

EM Algorithm - Exponential Families

Exponential Family $\exp(A(\theta) + h(x) + t(\theta)s(x))$

normalizing constant $\int \exp(A(\theta) + h(x) + t(\theta)s(x)) dx$

t suff. stat.

t corresponding parameter for suff. stat.

- The computation of these two steps simplify a great deal when it can be shown that the log-likelihood is linear in the sufficient statistic for θ (Exponential Family). *For EM algorithm, if your dist belongs to exponential family derivation of EM algorithm is much easier.*
- The E-step reduces to computing the expectation of the complete-data sufficient statistic given the observed data. *E-Step : Calculate the expectation of complete data suff. stat given the observed data*
- In M-step, the conditional expectations of the sufficient statistics computed in the E-step can be directly substituted for the sufficient statistics that occur in the expressions obtained for the complete-data maximum likelihood estimators of θ , to obtain the next iterate.

Exponential family : MLE is a function of suff. stat.

We plug-in the result of E-step into suff. stat.

If density cannot find MLE analytically, we have to use optimization
 \Rightarrow density belongs to exponential family. (we don't need)

Example: Univariate Normal Distribution

Complete data $y = (y_{\text{obs}}, y_{\text{miss}})$ $y \sim N(\mu, \sigma^2)$

$$L(\mu, \sigma^2) = \left(\frac{1}{2\pi\sigma^2} \right)^{n/2} \exp \left(-\frac{1}{2\sigma^2} \left(\sum y_i^2 - 2\mu \sum y_i + n\mu^2 \right) \right)$$

Let the complete-data vector $y = (y_1, \dots, y_n)$ be a random sample from $N(\mu, \sigma^2)$. Then,

$$f(y|\mu, \sigma^2) = \left(\frac{1}{2\pi\sigma^2} \right)^{n/2} \exp \left\{ -\frac{1}{2\sigma^2} \left(\sum y_i^2 - 2\mu \sum y_i + n\mu^2 \right) \right\}.$$

which implies that $(\sum y_i, \sum y_i^2)$ are sufficient statistics for $\theta = (\mu, \sigma^2)$. The complete-data log-likelihood function is

$$l(\mu, \sigma^2) = -\frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum y_i^2 + \frac{\mu}{\sigma^2} \sum y_i - \frac{n\mu^2}{\sigma^2} + \text{constant.}$$

Example: Univariate Normal Distribution

$y_i \quad i = 1, \dots, n. \quad (1, \dots, m) \quad (m+1, \dots, n) \quad m < n$

\hookrightarrow observed \hookrightarrow missing.

exponential family

Suppose $y_i, i = 1, \dots, m$ are observed and $y_i, i = m+1, \dots, n$ are missing.

Denote the observed data by $y_{obs} = (y_1, \dots, y_m)$.

$\sum y_i, \sum y_i^2$ suff. stat.

- The E-step requires the computation of

$$\text{Var } X = EX^2 - (EX)^2$$

$$EX^2 = \underbrace{\text{Var } X}_{\text{sample size}} + (EX)^2. \quad \hookrightarrow (\mu^{(t)})^2 E_\theta \left(\sum y_i | y_{obs} \right) \quad \text{and} \quad E_\theta \left(\sum y_i^2 | y_{obs} \right).$$

est. : $\mu^{(t)}$
sample size : $n-m$

- At the t -th iteration of the E-step, compute

$$E \left(\sum_{i=1}^n y_i | y_{obs} \right) = E \left(\sum_{i=1}^m y_i + \sum_{i=m+1}^n y_i | y_{obs} \right) = E \left(\sum_{i=1}^m y_i | y_{obs} \right) + E \left(\sum_{i=m+1}^n y_i | y_{obs} \right)$$

$$n-m$$

$$s_1^{(t)} = E_\theta \left(\sum y_i | y_{obs} \right) = \sum_{i=1}^m y_i + (n-m)\mu^{(t)} = \sum_{i=1}^m y_i + (n-m)\mu^{(t)}$$

\hookrightarrow involve $\mu^{(t)}$

$$\begin{pmatrix} \mu^{(t)} \\ \sigma^{2(t)} \end{pmatrix} \leftarrow s_1^{(t)} = E_\theta \left(\sum y_i^2 | y_{obs} \right) = \sum_{i=1}^m y_i^2 + (n-m) [\sigma^{2(t)} + \mu^{(t)2}]$$

$$E_\theta \left(\sum y_i^2 | y_{obs} \right) = E \left(\sum_{i=1}^m y_i^2 + \sum_{i=m+1}^n y_i^2 | y_{obs} \right) = \sum_{i=1}^m y_i^2 + (n-m) \left(\mu^{(t)2} + \sigma^{2(t)} \right)$$

$$= E \left(\sum_{i=1}^m y_i^2 | y_{obs} \right) + E \left(\sum_{i=m+1}^n y_i^2 | y_{obs} \right)$$

Example: Univariate Normal Distribution

M-Step : Normal ~ Exponential family.

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n y_i \quad \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n y_i^2 - \left(\frac{1}{n} \sum_{i=1}^n y_i \right)^2 = \frac{1}{n} \sum_{i=1}^n y_i^2 - \hat{\mu}^2.$$

- The complete-data maximum likelihood estimation of μ and σ^2 are

$$\hat{\mu} = \frac{\sum_{i=1}^n y_i}{n} \quad \text{and} \quad \hat{\sigma}^2 = \frac{\sum_{i=1}^n y_i^2}{n} - \left(\frac{\sum_{i=1}^n y_i}{n} \right)^2.$$

- At the t -th iteration of the M-step, compute

$$\mu^{(t+1)} = \frac{s_1^{(t)}}{n}$$

$$\mu^{(t+1)} = \frac{1}{n} s_1^{(t)}.$$

$$\sigma^{2(t+1)} = \frac{s_2^{(t)}}{n} - \mu^{(t+1)2}.$$

$$\sigma^{2(t+1)} = \frac{1}{n} s_2^{(t)} - (\mu^{(t+1)})^2.$$

Example: Multinomial Distribution

EM (1979)

$$y = (y_1, y_2, y_3, y_4)$$

$$(p_1, p_2, p_3, p_4) = \left(\frac{1}{2} - \frac{1}{2}\theta, \frac{1}{4}\theta, \frac{1}{4}\theta, \frac{1}{2}\right)$$

We want to estimate θ .

Complete-data log likelihood

not completely observed.

- Define $y = (y_1, y_2, y_3, y_4)$ with multinomial probability $(\frac{1}{2} - \frac{1}{2}\theta, \frac{1}{4}\theta, \frac{1}{4}\theta, \frac{1}{2}) = (p_1, p_2, p_3, p_4)$. What we observed, $(y_1, y_2, y_3 + y_4) = (38, 34, 125)$
- Let $y_{obs} = (y_1, y_2, y_3 + y_4) = (38, 34, 125)$ be observed counts from a multinomial population, and then, observed data vector is a function of the complete-data vector. genotype phenotype
- Since only $y_3 + y_4$ is observed and y_3 and y_4 are not, the observed data is considered incomplete. $\ell(\theta) = y_1 \log p_1 + y_2 \log p_2 + y_3 \log p_3 + y_4 \log p_4 + \text{constant}$
- The complete-data log-likelihood is $P_1 + P_2 + P_3 + P_4 = 1$

$$\ell(\theta) = y_1 \log p_1 + y_2 \log p_2 + y_3 \log p_3 + y_4 \log p_4 + \text{constant}$$

$$= y_1 \log \left(\frac{1}{2} - \frac{1}{2}\theta \right) + y_2 \log \frac{1}{4}\theta + y_3 \log \frac{1}{4}\theta + y_4 \log \frac{1}{2} + \text{constant}$$

where y_1, y_2, y_3, y_4 are sufficient statistics.

$$\frac{\partial \ell(\theta)}{\partial \theta} = -\frac{y_1}{1-\theta} + \frac{y_2}{\theta} + \frac{y_3}{\theta} = 0 \Rightarrow \frac{y_1}{1-\theta} = \frac{y_2 + y_3}{\theta} \Rightarrow (y_2 + y_3)(1-\theta) = y_1 \theta \\ (y_1 + y_2 + y_3) \theta = y_2 + y_3$$

Example: Multinomial Distribution

↳ belongs to exponential family.

- E-step:
 $E(y_1|y_{obs}) = y_1 = 38$
 $E(y_2|y_{obs}) = y_2 = 34$

$$E_\theta(y_1|y_{obs}) = y_1 = 38 \quad \text{and} \quad E_\theta(y_2|y_{obs}) = y_2 = 34,$$

$$E_\theta(y_3|y_{obs}) = E_\theta(y_3|y_3 + y_4) = 125 \frac{\frac{1}{4}\theta}{\frac{1}{2} + \frac{1}{4}\theta}, \quad y_3^{(t)} = 125 \frac{\frac{1}{4}\theta^{(t)}}{\frac{1}{2} + \frac{1}{4}\theta^{(t)}}$$

- M-step: The complete-data maximum likelihood estimate of θ is

$$\hat{\theta} = \frac{y_2 + y_3}{y_1 + y_2 + y_3} = \frac{34 + 125 \cdot \frac{\frac{1}{4}\theta^{(t)}}{\frac{1}{2} + \frac{1}{4}\theta^{(t)}}}{72 + 125 \frac{\frac{1}{4}\theta^{(t)}}{\frac{1}{2} + \frac{1}{4}\theta^{(t)}}}$$

Then,

$$\theta^{(t+1)} = \frac{34 + y_3^{(t)}}{72 + y_3^{(t)}} \quad (\text{HW})$$

Example: Peppered Moths

C $\xrightarrow{\text{dominant}}$ I $\xrightarrow{\text{dominant}}$ T $3P_2 = 6$

- Three possible alleles: C, I, and T.
- C is dominant to I and T is recessive to I
 - Genotype CC, CI, and CT \Rightarrow carbonaria phenotype. $n_C \dots p_C$
 - Genotype TT \Rightarrow typica phenotype. $n_T \dots p_T$
 - Genotype II and IT \Rightarrow intermediate phenotype called insularia.
- If the allele frequencies in the population are p_C , p_I , and p_T , then the genotype frequencies should be p_C^2 , $2p_C p_I$, $2p_C p_T$, p_I^2 , $2p_I p_T$, and p_T^2 for genotypes CC, CI, CT, II, IT, and TT, respectively. Note that $p_C + p_I + p_T = 1$.
- Capture n moths, of which there are n_C , n_I , and n_T of the carbonaria, insularia, and typica phenotypes, respectively. Thus, $n = n_C + n_I + n_T$.

Example: Peppered Moths

- The observed data are $x = (n_C, n_I, n_T)$ and the complete data are $y = (n_{CC}, n_{CI}, n_{CT}, n_{II}, n_{IT}, n_{TT})$. The mapping from the complete data to the observed data is $x = M(y) = (n_{CC} + n_{CI} + n_{CT}, n_{II} + n_{IT}, n_{TT})$.
- Wish to estimate the allele probabilities, p_C , p_I , and p_T .
- The complete data log likelihood function is multinomial

$$\begin{aligned}\log L(p|y) &= n_{CC} \log \left\{ p_C^2 \right\} + n_{CI} \log \{2p_C p_I\} + n_{CT} \log \{2p_C p_T\} \\ &\quad + n_{II} \log \left\{ p_I^2 \right\} + n_{IT} \log \{2p_I p_T\} + n_{TT} \log \left\{ p_T^2 \right\} \\ &\quad + \log \binom{n}{n_{CC} \ n_{CI} \ n_{CT} \ n_{II} \ n_{IT} \ n_{TT}}.\end{aligned}$$

Example: Peppered Moths

- E-Step

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

$$E \left\{ N_{CI} | n_C, n_I, n_T, p^{(t)} \right\} = n_{CI}^{(t)} = \frac{2n_C p_C^{(t)} p_I^{(t)}}{\left(p_C^{(t)} \right)^2 + 2p_C^{(t)} p_I^{(t)} + 2p_C^{(t)} p_T^{(t)}}$$

$$E \left\{ N_{CT} | n_C, n_I, n_T, p^{(t)} \right\} = n_{CT}^{(t)} = \frac{2n_C p_C^{(t)} p_T^{(t)}}{\left(p_C^{(t)} \right)^2 + 2p_C^{(t)} p_I^{(t)} + 2p_C^{(t)} p_T^{(t)}}$$

Example: Peppered Moths

- E-Step

$$E \left\{ N_{II} | n_C, n_I, n_T, p^{(t)} \right\} = n_{II}^{(t)} = \frac{n_I \left(p_I^{(t)} \right)^2}{\left(p_I^{(t)} \right)^2 + 2p_I^{(t)}p_T^{(t)}}$$

$$E \left\{ N_{IT} | n_C, n_I, n_T, p^{(t)} \right\} = n_{IT}^{(t)} = \frac{2n_I p_I^{(t)} p_T^{(t)}}{\left(p_I^{(t)} \right)^2 + 2p_I^{(t)}p_T^{(t)}}$$

Then,

$$\begin{aligned} Q(p|p^{(t)}) &= n_{CC}^{(t)} \log\{p_C^2\} + n_{CI}^{(t)} \log\{2p_C p_I\} + n_{CT}^{(t)} \log\{2p_C p_T\} \\ &\quad + n_{II}^{(t)} \log\{p_I^2\} + n_{IT}^{(t)} \log\{2p_I p_T\} + n_{TT} \log\{p_T^2\} + \kappa(n_C, n_I, n_T, p^{(t)}). \end{aligned}$$

Example: Peppered Moths

- M-Step: Recall $p_T = 1 - p_C - p_I$ and differentiate w.r.t p_C and p_I .

$$\frac{dQ(p|p^{(t)})}{dp_C} = \frac{2n_{CC}^{(t)} + n_{CI}^{(t)} + n_{CT}^{(t)}}{p_C} + \frac{2n_{TT}^{(t)} + n_{CT}^{(t)} + n_{IT}^{(t)}}{1 - p_C - p_I}$$
$$\frac{dQ(p|p^{(t)})}{dp_I} = \frac{2n_{II}^{(t)} + n_{IT}^{(t)} + n_{CI}^{(t)}}{p_C} + \frac{2n_{TT}^{(t)} + n_{CT}^{(t)} + n_{IT}^{(t)}}{1 - p_C - p_I}$$

Then,

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

- Suppose the observed phenotype counts are $n_C = 85$, $n_I = 196$, and $n_T = 341$.

Observed (n_c, n_I, n_T)

we want to estimate (P_c, P_I, P_T)

Direct estimation of P_c, P_I, P_T is not possible.

$$\Rightarrow C = (CC, CI, CT) \quad I = (II, IT)$$

Complete data: \sim Multinomial dist $(P_{CC}, P_{CI}, P_{CT}, P_{II}, P_{IT}, P_{TT})$

$$y = (n_{CC}, n_{CI}, n_{CT}, n_{II}, n_{IT}, n_{TT})$$

Observed data

$$x = (n_{CC} + n_{CI} + n_{CT}, n_{II} + n_{IT}, n_{TT})$$

:

P_c

:

P_I

:

P_T

$$(P_c^2, 2P_cP_I, 2P_cP_T, P_I^2, 2P_I P_T, P_T^2)$$

Complete-data log likelihood

$$\begin{aligned} l(p) = & n_{CC} \log P_c^2 + n_{CI} \log 2P_cP_I + n_{CT} \log 2P_cP_T \\ & + n_{II} \log P_I^2 + n_{IT} \log 2P_I P_T + n_{TT} \log P_T^2 \\ & + \text{Constant.} \end{aligned}$$

$$\text{E-Step } E\{N_{CC} | n_c, n_I, n_T, p^{(t)}\} = n_{CC}^{(t)} = \frac{n_c}{P_c^{(t)2} + 2P_c^{(t)}P_I^{(t)} + 2P_c^{(t)}P_T^{(t)}}$$

$$E\{N_{CI} | n_c, n_I, n_T, p^{(t)}\} = n_{CI}^{(t)} = n_c \frac{2P_c^{(t)}P_I^{(t)}}{P_c^{(t)2} + 2P_c^{(t)}P_I^{(t)} + 2P_c^{(t)}P_T^{(t)}}$$

$$E\{N_{CT} | n_c, n_I, n_T, p^{(t)}\} = n_{CT}^{(t)} = n_c \frac{2P_c^{(t)}P_T^{(t)}}{P_c^{(t)2} + 2P_c^{(t)}P_I^{(t)} + 2P_c^{(t)}P_T^{(t)}}$$

$$E\{N_{II} | n_c, n_I, n_T, P^{(t)}\} = n_{II}^{(t)} = n_I \times \frac{P_I^{(t)2}}{P_I^{(t)2} + 2P_I^{(t)}P_T^{(t)}}$$

$$E\{N_{IT} | n_c, n_I, n_T, P^{(t)}\} = n_{IT}^{(t)} = n_I \times \frac{2P_I^{(t)}P_T^{(t)}}{P_I^{(t)2} + 2P_I^{(t)}P_T^{(t)}}$$

$$E\{N_{II} | n_c, n_I, n_T, P^{(t)}\} = ? \times .$$

C → I → T

estimation: P

$$Q(P | P^{(t)}) = n_{cc}^{(t)} \log P_c^2 + n_{cI}^{(t)} \log 2P_cP_I + n_{cT}^{(t)} \log 2P_cP_T$$

Expectation $P^{(t)}$

use observed value.

$$+ n_{II}^{(t)} \log P_I^2 + n_{IT}^{(t)} \log 2P_I P_T + n_{TT} \log P_T^2 + \text{constant.}$$

$$P_T = (1 - P_c - P_I)$$

M-Step

$$\frac{dQ(P | P^{(t)})}{dP_c} = \frac{2n_{cc}^{(t)} + n_{cI}^{(t)} + n_{cT}^{(t)}}{P_c} + \frac{2n_T + n_{cT}^{(t)} + n_{IT}^{(t)}}{1 - P_c - P_I} = 0$$

$$\frac{dQ(P | P^{(t)})}{dP_I} = \frac{2n_{II}^{(t)} + n_{cI}^{(t)} + n_{IT}^{(t)}}{P_I} + \frac{2n_T + n_{cT}^{(t)} + n_{IT}^{(t)}}{1 - P_c - P_I} = 0$$

$$P_c^{(t+1)} = \frac{2n_{cc}^{(t)} + n_{cI}^{(t)} + n_{cT}^{(t)}}{2n.}$$

$$P_I^{(t+1)} = \frac{2n_{II}^{(t)} + n_{cI}^{(t)} + n_{IT}^{(t)}}{2n.}$$

$$(Ex) \quad \begin{aligned} n_c &= 85 \\ n_I &= 196 \\ n_T &= 341. \end{aligned} \quad \Rightarrow \text{Implement EM algorithm.} \quad (HW)$$

Due : Nov. 11

Class Material

- Optimization
- Combinatorial
- EM.
- Numerical Integration
- Bootstrap.

(Ex) Normal Mixture Model



(J Group)

$$x_1, \dots, x_n \sim \sum_{j=1}^J p_j N(\mu_j, \sigma_j^2) \quad \dots \text{ (Observed Data)}$$

↗ prop. of distn.
 ↘ mean ↘ variance.

Define $(y_{11}, \dots, y_{1J}) \sim \text{Multinomial}(1; p_1, \dots, p_J)$ where $\sum_j y_{1j} = 1$.

↳ Indicator variable.
(which distribution x_i belongs to)

Given $y_{1j^*} = 1$ and $y_{1j} = 0$ for $j \neq j^*$,

we assume $x_i \sim N(\mu_{j^*}, \sigma_{j^*}^2)$

Complete data : $\{x_i, y_{11}, \dots, y_{1J}\}_i$

Observed-data log-likelihood

$$l(\mu, \sigma^2, p | x) = \sum_i \log \left\{ \sum_{j=1}^J p_j N(x_i; \mu_j, \sigma_j^2) \right\}$$

$$L(x_i; y_{11}, \dots, y_{1J}) = \prod_{j=1}^J p_j^{y_{1j}} N(x_i; \mu_j, \sigma_j^2)^{y_{1j}}$$

$$L(x; y_1, \dots, y_P) = \prod_{i=1}^n \prod_{j=1}^J p_j^{y_{ij}} N(x_i; \mu_j, \sigma_j^2)^{y_{ij}}$$

Complete-data log-likelihood

$$l(\mu, \sigma^2, p | x, y) = \sum_i \sum_j y_{ij} \{ \log p_j + \log N(x_i; \mu_j, \sigma_j^2) \}$$

$$\propto \sum_{i,j} y_{ij} \{ \log p_j - \frac{1}{2\sigma_j^2} (x_i - \mu_j)^2 - \log \sigma_j \}$$

Σ -Step. : Need to calculate the expected proportion of mixture component

$$x_i \rightarrow P_{\hat{i}} \quad \begin{matrix} \vdots \\ \text{Need to estimate. ... } w_{\hat{i}j}^{(t)} \end{matrix}$$

$$(P_{\hat{i}j}) \quad \text{indicator variable.}$$

$$w_{\hat{i}j}^{(t)} = E(y_{\hat{i}j} | x_i, \mu^{(t)}, \sigma^2, p^{(t)})$$

$$\text{for individual weight} = P(y_{\hat{i}j} = 1 | x_i, \mu^{(t)}, \sigma^2, p^{(t)})$$

$$= \frac{P_j^{(t)} N(x_i; \mu_j^{(t)}, \sigma_j^{2(t)})}{\sum_{j=1}^J P_j^{(t)} N(x_i; \mu_j^{(t)}, \sigma_j^{2(t)})}$$

$$M\text{-Step} : P_j^{(t+1)} = \frac{1}{n} \sum_{i=1}^n w_{\hat{i}j}^{(t)}$$

$$\mu_j^{(t+1)} = \frac{\sum_i w_{\hat{i}j}^{(t)} \cdot x_i}{\sum_i w_{\hat{i}j}^{(t)}}$$

$$\sigma_j^{2(t+1)} = \frac{\sum_i w_{\hat{i}j}^{(t)} (x_i - \mu_j^{(t+1)})^2}{\sum_i w_{\hat{i}j}^{(t)}}$$

(Ex2) Mixed - Effect Model

$$Y_i = X_i \beta + Z_i b_i + \varepsilon_i$$

fixed effect

random effect

$$b_i \sim N_q(0, D)$$

$$\varepsilon_i \sim N_{n_i}(0, \sigma^2 I_{n_i})$$

Our Object: Estimate β, D, σ^2

covariance matrix for random effect.

Observed-data log-likelihood

$$l(\beta, D, \sigma^2 | Y_1, \dots, Y_n) = \sum_{i=1}^n \left\{ -\frac{1}{2} (Y_i - X_i \beta)^T \Sigma_i^{-1} (Y_i - X_i \beta) - \frac{1}{2} \log |\Sigma_i| \right\}$$

We can directly maximize (β, D, σ^2) using Newton-Raphson method

$$\Sigma_i = Z_i D Z_i^T + \sigma^2 I$$

IWLS