$

$$= \sum_{i=1}^5 \log[\theta_A^{X_i \cdot E[Z_{i,A}]} \cdot (1 - \theta_A)^{(10-X_i) \cdot E[Z_{i,A}]} \cdot \theta_B^{X_i \cdot E[Z_{i,B}]} \cdot (1 - \theta_B)^{(10-X_i) \cdot E[Z_{i,B}]} \cdot \frac{1}{2}]$$

$$= \log(\theta_A^{2.2} \cdot (1 - \theta_A)^{2.2} \cdot \theta_B^{2.8} \cdot (1 - \theta_B)^{2.8} \cdot \dots \cdot \theta_A^{4.5} \cdot (1 - \theta_A)^{1.9} \cdot \theta_B^{2.5} \cdot (1 - \theta_B)^{1.1} \cdot \frac{1}{2}).$$

La ultima egalitate de mai sus, cantitățile fracționare provin din calculele simple  $X_1 \cdot E[Z_{1,A} \mid X_1, \theta] \approx 2.2$ ,  $X_1 \cdot E[Z_{1,B} \mid X_1, \theta] \approx 2.8$ , ...,  $X_5 \cdot E[Z_{1,A} \mid X_5, \theta] \approx 4.5$ ,  $X_5 \cdot E[Z_{1,B} \mid X_5, \theta] \approx 2.5$  (a se vedea tabelele din cadrul *Notei* precedente).

**vi. Calculați**  $\theta_A^{(1)} \stackrel{\text{not.}}{=} \arg \max_{\theta_A} L_2(\theta_A, \theta_B)$  **și**  $\theta_B^{(1)} \stackrel{\text{not.}}{=} \arg \max_{\theta_B} L_2(\theta_A, \theta_B)$ .

Răspuns:

Valorile parametrilor  $\theta_A$  și  $\theta_B$  pentru care se atinge maximul mediei funcției de log-verosimilitate a datelor complete se obțin cu ajutorul derivatelor parțiale de ordinul întâi:<sup>a</sup>

$$\begin{aligned} \frac{\partial L_2(\theta_A, \theta_B)}{\partial \theta_A} &= 0 \\ \Rightarrow \frac{\partial}{\partial \theta_A} (2.2 \log \theta_A + 2.2 \log(1 - \theta_A) + \dots + 4.5 \log \theta_A + 1.9 \log(1 - \theta_A)) &= 0 \\ \Rightarrow \frac{2.2}{\theta_A} - \frac{2.2}{1 - \theta_A} + \dots + \frac{4.5}{\theta_A} - \frac{1.9}{1 - \theta_A} &= 0 \Rightarrow \dots \Rightarrow \theta_A^{(1)} \approx 0.71. \end{aligned}$$

**Similar, vom obține**  $\theta_B^{(1)} \approx 0.58$ .

---

<sup>a</sup>Este imediat că derivatele de ordinul al doilea au valori negative pe tot domeniul de definiție.



C. Formalizați pașii E și M ai algoritmului EM pentru estimarea parametrilor  $\theta_A$  și  $\theta_B$  în condițiile de la punctul B.

Răspuns:

Formulele care se folosesc în cadrul algoritmului EM pentru rezolvarea problemei date — i.e. estimarea parametrilor  $\theta_A$  și  $\theta_B$  atunci când variabilele  $Z_i$  sunt neobservabile —, se elaborează/deduc astfel:

***Pasul E:***

$$\begin{aligned}
 E[Z_{i,A} \mid X, \theta] &= P(Z_{i,A} = 1 \mid X, \theta) = P(Z_{i,A} = 1 \mid X_i, \theta) \\
 &\stackrel{T. B.}{=} \frac{P(X_i \mid Z_{i,A} = 1; \theta) \cdot \overbrace{P(Z_{i,A} = 1 \mid \theta)}^{1/2}}{P(X_i \mid Z_{i,A} = 1; \theta) \cdot P(Z_{i,A} = 1 \mid \theta) + P(X_i \mid Z_{i,B} = 1; \theta) \cdot \underbrace{P(Z_{i,B} = 1 \mid \theta)}_{1/2}} \\
 &= \frac{P(X_i \mid Z_{i,A} = 1; \theta)}{P(X_i \mid Z_{i,A} = 1; \theta) + P(X_i \mid Z_{i,B} = 1; \theta)} \\
 &= \frac{\theta_A^{X_i} (1 - \theta_A)^{10-X_i}}{\theta_A^{X_i} (1 - \theta_A)^{10-X_i} + \theta_B^{X_i} (1 - \theta_B)^{10-X_i}}
 \end{aligned}$$

**Similar, vom obține:**

$$E[Z_{i,B} \mid X, \theta] = \frac{\theta_B^{X_i} (1 - \theta_B)^{10-X_i}}{\theta_A^{X_i} (1 - \theta_A)^{10-X_i} + \theta_B^{X_i} (1 - \theta_B)^{10-X_i}}$$

Notând cu

- $x_i$  valoarea variabilei  $X_i$ ,
- $\theta_A^{(t)}$  și respectiv  $\theta_B^{(t)}$  estimările parametrilor  $\theta_A$  și  $\theta_B$  la iterația  $t$  a algoritmului EM,
- $p_{i,A}^{(t+1)}$  și respectiv  $p_{i,B}^{(t+1)}$ , mediile  $E[Z_{i,A} \mid X_i, \theta_A^{(t)}]$  și  $E[Z_{i,B} \mid X_i, \theta_B^{(t)}]$ ,

vom avea:

$$p_{i,A}^{(t+1)} = \frac{(\theta_A^{(t)})^{x_i} (10 - \theta_A^{(t)})^{10-x_i}}{(\theta_A^{(t)})^{x_i} (10 - \theta_A^{(t)})^{10-x_i} + (\theta_B^{(t)})^{x_i} (10 - \theta_B^{(t)})^{10-x_i}}$$

$$p_{i,B}^{(t+1)} = \frac{(\theta_B^{(t)})^{x_i} (10 - \theta_B^{(t)})^{10-x_i}}{(\theta_A^{(t)})^{x_i} (10 - \theta_A^{(t)})^{10-x_i} + (\theta_B^{(t)})^{x_i} (10 - \theta_B^{(t)})^{10-x_i}}$$

*Pasul M:*

Ca și mai înainte, în formulele de mai jos vom folosi notațiile simplificate  $E[Z_{i,A}] \stackrel{not.}{=} E[Z_{i,A} | X_i, \theta^{(t)}]$  și  $E[Z_{i,B}] \stackrel{not.}{=} E[Z_{i,B} | X_i, \theta^{(t)}]$ .

Cu aceste notații, procedând similar cu calculul de la partea B, punctul  $v$ , vom avea:

$$L_2(\theta_A, \theta_B) = \log \prod_{i=1}^5 \theta_A^{x_i E[Z_{i,A}]} (1 - \theta_A)^{(10-x_i) E[Z_{i,A}]} \theta_B^{x_i E[Z_{i,B}]} (1 - \theta_B)^{(10-x_i) E[Z_{i,B}]}$$

Prin urmare,

$$\begin{aligned}
 \frac{\partial}{\partial \theta_A} L_2(\theta_A, \theta_B) = 0 &\Rightarrow \frac{1}{\theta_A} \sum_{i=1}^5 x_i E[Z_{i,A}] = \frac{1}{1 - \theta_A} \sum_{i=1}^5 (10 - x_i) E[Z_{i,A}] \\
 &\Rightarrow (1 - \theta_A) \sum_{i=1}^5 x_i E[Z_{i,A}] = \theta_A \sum_{i=1}^5 (10 - x_i) E[Z_{i,A}] \\
 &\Rightarrow \sum_{i=1}^5 x_i E[Z_{i,A}] = 10 \theta_A \sum_{i=1}^5 E[Z_{i,A}] \\
 &\Rightarrow \theta_A = \frac{\sum_{i=1}^5 x_i E[Z_{i,A}]}{10 \sum_{i=1}^5 E[Z_{i,A}]} \quad \text{și, similar, } \theta_B = \frac{\sum_{i=1}^5 x_i E[Z_{i,B}]}{10 \sum_{i=1}^5 E[Z_{i,B}]}
 \end{aligned}$$

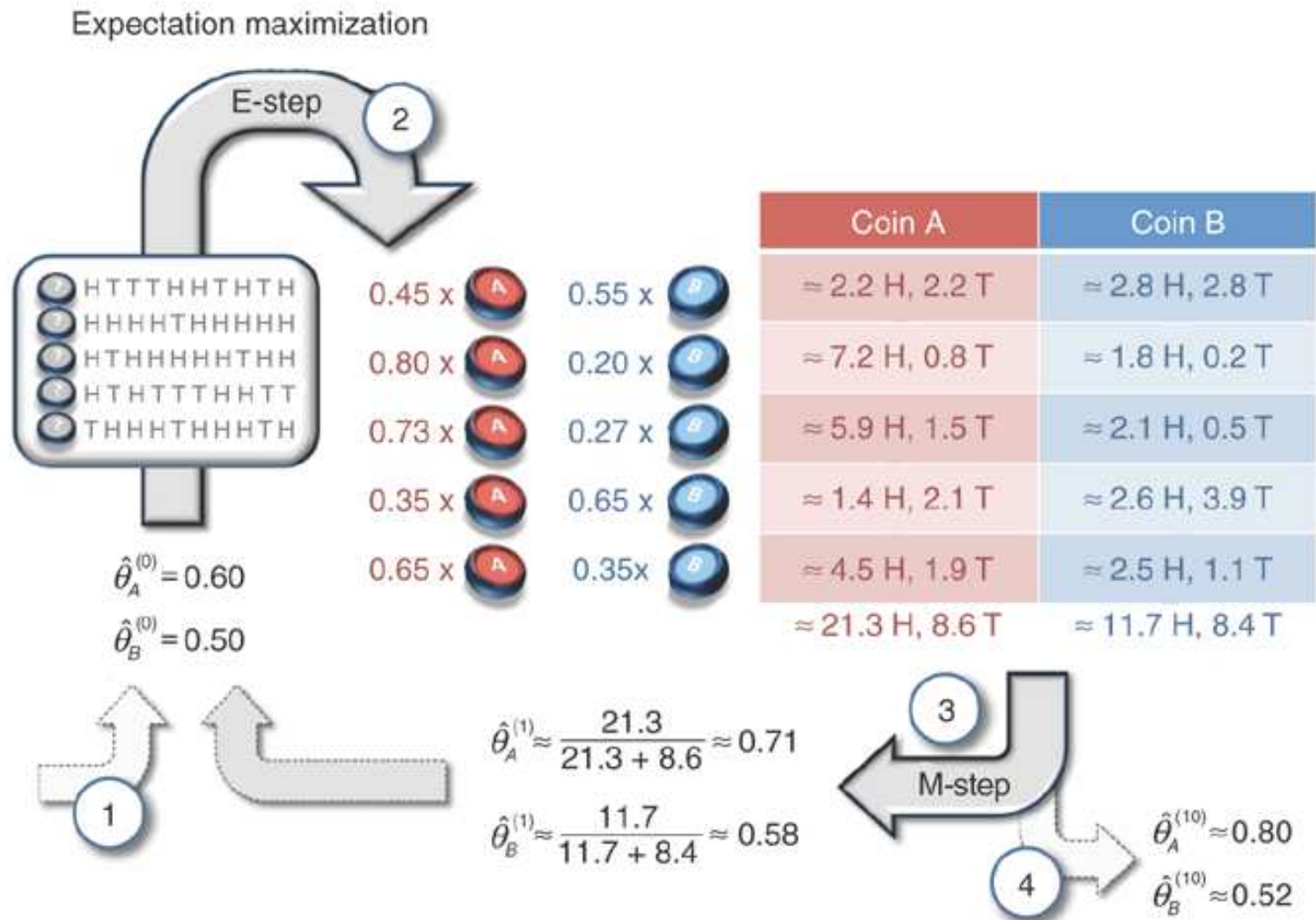
Așadar, la pasul M al algoritmului EM vom avea:

$$\theta_A^{(t+1)} = \frac{\sum_{i=1}^5 x_i p_{i,A}^{(t+1)}}{10 \sum_{i=1}^5 p_{i,A}^{(t+1)}} \quad \text{și} \quad \theta_B^{(t+1)} = \frac{\sum_{i=1}^5 x_i p_{i,B}^{(t+1)}}{10 \sum_{i=1}^5 p_{i,B}^{(t+1)}}$$

## Observație

Implementând algoritmul EM cu relațiile obținute pentru pasul E și pasul M, după execuția a 10 iterații se vor obține valorile  $\theta_A^{(10)} \approx 0.80$  și  $\theta_B^{(10)} \approx 0.52$ .

Este interesant de observat că estimarea obținută pentru parametrul  $\theta_A$  este acum la același nivel cu cea obținută prin metoda verosimilității maxime (MLE) în cazul observării tuturor variabilelor (0.80, vezi rezolvarea de la partea A, punctul *i*), iar estimarea obținută pentru parametrul  $\theta_B$  a coborât de la valoarea 0.58 care a fost obținută la prima iterație a algoritmului EM la o valoare (0.52) care este considerabil mai apropiată de estimarea prin metoda MLE (0.45).



The EM algorithm for solving  
a mixture of  $K$  categorical distributions

CMU, 2015 spring, Tom Mitchell, Nina Balcan, HW6, pr. 1



*Mixture models* are helpful for modelling unknown subpopulations in data. If we have a collection of data points  $X = \{X_1, \dots, X_n\}$ , where each  $X_i$  is [independently] drawn from one of  $K$  possible distributions, we can introduce a discrete-valued random variable  $Z_i \in \{1, \dots, K\}$  that indicates which distribution  $X_i$  is drawn from.

This exercise deals with the *categorical mixture model*, where each observation  $X_i$  is a discrete value drawn from a categorical distribution. The *parameter* for a categorical distribution is a  $K$ -dimensional vector  $\pi$  that lists the probability of each of  $K$  possible values (therefore,  $\sum_k \pi_k = 1$ ).

For *example*, suppose our categorical mixture model has 3 underlying distributions. Then, each  $Z_i$  could take on one of three values,  $\{1, 2, 3\}$ , with respective probabilities  $\pi_1, \pi_2, \pi_3$ ; equivalently,  $Z_i \sim \text{Categorical}(\pi)$ , where  $\pi = (\pi_1, \pi_2, \pi_3) \in \mathbb{R}_+^3$ , with  $\pi_1 + \pi_2 + \pi_3 = 1$ . The *observation*  $X_i$  is then generated from another categorical distribution, depending on the value of  $Z_i$ .

The *generative process* for a *categorical mixture model* is summarized as follows:

$$\begin{aligned} Z_i &\sim \text{Categorical}(\pi) \\ X_i &\sim \text{Categorical}(\theta_{Z_i}) \end{aligned}$$

For this model, where we observe  $X$  but not  $Z$ , we want to *learn* the *parameters* of the  $K$  categorical components  $\Theta = \{\pi, \theta_1, \dots, \theta_K\}$ , where each  $\theta_k \in \mathbb{R}^M$  is the parameter for the categorical distribution associated with the  $k$ -th mixture component. (It implies that each  $X_i$  can take on one of  $M$  possible values). We will use the EM algorithm to accomplish this.

A *note* on notation and a *hint*:

When working with categorical distributions it is helpful to use *indicator functions*. The indicator function  $1_{\{x=j\}}$  has the value 1 when  $x = j$  and 0 otherwise. With this notation, we can express the probability that a random variable drawn from a categorical distribution (e.g.,  $Y \sim \text{Categorical}(\phi)$ , where  $\phi \in \mathbb{R}^N$ ) takes on a particular value as

$$P(Y) = \prod_{i=1}^N \phi_i^{1_{\{Y=i\}}}.$$

- a. What is the *joint distribution*  $P(X, Z; \Theta)$ ?
- b. What is the *posterior distribution* of the *latent variables*,  $P(Z|X; \Theta)$ ?
- c. Compute the *expectation of the log-likelihood*,

$$Q(\Theta|\Theta') = E_{Z|X;\Theta'}[\log P(X, Z; \Theta')]$$

- d. What is the update step for  $\Theta$ ? That is, what  $\Theta$  maximizes  $Q(\Theta|\Theta')$ ?

*Hint:* Make sure your solution enforces the constraint that parameters to categorical distributions must sum to 1. Lagrange multipliers are a great way to solve constrained optimization problems.

**Answer:**

**a.**

$$\begin{aligned}
 P(X, Z; \Theta) &= \prod_{i=1}^n P(X_i, Z_i; \Theta) = \prod_{i=1}^n P(X_i|Z_i; \Theta)P(Z_i; \Theta) \\
 &= \prod_{i=1}^n \prod_{k=1}^K \left[ \pi_k \prod_{j=1}^M \theta_{kj}^{1_{\{X_i=v_j\}}} \right]^{1_{\{Z_i=k\}}}
 \end{aligned}$$

**b.**

$$\begin{aligned}
 P(Z_i = k|X_i; \Theta) &\stackrel{\text{Bayes F.}}{=} \frac{P(X_i|Z_i = k; \Theta)P(Z_i = k; \Theta)}{P(X_i; \Theta)} \\
 &= \frac{P(X_i|Z_i = k; \Theta)P(Z_i = k; \Theta)}{\sum_{l=1}^K P(Z_i = l; \Theta)P(X_i|Z_i = l; \Theta)} = \frac{\prod_{k=1}^K \left[ \pi_k \prod_{j=1}^M \theta_{kj}^{1_{\{X_i=v_j\}}} \right]^{1_{\{Z_i=k\}}}}{\sum_{l=1}^K \pi_l \prod_{j=1}^M \theta_{lj}^{1_{\{X_i=v_j\}}}}
 \end{aligned}$$

c.

$$\begin{aligned}
Q(\Theta|\Theta') &\stackrel{def.}{=} E_{Z|X;\Theta'}[\log P(X, Z; \Theta)] \\
&\stackrel{a.}{=} E_{Z|X;\Theta'} \left[ \log \prod_{i=1}^n \prod_{k=1}^K \left[ \pi_k \prod_{j=1}^M \theta_{kj}^{1_{\{X_i=v_j\}}} \right]^{1_{\{Z_i=k\}}} \right] \\
&= E_{Z|X;\Theta'} \left[ \sum_{i=1}^n \sum_{k=1}^K 1_{\{Z_i=k\}} \left( \log \pi_k + \sum_{j=1}^M 1_{\{X_i=v_j\}} \log \theta_{kj} \right) \right] \\
&= \sum_{i=1}^n \sum_{k=1}^K \underbrace{E_{Z|X;\Theta'}[1_{\{Z_i=k\}}]}_{P(Z_i=k|X_i;\Theta)} \left( \log \pi_k + \sum_{j=1}^M 1_{\{X_i=v_j\}} \log \theta_{kj} \right)
\end{aligned}$$

Using

$$\gamma_{ik} \stackrel{not.}{=} E_{Z|X;\Theta'}[1_{\{Z_i=k\}}] = P(Z_i = k|X_i; \Theta') \stackrel{b.}{=} \frac{\pi'_k \prod_{j=1}^M (\theta'_{kj})^{1_{\{X_i=v_j\}}}}{\sum_{l=1}^K \pi'_l \prod_{j=1}^M (\theta'_{lj})^{1_{\{X_i=v_j\}}}},$$

leads to

$$Q(\Theta|\Theta') = \sum_{i=1}^n \sum_{k=1}^K \gamma_{ik} \left( \log \pi_k + \sum_{j=1}^M 1_{\{X_i=v_j\}} \log \theta_{kj} \right)$$

d. To find the updates rules for the parameters, we solve the following *constrained optimization problem*:

$$\begin{aligned} & \max_{\pi, \theta_{1:K}} Q(\Theta|\Theta') \\ \text{subject to } & \sum_{k=1}^K \pi_k = 1 \\ & \sum_{j=1}^M \theta_{kj} = 1, \forall k = 1, \dots, K \end{aligned}$$

We introduce a *Lagrange multiplier* for each constraint associated to this problem and form the *Lagrangian functional*:

$$\mathcal{L} = Q(\Theta|\Theta') + \lambda_{\pi} \left( 1 - \sum_{k=1}^K \pi_k \right) + \sum_{k=1}^K \lambda_{\theta_k} \left( 1 - \sum_{j=1}^M \theta_{kj} \right)$$

Now we will differentiate the Lagrangian functional with respect to each parameter we want to optimize, set the derivative to zero, and solve for the optimal value.

We'll start with **optimizing for  $\pi$** , the parameter to the categorical distribution for the latent variables  $Z$ .

$$\frac{\partial \mathcal{L}}{\partial \pi_k} = 0 \Leftrightarrow \sum_{i=1}^n \gamma_{ik} \cdot \frac{1}{\pi_k} - \lambda_\pi = 0 \Leftrightarrow \pi_k^* = \frac{1}{\lambda_\pi} \sum_{i=1}^n \gamma_{ik}$$

We can find the value of the Lagrange multiplier  $\lambda_\pi$  by enforcing the constraint:

$$\begin{aligned} \sum_{k=1}^K \pi_k^* = 1 &\Leftrightarrow \sum_{k=1}^K \frac{1}{\lambda_\pi} \sum_{i=1}^n \gamma_{ik} = 1 \Leftrightarrow \frac{1}{\lambda_\pi} \sum_{k=1}^K \sum_{i=1}^n \gamma_{ik} = 1 \Leftrightarrow \\ \lambda_\pi &= \sum_{k=1}^K \sum_{i=1}^n \gamma_{ik} \Leftrightarrow \lambda_\pi = \sum_{i=1}^n \underbrace{\sum_{k=1}^K \gamma_{ik}}_1 \Leftrightarrow \lambda_\pi = n \end{aligned}$$

Substituting this value for  $\lambda_\pi$  back into the expression of  $\pi_k^*$  gives the answer:

$$\pi_k^* = \frac{1}{n} \sum_{i=1}^n \gamma_{ik}$$

The process for **optimizing the observation parameters**  $\theta_{kj}$  is similar. First, we differentiate the Lagrangian  $\mathcal{L}$  with respect to  $\theta_{kj}$ , set to zero, and [finally] solve.

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial \theta_{kj}} = 0 &\Leftrightarrow \sum_{i=1}^n \gamma_{ik} 1_{\{X_i=v_j\}} \cdot \frac{1}{\theta_{kj}} - \lambda_{\theta_k} = 0 \\ \Rightarrow \theta_{kj}^* &= \frac{1}{\lambda_{\theta_k}} \sum_{i=1}^n \gamma_{ik} 1_{\{X_i=v_j\}}\end{aligned}$$

Again, we can find the value of the Lagrange multiplier  $\lambda_{\theta_k}$  by enforcing the summation constraint:

$$\begin{aligned}\sum_{j=1}^M \theta_{kj}^* = 1 &\Leftrightarrow \sum_{j=1}^M \frac{1}{\lambda_{\theta_k}} \sum_{i=1}^n \gamma_{ik} 1_{\{X_i=v_j\}} = 1 \Leftrightarrow \frac{1}{\lambda_{\theta_k}} \sum_{j=1}^M \sum_{i=1}^n \gamma_{ik} 1_{\{X_i=v_j\}} = 1 \\ \Leftrightarrow \lambda_{\theta_k} &= \sum_{j=1}^M \sum_{i=1}^n \gamma_{ik} 1_{\{X_i=v_j\}}\end{aligned}$$

Substituting this value for  $\lambda_{\theta_k}$  back into the expression of  $\theta_{kj}^*$  gives the answer:

$$\theta_{kj}^* = \frac{\sum_{i=1}^n \gamma_{ik} 1_{\{X_i=v_j\}}}{\sum_{l=1}^M \sum_{i=1}^n \gamma_{ik} 1_{\{X_i=v_l\}}} = \frac{\sum_{i=1}^n \gamma_{ik} 1_{\{X_i=v_j\}}}{\sum_{i=1}^n \sum_{l=1}^M \gamma_{ik} 1_{\{X_i=v_l\}}}$$



The EM algorithm for solving  
a mixture of  $K$  categorical distributions,  
applied to the problem of Word Sense Disambiguation,  
i.e. identifying the semantic domains associated to words in a  
text document

CMU, 2012 fall, Eric Xing, Aarti Singh, HW3, pr. 3

The objective of this exercise is to derive the update equations of the EM algorithm for optimizing the latent variables [designating the semantic domains] involved in generating a text document.

Each *word* will be seen as a random variable  $w$  that can take values  $1, \dots, V$  from the *vocabulary* of words. In fact, we will denote each  $w$  by an array of  $V$  components such that  $w(i) = 1$  if  $w$  takes the value of the  $i$ -th word in the vocabulary. Hence,  $\sum_{i=1}^V w(i) = 1$ .

Given a *document* containing words  $w_j$ ,  $j = 1, \dots, N$ , where  $N$  is the length of the document, we will assume that these words are generated from a mixture of  $K$  discrete *topics*:

$$P(w) = \sum_{m=1}^K \pi_m P(w|\mu_m) \text{ and } P(w|\mu_m) = \prod_{i=1}^V \mu_m(i)^{w(i)},$$

where

$\pi_m$  denotes the prior [probability] for the latent topic variable  $t = m$ ,

$\mu_k \stackrel{\text{not.}}{=} (\mu_k(1), \dots, \mu_k(i), \dots, \mu_k(V))$ , with  $\mu_k(i) \geq 0$  for  $i = 1, \dots, V$  and  $\sum_{i=1}^V \mu_m(i) = 1$  for each  $k = 1, \dots, K$ , and

$\mu_m(i) \stackrel{\text{not.}}{=} P(w(i) = 1 | t = m)$ .

a. In the expectation step, for each word  $w_j$ , compute

$$F_j(t) \stackrel{not.}{=} P(t|w_j; \theta),$$

the probability that  $w_j$  belongs to each of the  $K$  topics, where  $\theta$  is the set of parameters of this mixture model.

**Answer:**

$$\begin{aligned} F_j(t_j = m) &\stackrel{not.}{=} P(t_j = m|w_j; \theta) \stackrel{Bayes\ T.}{=} \frac{P(w_j|t_j = m; \theta) P(t_j = m|\theta)}{P(w_j|\theta)} \\ &= \frac{\pi_m P(w_j|\mu_m)}{\sum_{m'=1}^K \pi_{m'} P(w_j|\mu_{m'})} = \frac{\pi_m \prod_{l=1}^V \mu_m(l)^{w_j(l)}}{\sum_{m'=1}^K \pi_{m'} \prod_{l=1}^V \mu_{m'}(l)^{w_j(l)}} \end{aligned}$$

b. In the maximization step, compute  $\theta$  that maximizes the log-likelihood of the data

$$l(w; \theta) \stackrel{not.}{=} \log \prod_{j=1}^N P(w_j; \theta)$$

**Hint:** Summing over the latent topic variable, we can write  $l(w; \theta)$  as

$$l(w; \theta) = \sum_{j=1}^N \log \sum_t P(w_j, t; \theta) = \sum_{j=1}^N \log F_j(t) \frac{\sum_t P(w_j, t; \theta)}{F_j(t)}$$

Further on, using Jensen's inequality (see pr. EM-2) we get:

$$l(w; \theta) \geq \sum_{j=1}^N \sum_t F_j(t) \log \frac{P(w_j, t; \theta)}{F_j(t)} = \sum_{j=1}^N \sum_t F_j(t) \log P(w_j, t; \theta) + H(F_j)$$

Hence compute  $\theta$  as:

$$\operatorname{argmax}_{\theta} \sum_{j=1}^N \sum_t F_j(t) \log P(w_j, t; \theta)$$

**Answer:**

$$\begin{aligned}
\theta &= \operatorname{argmax}_{\theta} \sum_{j=1}^N \sum_{m=1}^K F_j(t_j = m) \log P(w_j, t_j = m | \theta) \\
&= \operatorname{argmax}_{\theta} \sum_{j=1}^N \sum_{m=1}^K F_j(t_j = m) \log P(w_j | t_j = m; \theta) P(t_j = m | \theta) \\
&= \operatorname{argmax}_{\theta} \sum_{j=1}^N \sum_{m=1}^K F_j(t_j = m) \log \left( \pi_m \prod_{l=1}^V \mu_m(l)^{w_j(l)} \right) \\
&= \operatorname{argmax}_{\theta} \sum_{j=1}^N \sum_{m=1}^K [F_j(t_j = m) \log \pi_m + F_j(t_j = m) \sum_{l=1}^V \log \mu_m(l)^{w_j(l)}] \\
&= \operatorname{argmax}_{\theta} \sum_{j=1}^N \sum_{m=1}^K [F_j(t_j = m) \log \pi_m + F_j(t_j = m) \sum_{l=1}^V w_j(l) \log \mu_m(l)] \quad (3)
\end{aligned}$$

To optimize  $\mu_m(l)$ :

After eliminating from (3) the terms which are constant with respect to  $\mu_m$  we get:

$$\sum_{j=1}^N F_j(t_j = m) \sum_{l=1}^V w_j(l) \log \mu_m(l)$$

We will use a Lagrangian to constrain  $\mu_m$  to be a probability distribution:

$$\mathcal{L}(\mu_m(l)) = \sum_{j=1}^N F_j(t_j = m) \sum_{l=1}^V w_j(l) \log \mu_m(l) + \beta \left( \sum_{l=1}^V \mu_m(l) - 1 \right)$$

Then, solving for  $\mu_m(l)$ :

$$\begin{aligned} \frac{\partial}{\partial \mu_m(l)} \mathcal{L}(\mu_m(l)) = 0 &\Leftrightarrow \sum_{j=1}^N F_j(t_j = m) \frac{w_j(l)}{\mu_m(l)} + \beta = 0 \\ \Leftrightarrow \frac{1}{\mu_m(l)} \sum_{j=1}^N F_j(t_j = m) w_j(l) + \beta &= 0 \Leftrightarrow \frac{1}{\mu_m(l)} = \frac{-\beta}{\sum_{j=1}^N F_j(t_j = m) w_j(l)} \\ \Leftrightarrow \mu_m(l) &= \frac{\sum_{j=1}^N F_j(t_j = m) w_j(l)}{-\beta} \end{aligned} \tag{4}$$

Knowing that  $\sum_{l=1}^V \mu_m(l) = 1$ , we have:

$$\sum_{l=1}^V \frac{\sum_{j=1}^N F_j(t_j = m) w_j(l)}{-\beta} = 1 \Leftrightarrow -\beta = \sum_{l=1}^V \sum_{j=1}^N F_j(t_j = m) w_j(l)$$

Hence, substituting for  $-\beta$  in (4), we get:

$$\begin{aligned} \mu_m(l) &= \frac{\sum_{j=1}^N F_j(t_j = m) w_j(l)}{\sum_{l=1}^V \sum_{j=1}^N F_j(t_j = m) w_j(l)} = \frac{\sum_{j=1}^N F_j(t_j = m) w_j(l)}{\sum_{j=1}^N \sum_{l=1}^V F_j(t_j = m) w_j(l)} \\ &= \frac{\sum_{j=1}^N F_j(t_j = m) w_j(l)}{\sum_{j=1}^N F_j(t_j = m) \underbrace{\sum_{l=1}^V w_j(l)}_1} = \frac{\sum_{j=1}^N F_j(t_j = m) w_j(l)}{\sum_{j=1}^N F_j(t_j = m)} \end{aligned}$$

**Note:** Intuitively the last expression can be interpreted as the portion [of word[ occurrence]s] that had  $w(l) = 1$  among all words [in the given document] which are deemed to belong to cluster  $m$ .

To optimize  $\pi_m$ , we proceed similarly:

We begin by removing from (3) the terms that are constant with respect to  $\pi_m$ , thus getting:

$$\sum_{j=1}^N F_j(t_j = m) \log \pi_m$$

So, using the Lagrangian with the constraint that  $\sum_{m=1}^K \pi_m = 1$

$$\mathcal{L}(\pi_m) = \sum_{j=1}^N F_j(t_j = m) \log \pi_m + \beta \left( \sum_{m=1}^K \pi_m - 1 \right),$$

and solving for  $\pi_m$ , we have:

$$\begin{aligned} \frac{\partial}{\partial \pi_m} \mathcal{L}(\pi_m) = 0 &\Leftrightarrow \sum_{j=1}^N \frac{F_j(t_j = m)}{\pi_m} + \beta = 0 \Leftrightarrow \frac{1}{\pi_m} \sum_{j=1}^N F_j(t_j = m) = -\beta \\ &\Leftrightarrow \pi_m = \frac{\sum_{j=1}^N F_j(t_j = m)}{-\beta} \end{aligned} \tag{5}$$



Since  $\sum_{m=1}^K \pi_m = 1$ , it gives us:

$$\begin{aligned} \sum_{m=1}^K \frac{\sum_{j=1}^N F_j(t_j = m)}{-\beta} = 1 &\Leftrightarrow \frac{1}{-\beta} \sum_{m=1}^K \sum_{j=1}^N F_j(t_j = m) = 1 \\ &\Leftrightarrow -\beta = \sum_{m=1}^K \sum_{j=1}^N F_j(t_j = m) \end{aligned}$$

Substituting for  $-\beta$  in (5), we get:

$$\pi_m = \frac{\sum_{j=1}^N F_j(t_j = m)}{\sum_{m=1}^K \sum_{j=1}^N F_j(t_j = m)} = \frac{\sum_{j=1}^N F_j(t_j = m)}{\underbrace{\sum_{j=1}^N \sum_{m=1}^K F_j(t_j = m)}_1} = \frac{\sum_{j=1}^N F_j(t_j = m)}{N}$$

**Note:** Intuitively the last expression can be interpreted as the portion [of word[ occurrence]s] that belongs to cluster  $m$  among the total of  $N$  words.

**To summarize:**

**E step:**

$$F_j(t_j = m) = \frac{\pi_m \prod_{l=1}^V \mu_m(l)^{w_j(l)}}{\sum_{m'=1}^K \pi_{m'} \prod_{l=1}^V \mu_{m'}(l)^{w_j(l)}} \text{ for } j = 1, \dots, N \text{ and } m = 1, \dots, K$$

**M step:**

$$\mu_m(l) = \frac{\sum_{j=1}^N F_j(t_j = m) w_j(l)}{\sum_{j=1}^N F_j(t_j = m)} \text{ for } m = 1, \dots, K \text{ and } l = 1, \dots, V$$

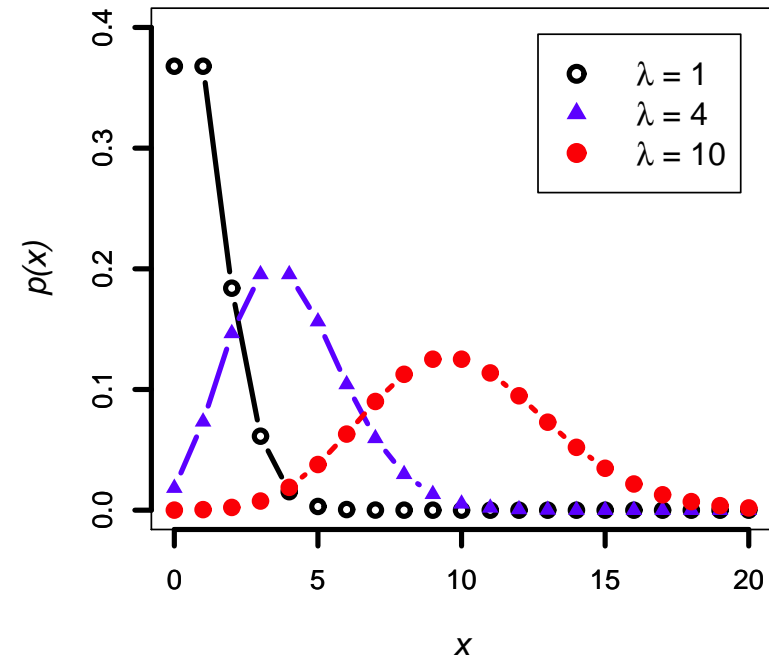
$$\pi_m = \frac{\sum_{j=1}^N F_j(t_j = m)}{N} \text{ for } m = 1, \dots, K$$

**The EM algorithm for  
solving mixtures of Poisson distributions**

University of Utah, 2008 fall, Hal Daumé III, HW9, pr. 1

The Poisson distribution is a distribution over positive count values; for a count  $x$  with parameter  $\lambda$ , the Poisson has the form  $Poisson(x|\lambda) = \frac{1}{e^\lambda} \cdot \frac{\lambda^x}{x!}$ .

It can be proven that the maximum likelihood estimation for  $\lambda$  given a sequence of counts  $x_1, \dots, x_n$  was simply  $\frac{1}{n} \sum_{i=1}^n x_i$  – the mean of the counts.



Let's consider a generalization of this: *the Poisson mixture model*. This is actually used in *web server monitoring*. The number of accesses to a web server in a minute typically follows a Poisson distribution.

Suppose we have  $n$  web servers we are monitoring and we monitor each for  $M$  minutes. Thus, we have  $n \times M$  counts; call  $x_{i,m}$  the number of hits to web server  $i$  in minute  $m$ . Our goal is to *cluster* the web servers according to their hit frequency.

Using the EM algorithm, construct a Poisson mixture model for this problem and compute the expectations and maximization steps for this model.

**Hint:** Suppose we want  $K$  clusters; let  $z_i$  be the latent variable telling us which cluster the web server  $i$  belongs to (from 1 to  $K$ ). Let  $\lambda_k$  denote the parameter for the Poisson distribution for cluster  $k$ , and  $\pi_k$  the . Then, the complete data likelihood is:

$$\begin{aligned}
 L(\bar{\lambda}, \bar{\pi}) &\stackrel{\text{def.}}{=} p(\bar{x}, \bar{z} | \bar{\lambda}, \bar{\pi}) \stackrel{i.i.d.}{=} \prod_{i=1}^n p(x_i, z_i | \bar{\lambda}, \bar{\pi}) = \prod_{i=1}^n p(x_i | z_i, \bar{\lambda}, \bar{\pi}) \cdot \underbrace{p(\underbrace{z_i}_{k} | \bar{\lambda}, \bar{\pi})}_{\pi_k} \\
 &= \prod_{i=1}^n \prod_{k=1}^K \left[ \pi_k \prod_{m=1}^M \text{Poisson}(x_{i,m} | \lambda_k) \right]^{1_{\{z_i=k\}}}
 \end{aligned}$$

where

$\bar{x} \stackrel{\text{not.}}{=} (x_1, \dots, x_n)$ , with  $x_i \stackrel{\text{not.}}{=} (x_{i,1}, \dots, x_{i,M})$  for  $i = 1, \dots, n$ ,  
 $\bar{z} \stackrel{\text{not.}}{=} (z_1, \dots, z_K)$ ,  $\bar{\pi} \stackrel{\text{not.}}{=} (\pi_1, \dots, \pi_K)$ ;  $\bar{\lambda} \stackrel{\text{not.}}{=} (\lambda_1, \dots, \lambda_K)$ , and  
 $1_{\{z_i=k\}}$  is one if  $z_i = k$  and zero otherwise.

## Solution

110.

The likelihood function:

$$L(\bar{\lambda}, \bar{\pi}) = \prod_{i=1}^n \prod_{k=1}^K \left[ \pi_k \prod_{m=1}^M \frac{1}{e^{\lambda_k}} \frac{\lambda_k^{x_{i,m}}}{x_{i,m}!} \right]^{1_{\{z_i=k\}}}$$

The likelihood function:

$$\begin{aligned} \ell(\bar{\lambda}, \bar{\pi}) &\stackrel{\text{def.}}{=} \ln L(\bar{\lambda}, \bar{\pi}) = \sum_{i=1}^n \sum_{k=1}^K 1_{\{z_i=k\}} \left[ \ln \pi_k + \sum_{m=1}^M \ln \frac{1}{e^{\lambda_k}} \frac{\lambda_k^{x_{i,m}}}{x_{i,m}!} \right] \\ &= \sum_{i=1}^n \sum_{k=1}^K 1_{\{z_i=k\}} \left[ \ln \pi_k + \sum_{m=1}^M (-\lambda_k + x_{i,m} \ln \lambda_k - \ln(x_{i,m}!)) \right] \end{aligned}$$

The auxiliary function:

$$\begin{aligned} Q(\bar{\lambda}, \bar{\pi} | \bar{\lambda}^{(t)}, \bar{\pi}^{(t)}) &\stackrel{\text{def.}}{=} E_{P(z_i=k|x_i, \bar{\lambda}^{(t)}, \bar{\pi}^{(t)})} [\ell(\bar{\lambda}, \bar{\pi})] \\ &= \sum_{i=1}^n \sum_{k=1}^K P(z_i = k | x_i, \bar{\lambda}^{(t)}, \bar{\pi}^{(t)}) \left[ \ln \pi_k + \sum_{m=1}^M (-\lambda_k + x_{i,m} \ln \lambda_k - \ln(x_{i,m}!)) \right] \end{aligned}$$

E-step:

$$\begin{aligned}
 p_{ik}^{(t)} &\stackrel{\text{not.}}{=} P(z_i = k | x_i, \bar{\lambda}^{(t)}, \bar{\pi}^{(t)}) \\
 &\stackrel{\text{Bayes F.}}{=} \frac{P(x_i | z_i = k, \bar{\lambda}^{(t)}, \bar{\pi}^{(t)}) \cdot P(z_i = k | \bar{\lambda}^{(t)}, \bar{\pi}^{(t)})}{\sum_{k'=1}^K P(x_i | z_i = k', \bar{\lambda}^{(t)}, \bar{\pi}^{(t)}) \cdot P(z_i = k' | \bar{\lambda}^{(t)}, \bar{\pi}^{(t)})} \\
 &= \frac{\pi_k^{(t)} \cdot P(x_i | z_i = k, \bar{\lambda}^{(t)}, \bar{\pi}^{(t)})}{\sum_{k'=1}^K \pi_{k'}^{(t)} \cdot P(x_i | z_i = k', \bar{\lambda}^{(t)}, \bar{\pi}^{(t)})} \\
 &= \frac{\pi_k^{(t)} \cdot \prod_{m=1}^M \text{Poisson}(x_{i,m} | \lambda_k)}{\sum_{k'=1}^K (\pi_{k'}^{(t)} \cdot \prod_{m=1}^M \text{Poisson}(x_{i,m} | \lambda_{k'}))}
 \end{aligned}$$



### M-step:

Since  $\sum_{k=1}^K \pi_k = 1$ , we have to introduce the Lagrange multiplier  $\lambda$ :

$$\begin{aligned} \frac{\partial}{\partial \pi_k} \left[ Q(\bar{\lambda}, \bar{\pi} | \bar{\lambda}^{(t)}, \bar{\pi}^{(t)}) + \lambda(1 - \sum_{k=1}^K \pi_k) \right] = 0 &\Leftrightarrow \\ \sum_{i=1}^n \frac{p_{ik}^{(t)}}{\pi_k} - \lambda = 0 &\Leftrightarrow \pi_k^{(t+1)} = \frac{1}{\lambda} \sum_{i=1}^n p_{ik}^{(t)} \end{aligned}$$

It is easy to see that  $\pi_k^{(t+1)}$  is indeed a maximum point for  $Q$  w.r.t. the parameter  $\pi_k$ .

Imposing the constraint  $\sum_{k=1}^K \pi_k^{(t+1)} = 1$  leads to:

$$\frac{1}{\lambda} \sum_{k=1}^K \sum_{i=1}^n p_{ik}^{(t)} = 1 \Leftrightarrow \frac{1}{\lambda} \sum_{i=1}^n \underbrace{\sum_{k=1}^K p_{ik}^{(t)}}_1 = 1 \Leftrightarrow \frac{n}{\lambda} = 1 \Leftrightarrow \lambda = n.$$

Therefore,  $\pi_k^{(t+1)} = \frac{1}{n} \sum_{i=1}^n p_{ik}^{(t)}$ , which is always a non-negative quantity.

Now, regarding the parameters  $\lambda_k$ :

$$\begin{aligned} \frac{\partial}{\partial \lambda_k} Q(\bar{\lambda}, \bar{\pi} | \bar{\lambda}^{(t)}, \bar{\pi}^{(t)}) = 0 &\Leftrightarrow \sum_{i=1}^n p_{ik} \left( \sum_{m=1}^M \left( -1 + x_{i,m} \cdot \frac{1}{\lambda_k} \right) \right) = 0 \Leftrightarrow \\ \sum_{i=1}^n p_{ik} \left( -M + \frac{1}{\lambda_k} \sum_{m=1}^M x_{i,m} \right) = 0 &\Leftrightarrow \frac{1}{\lambda_k} \sum_{i=1}^n p_{ik} \left( \sum_{m=1}^M x_{i,m} \right) = M \sum_{i=1}^n p_{ik} \Leftrightarrow \\ \lambda_k^{(t+1)} &= \frac{\sum_{i=1}^n p_{ik} \left( \sum_{m=1}^M x_{i,m} \right)}{M \sum_{i=1}^n p_{ik}} \end{aligned}$$

It can be easily shown that for each  $k \in \{1, \dots, K\}$ , the value  $\lambda_k^{(t+1)}$  is positive, and it maximizes  $Q$  w.r.t. the parameter  $\lambda_k$ .

# The EM algorithm for solving a Gamma Mixture Model

G. Vegas-Sánchez-Ferrero, M. Martín-Fernández, J. Miguel Sanches

*A Gamma Mixture Model for IVUS Imaging, 2014*

[ adapted by Liviu Ciortuz ]

## The Gamma Mixture Model

$$p(x|\theta) = \sum_{k=1}^K \pi_k \text{Gamma}(x|\theta_k),$$

with  $r \stackrel{\text{not.}}{=} (r_1, \dots, r_K)$ ,  $\beta \stackrel{\text{not.}}{=} (\beta_1, \dots, \beta_K)$ ,  $\theta = (r, \beta)$ , and  $\theta_k = (r_k, \beta_k)$  for  $k = 1, \dots, K$ , and  $\pi_k > 0$  for  $k = 1, \dots, K$  and  $\sum_{k=1}^K \pi_k = 1$ .

Consider  $x_1, \dots, x_n \in \mathbb{R}^+$ , each  $x_i$  being generated by one component of the above mixture, denoted by  $z_i \in \{1, \dots, K\}$ .

Find the maximum likelihood estimation of the parameters  $\theta$  by using the EM algorithm.

## Solution

- The verosimilarity of the “complete” data  $(x, z)$ , with  $x \stackrel{\text{not.}}{=} (x_1, \dots, x_n)$ , and  $z \stackrel{\text{not.}}{=} (z_1, \dots, z_n)$ , can be written as follows:

$$p(x, z|\theta) \stackrel{i.i.d.}{=} \prod_{i=1}^n p(x_i|z_i, \theta) = \prod_{i=1}^n \prod_{k=1}^K [p(x_i|z_i, \theta) \cdot \underbrace{p(z_i|\theta)}_{\pi_k}]^{1_{\{z_i=k\}}} = \prod_{i=1}^n \prod_{k=1}^K \left[ \underbrace{p(x_i|z_i, \theta)}_{\text{Gamma}(x_i|\theta_k)} \cdot \pi_k \right]^{1_{\{z_i=k\}}}$$

where  $1_{\{z_i=k\}}$  is the indicator function.

- The log-verosimilarity function:

$$\begin{aligned} \ell(\theta) &\stackrel{\text{def.}}{=} \ln p(x, z|\theta) = \sum_{i=1}^n \sum_{k=1}^K 1_{\{z_i=k\}} [\ln \pi_k + \ln \text{Gamma}(x_i|\theta_k)] \\ &\stackrel{\theta_k=(r_k, \beta_k)}{=} \sum_{i=1}^n \sum_{k=1}^K 1_{\{z_i=k\}} \left[ \ln \pi_k - r_k \ln \beta_k - \ln \Gamma(r_k) + (r_k - 1) \ln x_i - \frac{x_i}{\beta_k} \right] \end{aligned}$$

- The auxiliary function:

$$\begin{aligned} Q(\theta|\theta^{(t)}) &\stackrel{\theta^{(t)}=(r^{(t)}, \beta^{(t)})}{=} E_{P(Z|X, \theta^{(t)})} \left[ \sum_{i=1}^n \sum_{k=1}^K 1_{\{z_i=k\}} \left[ \ln \pi_k - r_k \ln \beta_k - \ln \Gamma(r_k) + (r_k - 1) \ln x_i - \frac{x_i}{\beta_k} \right] \right] \\ &\stackrel{\text{lin. of exp.}}{=} \sum_{i=1}^n \sum_{k=1}^K \underbrace{P(z_i = k|x_i, \theta^{(t)})}_{\text{not. } \gamma_{ik}} \left[ \ln \pi_k - r_k \ln \beta_k - \ln \Gamma(r_k) + (r_k - 1) \ln x_i - \frac{x_i}{\beta_k} \right] \\ &= \sum_{i=1}^n \sum_{k=1}^K \gamma_{ik} \left[ \ln \pi_k - r_k \ln \beta_k - \ln \Gamma(r_k) + (r_k - 1) \ln x_i - \frac{x_i}{\beta_k} \right] \end{aligned}$$

## E step

$$\begin{aligned}
 \gamma_{ik}^{(t+1)} &\stackrel{\text{not.}}{=} P(z_i = k | x_i, \theta^{(t)}) \stackrel{\text{Bayes F.}}{=} \frac{\text{Gamma}(x_i | z_i = k, \theta) \cdot \overbrace{P(z_i = k | \theta)}^{\pi_k}}{\sum_{k'=1}^K \text{Gamma}(x_i | z_i = k', \theta) \cdot \underbrace{P(z_i = k' | \theta)}_{\pi_{k'}}} \\
 &= \frac{\text{Gamma}(x_i | \theta_k) \cdot \pi_k}{\sum_{k'=1}^K \text{Gamma}(x_i | \theta_{k'}) \cdot \pi_{k'}}
 \end{aligned}$$

**Note that**  $\gamma_{ik}^{(t+1)} > 0$  **for**  $k = 1, \dots, K$  **because**  $x_i > 0$  **for**  $i = 1, \dots, n$ .

## M step: $\pi_k$

Because  $\sum_{k=1}^K \pi_k = 1$ , we will use a Lagrange multiplier  $\lambda$ , and solve as usually:

$$\frac{\partial}{\partial \pi_k} \left[ Q(\theta | \theta^{(t)}) + \lambda \left( 1 - \sum_{k=1}^K \pi_k \right) \right] = 0 \Leftrightarrow \sum_{i=1}^n \gamma_{ik} \frac{1}{\pi_k} - \lambda = 0 \Leftrightarrow \frac{1}{\pi_k} \sum_{i=1}^n \gamma_{ik} = \lambda \Leftrightarrow$$

$$\pi_k^{(t+1)} = \frac{1}{\lambda} \sum_{i=1}^n \gamma_{ik}$$

By substituting this expression into the constraint  $\sum_{k=1}^K \pi_k = 1$ , we will get:

$$\frac{1}{\lambda} \sum_{k=1}^K \sum_{i=1}^n \gamma_{ik} = 1 \Leftrightarrow \frac{1}{\lambda} \sum_{i=1}^n \underbrace{\sum_{k=1}^K \gamma_{ik}}_1 = 1 \Leftrightarrow \frac{1}{\lambda} n = 1 \Leftrightarrow \lambda = n$$

Therefore,

$$\pi_k^{(t+1)} = \frac{1}{n} \sum_{i=1}^n \gamma_{ik} k^{(t+1)}$$

Note that  $\pi_k^{(t+1)} > 0$  for  $k = 1, \dots, K$ .

**M step:**  $\beta_k$

$$\begin{aligned} \frac{\partial}{\partial \beta_k} Q(\theta | \theta^{(t)}) = 0 &\Leftrightarrow \sum_{i=1}^n \gamma_{ik} \left( -\frac{r_k}{\beta_k} + \frac{x_i}{\beta_k^2} \right) = 0 \Leftrightarrow \frac{1}{\beta_k^2} \left[ \sum_{i=1}^n \gamma_{ik} (x_i - r_k \beta_k) \right] = 0 \\ &\Leftrightarrow \sum_{i=1}^n \gamma_{ik} (x_i - r_k \beta_k) = 0 \Leftrightarrow \sum_{i=1}^n \gamma_{ik} r_k \beta_k = \sum_{i=1}^n \gamma_{ik} x_i \Leftrightarrow r_k \beta_k \sum_{i=1}^n \gamma_{ik} = \sum_{i=1}^n \gamma_{ik} x_i \end{aligned}$$

Therefore,

$$\beta_k^{(t+1)} = \frac{\sum_{i=1}^n \gamma_{ik}^{(t+1)} x_i^{(t+1)}}{r_k \sum_{i=1}^n \gamma_{ik}} \quad \text{for } k = 1, \dots, K.$$

It can be easily shown that  $\beta_k^{(t+1)}$  is always positive, and it maximizes  $Q(\theta | \theta^{(t)})$  w.r.t.  $\beta_k$ .



## M step: $r_k$

By substituting [the expression of]  $\beta_k^{(t+1)}$  into  $Q(\theta|\theta^{(t)})$ , we will get:

$$Q(r, \beta^{(t+1)}|\theta^{(t)}) = \sum_{i=1}^n \sum_{k=1}^K \gamma_{ik} \left[ \ln \pi_k - r_k \ln \beta_k^{(t+1)} - \ln \Gamma(r_k) + (r_k - 1) \ln x_i - \frac{x_i}{\beta_k^{(t+1)}} \right]$$

Then, solving for the partial derivative of  $Q(r, \beta^{(t+1)}|\theta^{(t)})$  w.r.t.  $r_k$ , amounts to:

$$\begin{aligned} \frac{\partial}{\partial r_k} Q(r, \beta^{(t+1)}|\theta^{(t)}) &= 0 \Leftrightarrow \\ \frac{\partial}{\partial r_k} \sum_{i=1}^n \gamma_{ik}^{(t+1)} \left[ \ln \pi_k - r_k \ln \beta_k^{(t+1)} - \ln \Gamma(r_k) + (r_k - 1) \ln x_i - \frac{x_i}{\beta_k^{(t+1)}} \right] &= 0 \Leftrightarrow \\ \sum_{i=1}^n \gamma_{ik}^{(t+1)} \left[ -\ln \beta_k^{(t+1)} - r_k \frac{\partial}{\partial r_k} \ln \beta_k^{(t+1)} - \frac{\Gamma'(r_k)}{\Gamma(r_k)} + \ln x_i - x_i \frac{\partial}{\partial r_k} \frac{1}{\beta_k^{(t+1)}} \right] &= 0 \end{aligned}$$

Since

$$\frac{\partial}{\partial r_k} \ln \beta_k^{(t+1)} = \frac{\partial}{\partial r_k} \ln \frac{\sum_{i=1}^n \gamma_{ik}^{(t+1)} x_i}{r_k \sum_{i=1}^n \gamma_{ik}^{(t+1)}} = \frac{\partial}{\partial r_k} \left( -\ln r_k + \ln \frac{\sum_{i=1}^n \gamma_{ik}^{(t+1)} x_i}{\sum_{i=1}^n \gamma_{ik}^{(t+1)}} \right) = -\frac{1}{r_k}$$

and

$$\frac{\partial}{\partial r_k} \frac{1}{\beta_k^{(t+1)}} = \frac{\partial}{\partial r_k} \frac{r_k \sum_{i=1}^n \gamma_{ik}^{(t+1)}}{\sum_{i=1}^n \gamma_{ik}^{(t+1)} x_i} = \frac{\sum_{i=1}^n \gamma_{ik}^{(t+1)}}{\sum_{i=1}^n \gamma_{ik}^{(t+1)} x_i},$$

it follows that

$$\begin{aligned} & \frac{\partial}{\partial r_k} Q(r, \beta^{(t+1)} | \theta^{(t)}) = 0 \Leftrightarrow \\ & \sum_{i=1}^n \gamma_{ik}^{(t+1)} \left[ \ln r_k - \ln \frac{\sum_{i'=1}^n \gamma_{i'k}^{(t+1)} x_{i'}}{\sum_{i'=1}^n \gamma_{i'k}^{(t+1)}} + 1 - \frac{\Gamma'(r_k)}{\Gamma(r_k)} + \ln x_i - x_i \frac{\sum_{i'=1}^n \gamma_{i'k}^{(t+1)}}{\sum_{i'=1}^n \gamma_{i'k}^{(t+1)} x_{i'}} \right] = 0 \end{aligned}$$

Therefore,

$$\begin{aligned} & \left( \ln r_k - \frac{\Gamma'(r_k)}{\Gamma(r_k)} \right) \sum_{i=1}^n \gamma_{ik}^{(t+1)} = \\ & \sum_{i=1}^n \gamma_{ik}^{(t+1)} \ln \frac{\sum_{i'=1}^n \gamma_{i'k}^{(t+1)} x_{i'}}{\sum_{i'=1}^n \gamma_{i'k}^{(t+1)}} - \sum_{i=1}^n \gamma_{ik}^{(t+1)} - \sum_{i=1}^n \gamma_{ik}^{(t+1)} \ln x_i + \sum_{i=1}^n \gamma_{ik}^{(t+1)} x_i \cdot \frac{\sum_{i'=1}^n \gamma_{i'k}^{(t+1)}}{\sum_{i'=1}^n \gamma_{i'k}^{(t+1)} x_{i'}} \Leftrightarrow \\ & \left( \ln r_k - \frac{\Gamma'(r_k)}{\Gamma(r_k)} \right) \sum_{i=1}^n \gamma_{ik}^{(t+1)} = \sum_{i=1}^n \gamma_{ik}^{(t+1)} \ln \frac{\sum_{i'=1}^n \gamma_{i'k}^{(t+1)} x_{i'}}{\sum_{i'=1}^n \gamma_{i'k}^{(t+1)}} - \sum_{i=1}^n \gamma_{ik}^{(t+1)} \ln x_i \end{aligned}$$

So,

$$\ln r_{ik} - \frac{\Gamma'(r_k)}{\Gamma(r_k)} = \ln \frac{\sum_{i=1}^n \gamma_{ik}^{(t+1)} x_i}{\sum_{i=1}^n \gamma_{ik}^{(t+1)}} - \frac{\sum_{i=1}^n \gamma_{ik}^{(t+1)} \ln x_i}{\sum_{i=1}^n \gamma_{ik}^{(t+1)}}$$

Unfortunately, there is no known close form of  $r_k$  satisfying this equation.

However, it can be shown that the function  $\ln r_{ik} - \frac{\Gamma'(r_k)}{\Gamma(r_k)}$  is well-behaved, yielding a unique solution to the above equation, which will be denoted  $r_{ik}^{(t+1)}$ ; it can be found through certain numerical methods.

**Note:** The function  $\psi(r) \stackrel{not}{=} \ln r_{ik} - \frac{\Gamma'(r_k)}{\Gamma(r_k)}$  is the so-called **di-gamma function**.

## Note

It is interesting to see that if in the *update equations* of this EM algorithm for solving the Gamma Mixture Model we substitute  $r = r_k$  and  $\gamma_k = \frac{1}{k}$  for  $k = 1, \dots, K$ , we obtain exactly the Maximum Likelihood estimations for the parameters  $r$  and  $\beta$  of the Gamma distribution.

Using the EM algorithm for  
estimating the *selection probability* for a mixture  
of two (arbitrary) distributions

CMU, 2006 spring, Carlos Guestrin, final exam, pr. 8

CMU, 2004 fall, Carlos Guestrin, HW2, pr. 2.1

We want to derive an EM algorithm for estimating the mixing parameter for a mixture of arbitrary probability densities  $f_1$  and  $f_2$ .

For *example*,  $f_1(x)$  could be a standard normal distribution centered at 0, and  $f_2(x)$  could be the uniform distribution between  $[0, 1]$ . You can think about such mixtures in the following way: First, you flip a coin. With probability  $\lambda$  (i.e., the coin comes up *heads*), you will sample  $x$  from density  $f_1$ , and with probability  $(1 - \lambda)$  you sample from density  $f_2$ .

More formally, let  $f_\lambda(x) = \lambda f_1(x) + (1 - \lambda) f_2(x)$ , where  $f_1$  and  $f_2$  are arbitrary probability density functions on  $\mathbb{R}$ , and  $\lambda \in [0, 1]$  is an unknown mixture parameter.

a.

Given a data point  $x$ , and a value for the mixture parameter  $\lambda$ , compute the probability that  $x$  was generated from density  $f_1$ .

b.

Now, suppose you are given a data set  $\{x_1, \dots, x_n\}$  drawn i.i.d. from the mixture density, and a set of coin flips  $\{z_1, z_2, \dots, z_n\}$ , such that  $z_i = 1$  means that  $x_i$  is a sample from  $f_1$ , and  $z_i = 0$  means that  $x_i$  was generated from density  $f_2$ .

For a fixed parameter  $\lambda$ , compute the complete log-likelihood of the data, i.e.,  $\ln P(x_1, z_1, x_2, z_2, \dots, x_n, z_n | \lambda)$ .

c.

Now, suppose you are given only a sample  $\{x_1, \dots, x_n\}$  drawn i.i.d. from the mixture density, without the knowledge about which component the samples were drawn from (i.e., the  $z_i$  are unknown).

Using your derivations from part *a* and *b*, derive the E- and M-steps for an EM algorithm to compute the maximum likelihood estimate (MLE) of the mixture parameter  $\lambda$ .

## Solution

$$\text{a. } P(Z = 1|X = x) = \frac{P(X = x|Z = 1) \cdot P(Z = 1)}{P(X = x)} = \frac{\lambda f_1(x)}{f_\lambda(x)}.$$

$$\text{b. } P(x_1, z_1, x_2, z_2, \dots, x_n, z_n|\lambda) = \prod_{i=1}^n P(x_i, z_i|\lambda) = \prod_{i=1}^n (P(x_i|z_i, \lambda) \cdot P(z_i|\lambda)).$$

$$\begin{aligned} P(x_i|z_i, \lambda) \cdot P(z_i|\lambda) &= \begin{cases} \lambda f_1(x_i) & \text{if } z_i = 1 \\ (1 - \lambda) f_2(x_i) & \text{if } z_i = 0 \end{cases} \\ &= \lambda^{z_i} f_1(x_i)^{z_i} (1 - \lambda)^{1-z_i} f_2(x_i)^{1-z_i} \end{aligned}$$

Therefore,

$$\begin{aligned} \ln P(x_1, z_1, x_2, z_2, \dots, x_n, z_n|\lambda) &= \sum_{i=1}^n \ln P(x_i, z_i|\lambda) \\ &= \sum_{i=1}^n \ln (\lambda^{z_i} f_1(x_i)^{z_i} (1 - \lambda)^{1-z_i} f_2(x_i)^{1-z_i}) \\ &= \sum_{i=1}^n z_i ((\ln \lambda + \ln f_1(x_i))) + (1 - z_i) (\ln(1 - \lambda) + \ln f_2(x_i)) \end{aligned}$$



**c.**

**E-step:**  $q(z_i) \stackrel{not.}{=} P(z_i = 1|x_i, \lambda^{(t)}) \stackrel{B.Th.}{=} \frac{\lambda^{(t)} f_1(x_i)}{f_{\lambda^{(t)}}(x_i)}$

**M-step:**

$$\begin{aligned}\lambda^{(t+1)} &= \operatorname{argmax}_{\lambda} E_q \left[ \ln \prod_{i=1}^n P(x_1, z_1, x_2, z_2, \dots, x_n, z_n | \lambda) \right] \\ &= \operatorname{argmax}_{\lambda} \left( \ln \lambda \cdot \underbrace{\sum_{i=1}^n q(z_i)}_c + \ln(1 - \lambda) \cdot \underbrace{\sum_{i=1}^n (1 - q(z_i))}_{n-c} \right)\end{aligned}$$

$$\begin{aligned}\frac{\partial}{\partial \lambda} (c \ln \lambda + (n - c) \ln(1 - \lambda)) &= 0 \Leftrightarrow \frac{c}{\lambda} = \frac{n - c}{1 - \lambda} \\ \Leftrightarrow c(1 - \lambda) &= \lambda(n - c) \Leftrightarrow c = n\lambda \Leftrightarrow \lambda = \frac{c}{n} = \frac{\sum_{i=1}^n q(z_i)}{n} \\ \Rightarrow \lambda^{(t+1)} &= \frac{\lambda^{(t)}}{n} \sum_{i=1}^n \frac{f_1(x_i)}{f_{\lambda^{(t)}}(x_i)}\end{aligned}$$

**Note:**  $\frac{\partial^2}{\partial \lambda^2} (c \ln \lambda + (n - c) \ln(1 - \lambda)) = -\frac{c}{\lambda} - (n - c) \frac{1}{(1 - \lambda)^2} \leq 0 \Rightarrow \lambda^{(t+1)}$  corresponds to a maximum.

## Note

The EM algorithm can [be naturally extended so as to] handle mixtures of an arbitrary (although fixed) number of probabilistic distributions.

## Exemplification [for a mixture of 3 densities]

Source: Brani Vidakovic, *EM Algorithm and Mixtures* [Handout 12]

A sample of size 150 is generated from the mixture

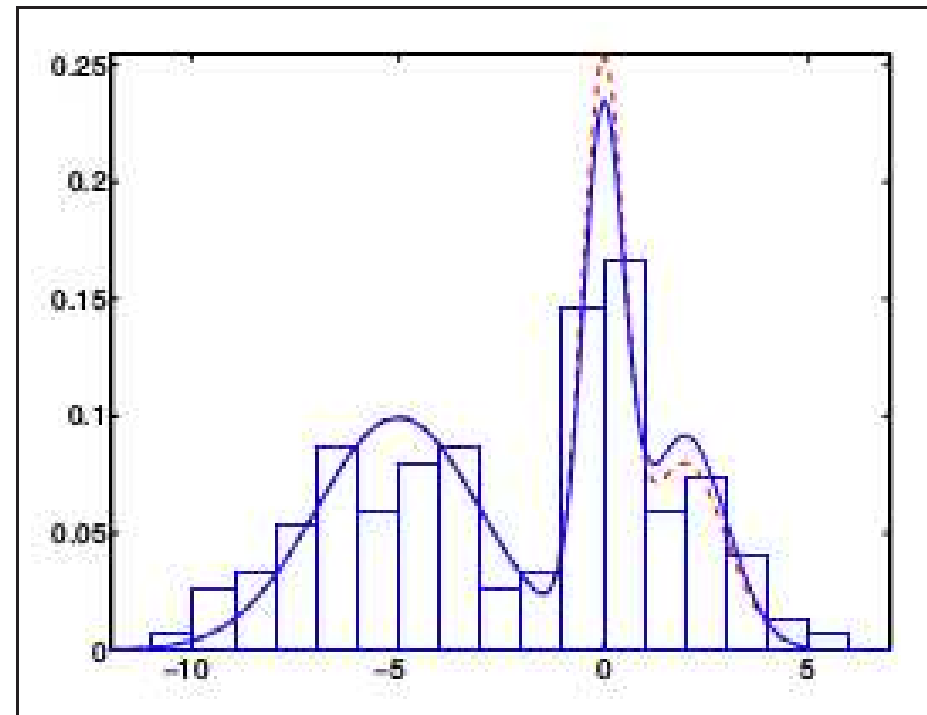
$$f(x) = 0.5\mathcal{N}(-5, 2^2) + 0.3\mathcal{N}(0, 0.5^2) + 0.2\mathcal{N}(2, 1),$$

where  $\mathcal{N}$  denotes the normal distribution.

The mixing weights are estimated by the EM algorithm.

$M = 20$  iterations of the EM algorithm yielded the weights (0.4977, 0.2732 and 0.2290).

The nearby figure shows the histogram of data, the theoretical mixture (dotted line) and the EM estimate (solid line).



Deriving the EM algorithm for  
estimating the parameters of two random variables  
of exponential distribution,  
given a set of instances generated by their sum

CMU, 2004 fall, T. Mitchell, Z. Bar-Joseph, HW2, pr. 2.2

Suppose that  $Z_1 \sim \exp(1/\lambda_1)$  and  $Z_2 \sim \exp(1/\lambda_2)$ , and  $\lambda_1 \neq \lambda_2$ .  $Z_1$  and  $Z_2$  are independent. Let  $X = Z_1 + Z_2$  denote the sum of  $Z_1$  and  $Z_2$ . Suppose that  $\{x_1, x_2, \dots, x_n\}$  are i.i.d. samples from the distribution of  $X$ .

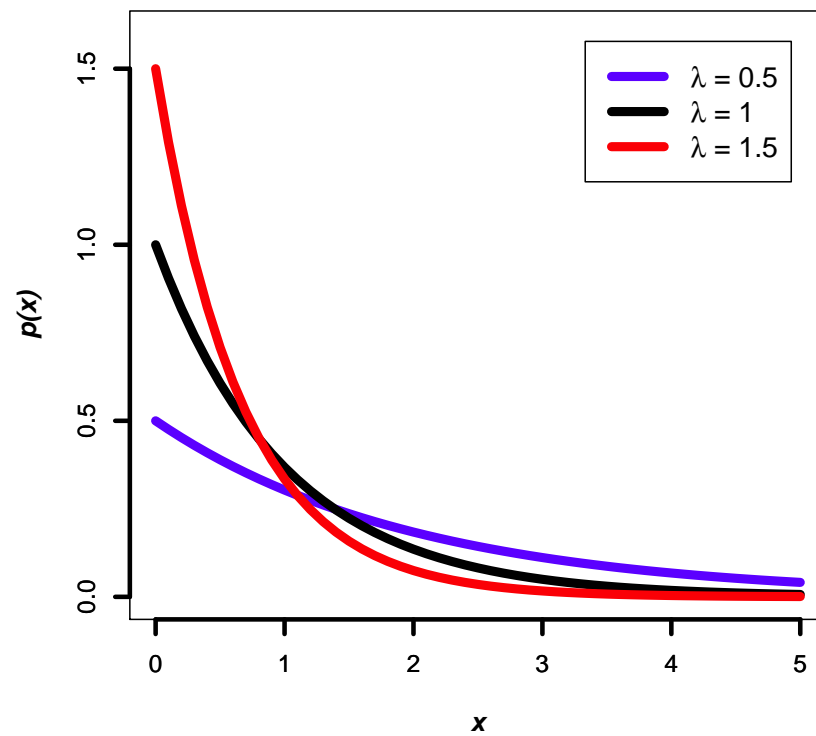
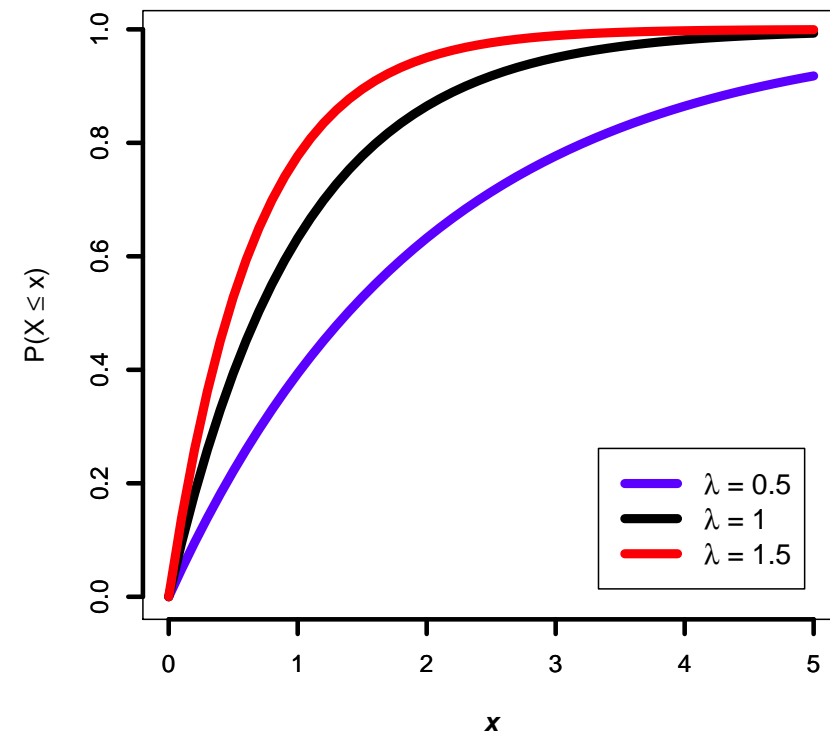
a. Derive an expression for the density of  $X$  in terms of  $\lambda_1$  and  $\lambda_2$ .

*Hint 1:* The density of  $Z_1$  is  $f_{\lambda_1}(z) = \lambda_1 \cdot e^{-\lambda_1 z}$  for  $z \geq 0$ , and 0 otherwise. The density of  $Z_2$  is defined in a similar way.

*Hint 2:* You could first derive c.d.f. (cumulative distribution function) of  $X$ , which is defined as  $F(x) = P(Z_1 + Z_2 < x) = \int_0^x \int_0^{x-z_1} f_{\lambda_1}(z_1) \cdot f_{\lambda_2}(z_2) dz_2 dz_1$ .

*Explanation:*

$$\begin{aligned}
 F(x) &\stackrel{\text{def.}}{=} P(Z_1 + Z_2 < x) \\
 &= \int_0^x \left( \int_0^{x-z_2} f_{\lambda_1}(z_1) dz_1 \right) \cdot f_{\lambda_2}(z_2) dz_2 \\
 &= \int_0^x \left( \int_0^{x-z_1} f_{\lambda_2}(z_2) dz_2 \right) \cdot f_{\lambda_1}(z_1) dz_1 = \int_0^x \int_0^{x-z_1} f_{\lambda_1}(z_1) \cdot f_{\lambda_2}(z_2) dz_2 dz_1.
 \end{aligned}$$

**Exponential probability density function****Exponential cumulative distribution function**

## Answer

The c.d.f. of  $X$ :

$$\begin{aligned}
 F(x) &\stackrel{\text{def.}}{=} P(Z_1 + Z_2 < x) = \int_0^x \int_0^{x-z_1} f_{\lambda_1}(z_1) \cdot f_{\lambda_2}(z_2) dz_2 dz_1 = \int_0^x \int_0^{x-z_1} \lambda_1 e^{-\lambda_1 z_1} \cdot \lambda_2 e^{-\lambda_2 z_2} dz_2 dz_1 \\
 &= \int_0^x (-\lambda_1) e^{-\lambda_1 z_1} \left( \int_0^{x-z_1} (-\lambda_2) e^{-\lambda_2 z_2} dz_2 \right) dz_1 = \int_0^x (-\lambda_1) e^{-\lambda_1 z_1} \left( e^{-\lambda_2 z_2} \Big|_0^{x-z_1} \right) dz_1 \\
 &= \int_0^x (-\lambda_1) e^{-\lambda_1 z_1} \left( e^{-\lambda_2(x-z_1)} - \underbrace{e^{-\lambda_2 \cdot 0}}_1 \right) dz_1 \\
 &= \int_0^x (-\lambda_1) e^{-\lambda_1 z_1} e^{-\lambda_2(x-z_1)} dz_1 - \int_0^x (-\lambda_1) e^{-\lambda_1 z_1} dz_1 \\
 &= \frac{-\lambda_1}{\lambda_2 - \lambda_1} e^{-\lambda_2 x} \int_0^x (\lambda_2 - \lambda_1) e^{(\lambda_2 - \lambda_1) z_1} dz_1 - \int_0^x (-\lambda_1) e^{-\lambda_1 z_1} dz_1 \\
 &= -\frac{\lambda_1}{\lambda_2 - \lambda_1} e^{-\lambda_2 x} \left( e^{(\lambda_2 - \lambda_1) z_1} \Big|_0^x \right) - e^{-\lambda_1 z_1} \Big|_0^x \\
 &= -\frac{\lambda_1}{\lambda_2 - \lambda_1} e^{-\lambda_2 x} \left( e^{(\lambda_2 - \lambda_1)x} - 1 \right) - \left( e^{-\lambda_1 x} - 1 \right) = -\frac{\lambda_1}{\lambda_2 - \lambda_1} e^{-\lambda_1 x} + \frac{\lambda_1}{\lambda_2 - \lambda_1} e^{-\lambda_2 x} - e^{-\lambda_1 x} + 1 \\
 &= 1 - \frac{1}{\lambda_2 - \lambda_1} (\lambda_1 e^{-\lambda_1 x} - \lambda_1 e^{-\lambda_2 x} + \lambda_2 e^{-\lambda_1 x} - \lambda_1 e^{-\lambda_1 x}) = 1 - \frac{\lambda_2 e^{-\lambda_1 x} - \lambda_1 e^{-\lambda_2 x}}{\lambda_2 - \lambda_1}
 \end{aligned}$$

The p.d.f. of  $X$ :

$$p(x) = \frac{\partial F(x)}{\partial x} = -\frac{1}{\lambda_2 - \lambda_1} (-\lambda_1 \lambda_2 e^{-\lambda_1 x} + \lambda_1 \lambda_2 e^{-\lambda_2 x}) = \frac{\lambda_1 \lambda_2}{\lambda_2 - \lambda_1} (e^{-\lambda_1 x} - e^{-\lambda_2 x})$$



b. Derive the E-step and M-step of the EM algorithm, and give explicit expressions for the parameter updates in the EM process for computing the MLE of  $\lambda_1$  and  $\lambda_2$ .

*Pasul E (Expectation):* Se calculează funcția „auxiliară“ (media log-verosimilității datelor complete) pentru iterația  $t$ :

$$Q(\lambda \mid \lambda^{(t)}) = E_{P(Z \mid X, \lambda^{(t)})}[\log P(X, Z \mid \lambda)]$$

*Pasul M (Maximization):* Se maximizează media log-verosimilității datelor complete, calculată la pasul E, în raport cu  $\lambda$ :

$$\lambda^{(t+1)} = \underset{\lambda}{\operatorname{argmax}} Q(\lambda \mid \lambda^{(t)})$$

*Notățiile* de mai sus au următoarele semnificații:

- $\lambda$  =  $(\lambda_1, \lambda_2)$
- $\lambda^{(t)}$  = valoarea parametrului  $\lambda$  la iterația  $t$
- $X$  = variabila observabilă, cu instanțele  $x_1, x_2, \dots, x_n$
- $Z$  =  $(Z_1, Z_2)$  variabilele ascunse/neobservabile  
având valorile  $z_{1j}, z_{2j}, \dots, z_{nj}$  cu  $j \in \{1, 2\}$ ,  
așa încât  $x_i = z_{i1} + z_{i2}$ .

**Solution: E step; computation of the “auxiliary” function**

$$\begin{aligned}
 Q(\lambda \mid \lambda^{(t)}) &= E_{P(Z|X, \lambda^{(t)})} \left[ \log \prod_{i=1}^n p(x_i, z_{i1}, z_{i2} \mid \lambda) \right] = E_{P(Z|X, \lambda^{(t)})} \left[ \log \prod_{i=1}^n f_{\lambda_1}(z_{i1}) \cdot f_{\lambda_2}(z_{i2}) \right] \\
 &= E_{P(Z|X, \lambda^{(t)})} \left[ \sum_{i=1}^n \log (\lambda_1 e^{-\lambda_1 z_{i1}} \cdot \lambda_2 e^{-\lambda_2 z_{i2}}) \right] \\
 &= E_{P(Z|X, \lambda^{(t)})} \left[ \sum_{i=1}^n (\log \lambda_1 - \lambda_1 z_{i1} + \log \lambda_2 - \lambda_2 z_{i2}) \right] \\
 &= E_{P(Z|X, \lambda^{(t)})} \left[ n \log \lambda_1 + n \log \lambda_2 - \sum_{i=1}^n (\lambda_1 z_{i1} + \lambda_2 z_{i2}) \right] \\
 &= n \log \lambda_1 + n \log \lambda_2 - \lambda_1 \sum_{i=1}^n E_{p(z_{i1}|x_i, \lambda^{(t)})}[z_{i1}] - \lambda_2 \sum_{i=1}^n E_{p(z_{i2}|x_i, \lambda^{(t)})}[z_{i2}]
 \end{aligned}$$

## Solution: E step; computation of expectations

$$\begin{aligned}
 E_{p(z_{i1}|x_i, \lambda^{(t)})}[z_{i1}] &\stackrel{\text{def.}}{=} \int_0^{x_i} z_{i1} \cdot p(z_{i1} | x_i, \lambda^{(t)}) dz_{i1} \\
 p(z_{i1} | x_i, \lambda^{(t)}) &\stackrel{\text{def.}}{=} \frac{p(z_{i1}, x_i | \lambda^{(t)})}{p(x_i | \lambda^{(t)})} \\
 &= \frac{p(z_{i1}, z_{i2} | \lambda^{(t)})}{p(x_i | \lambda^{(t)})} = \frac{p(z_{i1} | \lambda_1^{(t)}) \cdot p(z_{i2} | \lambda_2^{(t)})}{p(x_i | \lambda^{(t)})} \stackrel{a.}{=} \frac{f_{\lambda_1^{(t)}}(z_{i1}) \cdot f_{\lambda_2^{(t)}}(x_i - z_{i1})}{\frac{\lambda_1^{(t)} \lambda_2^{(t)}}{\lambda_2^{(t)} - \lambda_1^{(t)}} \left( e^{-\lambda_1^{(t)} x_i} - e^{-\lambda_2^{(t)} x_i} \right)} \\
 &= (\lambda_2^{(t)} - \lambda_1^{(t)}) \cdot \frac{\lambda_1^{(t)} e^{-\lambda_1^{(t)} z_{i1}} \cdot \lambda_2^{(t)} e^{-\lambda_2^{(t)} (x_i - z_{i1})}}{\lambda_1^{(t)} \lambda_2^{(t)} \left( e^{-\lambda_1^{(t)} x_i} - e^{-\lambda_2^{(t)} x_i} \right)} \\
 &= (\lambda_2^{(t)} - \lambda_1^{(t)}) \cdot \frac{e^{(\lambda_2^{(t)} - \lambda_1^{(t)}) z_{i1}}}{\left( e^{-\lambda_1^{(t)} x_i} - e^{-\lambda_2^{(t)} x_i} \right) \cdot e^{\lambda_2^{(t)} x_i}} \\
 &= (\lambda_2^{(t)} - \lambda_1^{(t)}) \cdot \frac{e^{(\lambda_2^{(t)} - \lambda_1^{(t)}) z_{i1}}}{e^{(\lambda_2^{(t)} - \lambda_1^{(t)}) x_i} - 1}
 \end{aligned}$$

$$\begin{aligned}
\Rightarrow E_{p(z_{i1}|x_i,\lambda^{(t)})}[z_{i1}] &\stackrel{\text{def.}}{=} \int_0^{x_i} z_{i1} \cdot p(z_{i1} | x_i, \lambda^{(t)}) dz_{i1} = \int_0^{x_i} z_{i1} \cdot (\lambda_2^{(t)} - \lambda_1^{(t)}) \cdot \frac{e^{(\lambda_2^{(t)} - \lambda_1^{(t)})z_{i1}}}{e^{(\lambda_2^{(t)} - \lambda_1^{(t)})x_i} - 1} dz_{i1} \\
&= \frac{1}{e^{(\lambda_2^{(t)} - \lambda_1^{(t)})x_i} - 1} \int_0^{x_i} z_{i1} \cdot (\lambda_2^{(t)} - \lambda_1^{(t)}) \cdot e^{(\lambda_2^{(t)} - \lambda_1^{(t)})z_{i1}} dz_{i1}
\end{aligned}$$

Vom rezolva ultima integrală de mai sus utilizând formula de integrare prin părți.

$$\begin{aligned}
\int_0^{x_i} z_{i1} \cdot (\lambda_2^{(t)} - \lambda_1^{(t)}) \cdot e^{(\lambda_2^{(t)} - \lambda_1^{(t)})z_{i1}} dz_{i1} &= \int_0^{x_i} z_{i1} \cdot \frac{\partial}{\partial z_{i1}} \left( e^{(\lambda_2^{(t)} - \lambda_1^{(t)})z_{i1}} \right) dz_{i1} \\
&= \left( z_{i1} \cdot e^{(\lambda_2^{(t)} - \lambda_1^{(t)})z_{i1}} \right) \Big|_0^{x_i} - \int_0^{x_i} e^{(\lambda_2^{(t)} - \lambda_1^{(t)})z_{i1}} dz_{i1} \\
&= \left( x_i \cdot e^{(\lambda_2^{(t)} - \lambda_1^{(t)})x_i} - 0 \right) - \frac{1}{\lambda_2^{(t)} - \lambda_1^{(t)}} \cdot e^{(\lambda_2^{(t)} - \lambda_1^{(t)})z_{i1}} \Big|_0^{x_i} = x_i \cdot e^{(\lambda_2^{(t)} - \lambda_1^{(t)})x_i} - \frac{e^{(\lambda_2^{(t)} - \lambda_1^{(t)})x_i} - 1}{\lambda_2^{(t)} - \lambda_1^{(t)}}
\end{aligned}$$

Prin urmare,

$$\begin{aligned}
E_{p(z_{i1}|x_i,\lambda^{(t)})}[z_{i1}] &= \frac{1}{e^{(\lambda_2^{(t)} - \lambda_1^{(t)})x_i} - 1} \cdot \left( x_i \cdot e^{(\lambda_2^{(t)} - \lambda_1^{(t)})x_i} - \frac{e^{(\lambda_2^{(t)} - \lambda_1^{(t)})x_i} - 1}{\lambda_2^{(t)} - \lambda_1^{(t)}} \right) \\
&= \frac{x_i \cdot e^{(\lambda_2^{(t)} - \lambda_1^{(t)})x_i}}{e^{(\lambda_2^{(t)} - \lambda_1^{(t)})x_i} - 1} - \frac{1}{\lambda_2^{(t)} - \lambda_1^{(t)}}
\end{aligned}$$

$$\begin{aligned}
x_i &= z_{i1} + z_{i2} \\
\Rightarrow E_{p(z_{i2}|x_i,\lambda^{(t)})}[z_{i2}] &= E_{p(z_{i2}|x_i,\lambda^{(t)})}[x_i - z_{i1}] = x_i - E_{p(z_{i2}|x_i,\lambda^{(t)})}[z_{i1}] \\
&= x_i - E_{p(z_{i1}|x_i,\lambda^{(t)})}[z_{i1}] = x_i - E_{p(z_{i1}|x_i,\lambda^{(t)})}[z_{i1}] \\
&= x_i - \frac{x_i \cdot e^{(\lambda_2^{(t)} - \lambda_1^{(t)})x_i}}{e^{(\lambda_2^{(t)} - \lambda_1^{(t)})x_i} - 1} + \frac{1}{\lambda_2^{(t)} - \lambda_1^{(t)}} \\
&= \frac{1}{\lambda_2^{(t)} - \lambda_1^{(t)}} - \frac{x_i}{e^{(\lambda_2^{(t)} - \lambda_1^{(t)})x_i} - 1}
\end{aligned}$$

## Solution: M step

$$\begin{aligned}\lambda^{(t+1)} &= \operatorname{argmax}_{\lambda} Q(\lambda \mid \lambda^{(t)}) \\ &= \operatorname{argmax}_{\lambda_1 > 0, \lambda_2 > 0} \left( n \log \lambda_1 + n \log \lambda_2 - \lambda_1 \sum_{i=1}^n E_{p(z_{i1} \mid x_i, \lambda^{(t)})}[z_{i1}] - \lambda_2 \sum_{i=1}^n E_{p(z_{i2} \mid x_i, \lambda^{(t)})}[z_{i2}] \right)\end{aligned}$$

$$\begin{aligned}\frac{\partial}{\partial \lambda_1} Q(\lambda \mid \lambda^{(t)}) = 0 &\Leftrightarrow \frac{n}{\lambda_1} = \sum_{i=1}^n E_{p(z_{i1} \mid x_i, \lambda^{(t)})}[z_{i1}] \\ \Rightarrow \lambda_1^{(t+1)} &= \frac{n}{\sum_{i=1}^n E_{p(z_{i1} \mid x_i, \lambda^{(t)})}[z_{i1}]} > 0\end{aligned}$$

$$\begin{aligned}\frac{\partial}{\partial \lambda_2} Q(\lambda \mid \lambda^{(t)}) = 0 &\Leftrightarrow \frac{n}{\lambda_2} = \sum_{i=1}^n E_{p(z_{i2} \mid x_i, \lambda^{(t)})}[z_{i2}] \\ \Rightarrow \lambda_2^{(t+1)} &= \frac{n}{\sum_{i=1}^n E_{p(z_{i2} \mid x_i, \lambda^{(t)})}[z_{i2}]} > 0\end{aligned}$$

**Observație:**

Se verifică imediat că aceste soluții într-adevăr maximizează  $Q(\lambda \mid \lambda^{(t)})$ .

Using the EM algorithm for  
learning a Poisson distribution  
*when some data are missing*

Liviu Ciortuz, following

Anirban DasGupta, *Probability for Statistics and Machine Learning*, Springer, 2011, Ex. 20.8

Presupunem că vrem să modelăm statistic numărul de accidente ușoare care s-au produs în  $n$  locații într-un anumit interval de timp, să zicem o săptămână. În acest scop, vom folosi o distribuție Poisson de parametru  $\lambda$ . La sârșitul perioadei de timp respective, ni se transmite de la  $m$  din cele  $n$  locații câte o „înregistrare” (notată cu  $x_i$ ), reprezentând numărul de accidente ușoare produse în locația  $i$ .

În mod implicit, ar trebui să considerăm că în cele  $n - m$  locații de la care n-am primit înregistrări nu s-a produs niciun accident. Însă, în urma „inspectării” datelor suntem determinați să luăm în considerare *presupunerea* că în unele din aceste  $k$  locații s-a produs câte un [singur] accident ușor, care a fost „trecut sub tăcere” la raportare.

Așadar, formalizând, vom considera datele „neobservabile”  $z_1 = 0, \dots, z_{n_0} = 0, z_{n_0+1} = 1, \dots, z_{n_0+n_1} = 1$ , cu  $n_0 + n_1 = n - m$ , alături de datele observabile  $x_1, \dots, x_m \in \mathbb{N}$ . Toate aceste date sunt produse de variabile aleatoare urmând distribuția *Poisson* de [același] parametru  $\lambda$ .

Elaborați algoritmul EM (în speță pasul E și pasul M) pentru estimarea parametrului  $\lambda$ .

*Sugestie:* În loc să se lucreze cu  $z_j$ , cu  $j = 1, \dots, m$ , va fi suficient să considerați ca date neobservabile  $n_0$  și  $n_1$ . Mai mult, considerând numerele  $n$  și  $m$  cunoscute, va fi suficient să lucrați doar cu  $n_1$  ca dată neobservabilă (bineînțeles, pe lângă datele observabile  $x_i$ ).



## Answer

$$z_1, \dots, z_{n+0}, z_{n_0+1}, \dots, z_{n_0+n_1}, x_1, \dots, x_m \sim \text{Poisson}(\lambda)$$

**Poisson p.m.f.:**  $P(x|\lambda) = \frac{1}{e^\lambda} \cdot \frac{\lambda^x}{x!}$

**The verosimilarity of complete data:**

$$\begin{aligned} L(\lambda) &\stackrel{\text{def.}}{=} P(z_1, \dots, z_{n+0}, z_{n_0+1}, \dots, z_{n_0+n_1}, x_1, \dots, x_m | \lambda) \\ &\stackrel{i.i.d.}{=} \prod_{j=1}^{n_0+n_1} P(z_j | \lambda) \cdot \prod_{i=1}^m P(x_i | \lambda) = \prod_{j=1}^{n_0} \frac{1}{e^\lambda} \frac{\lambda^0}{0!} \cdot \prod_{j=n_0+1}^{n_1} \frac{1}{e^\lambda} \frac{\lambda}{1!} \cdot \prod_{i=1}^m \frac{1}{e^\lambda} \frac{\lambda^{x_i}}{x_i!} \\ &= \frac{1}{(e^\lambda)^{n_0+n_1+m}} \cdot \lambda^{n_1} \cdot \prod_{i=1}^m \frac{\lambda^{x_i}}{x_i!} = \frac{1}{(e^\lambda)^n} \cdot \lambda^{n_1} \cdot \prod_{i=1}^m \frac{\lambda^{x_i}}{x_i!} \\ &= e^{-n\lambda} \cdot \lambda^{n_1} \cdot \lambda^{\sum_{i=1}^m x_i} \cdot \frac{1}{\prod_{i=1}^m x_i!} = e^{-n\lambda} \cdot \lambda^{n_1 + \sum_{i=1}^m x_i} \cdot \frac{1}{\prod_{i=1}^m x_i!} \end{aligned}$$

**The log-verosimilarity of complete data:**

$$\ell(\lambda) \stackrel{\text{def.}}{=} \ln L(\lambda) = -n\lambda + \left( n_1 + \sum_{i=1}^m x_i \right) \ln \lambda - \sum_{i=1}^m \ln x_i!$$

The auxiliary function:

$$Q(\lambda|\lambda^{(t)}) \stackrel{def.}{=} E_{n_1|x_i, \lambda^{(t)}}[\ell(\lambda)] = -n\lambda + \left( E[n_1|x_i, \lambda^{(t)}] + \sum_{i=1}^m x_i \right) \ln \lambda - \sum_{i=1}^m \ln x_i!.$$

**E-step:**

$E_{n_1|x_i, \lambda^{(t)}}[\ell(\lambda)]$  is the expected number of unobservable instances  $z_j$  having value 1, out of the total of  $n - m$  unobservable instances.

The definition of the density function for the Poisson distribution implies

$$P(x = 1|\lambda) = \frac{1}{e^\lambda} \cdot \frac{\lambda^1}{1!} = \frac{1}{e^\lambda} \cdot \lambda \text{ and } P(x = 0|\lambda) = \frac{1}{e^\lambda} \cdot \frac{\lambda^0}{0!} = \frac{1}{e^\lambda}.$$

Therefore  $E[n_1|x_i, \lambda^{(t)}]$  is exactly the expectation of the binomial distribution

$$\text{Bin}\left(n - m; \frac{\lambda^{(t)}}{1 + \lambda^{(t)}}\right). \text{ So,}$$

$$E[n_1|x_i, \lambda^{(t)}] = (n - m) \frac{\lambda^{(t)}}{1 + \lambda^{(t)}}.$$

## M-step:

Now we can write

$$Q(\lambda|\lambda^{(t)}) = -n\lambda + \left[ (n-m)\frac{\lambda^{(t)}}{1+\lambda^{(t)}} + \sum_{i=1}^m x_i \right] \ln \lambda - \sum_{i=1}^m \ln x_i!$$

The first two derivatives of this function w.r.t.  $\lambda$  are:

$$\frac{\partial}{\partial \lambda} Q(\lambda|\lambda^{(t)}) = -n + \left[ (n-m)\frac{\lambda^{(t)}}{1+\lambda^{(t)}} + \sum_{i=1}^m x_i \right] \frac{1}{\lambda}$$

$$\frac{\partial^2}{\partial \lambda^2} Q(\lambda|\lambda^{(t)}) = - \left[ (n-m)\frac{\lambda^{(t)}}{1+\lambda^{(t)}} + \sum_{i=1}^m x_i \right] \frac{1}{\lambda^2}$$

From the problem's statement we know that  $n > m$  and  $\sum_{i=1}^m x_i \geq 0$ .

At the initialization step of the EM algorithm the parameter  $\lambda$  is assigned a positive value ( $\lambda^{(0)}$ ), and we will soon see that  $\lambda^{(t)} > 0$  at each iteration  $t > 0$ .

Therefore, the second derivative of  $Q$  is always negative, meaning that  $Q$  is concave and therefore it has a maximum.

In order to find this maximum we equate the first derivative to 0, which gives us the *update equation*:

$$\lambda^{(t+1)} = \frac{1}{n} \left[ (n - m) \frac{\lambda^{(t)}}{1 + \lambda^{(t)}} + \sum_{i=1}^m x_i \right]$$

One can easily see now (by induction) that  $\lambda^{(t+1)} > 0$ .