

Empirical Bayes methods & the EM algorithm

- (I) Setup & Notation
- (II) Introductory Example
- (III) Basic Ideas
- (IV) EM Algorithm

Efron - Large-scale inference: empirical Bayes methods for estimation, testing and prediction (2012)

McLachlan, Krishnan - The EM algorithm and extensions (2007)

(I) Setup & Notation

$X \in \mathbb{R}^d$ unknown parameter, $Z \in \mathbb{R}^n$ measurement

Last time Input: • prior p_X of X

- likelihood model $\{f_Z(\cdot | X=x) | x \in \mathbb{R}^d\}$

- $Z_1, \dots, Z_M \sim Z_{\text{True}}$, $Z_{\text{True}} \sim f_Z(\cdot | X=x_{\text{True}})$ measurements for the same person / individual with true parameter value x_{True}

Output: • Estimator \hat{x}_M of $X=x_{\text{True}}$ & its uncertainty

Today

Input: • ~~prior p_X~~

- likelihood model $\{f_Z(\cdot | X=x) | x \in \mathbb{R}^d\}$

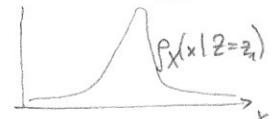
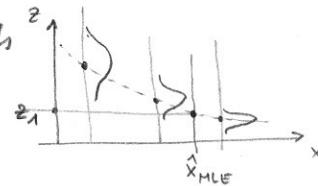
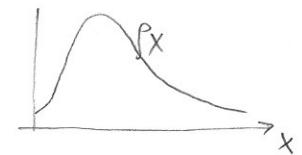
- $\vec{Z} = (Z_1, \dots, Z_M)$ measurements of several (M) individuals with M parametrizations X_j , i.e.

$$Z_m \sim f_Z(\cdot | X=X_m) \text{ indep.}, \quad X_m \stackrel{\text{iid}}{\sim} p_X \text{ unknown}$$

$$\vec{X} := (X_1, \dots, X_M)$$

Output: • Estimation of p_X

- Estimators $\hat{x}_1, \dots, \hat{x}_M$ of X_1, \dots, X_M (& their uncertainty)



(II) Introductory Example

- (a) Estimating the weight X of a person after measuring it with inaccurate weighing scales:

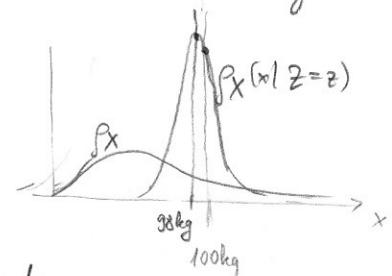
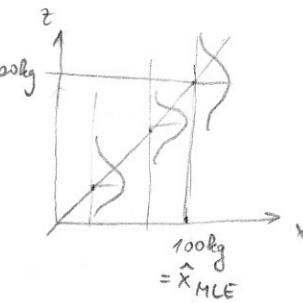
Assume that we measure $z = 100\text{kg}$.

Is $\hat{x}_{MLE} = 100\text{ kg}$ a good estimator? Yes, unless we have additional knowledge:

- (b) If we know the weight distribution p_X within the population,

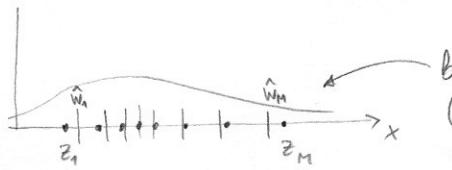
Bayes' rule gives us better criteria for estimating $X|Z=z$;

e.g. $\hat{x}_{MAP} = \max_x p_X(x|Z=z)$ maybe 98 kg.



- (c) Usually, p_X is unknown, but assume we have measurements

$z_1 = z_1, \dots, z_M = z_M$ for several persons (with unknown weights $X_1, \dots, X_M \sim p_X$)



better estimators than $\hat{x}_{MLE} = z_M$ for $M > 3$
(in the mean squared error sense), see also
"James Stein Estimator"

(III) Basic ideas

Roughly speaking,

- (1) MLE treats every person separately
- (2) Bayesian approach
- (3) Empirical Bayes methods
 - (i) first gather all measurements to estimate the prior p_X
 - (ii) then use this information to infer better individual estimators (usually by Bayes' rule)

[(i) and (ii) are sometimes hidden within one formula for the estimators]

We will concentrate on (i), approximations of p_X will be denoted by f .

Main observation: Viewing the prior as a parameter ("hyper-parameter") $\mathcal{F} = f$ (with $f_{\text{true}} = p_X$), the measurements z_1, \dots, z_M are iid realisations of the same experiment (last time's case) with likelihood model

$$p_Z(z | \mathcal{F} = f) = \int p_Z(z | X=x) \underbrace{p_X(x | \mathcal{F} = f)}_{(=f(x) \text{ in non-parametric case})} dx \quad \text{"(marginal) likelihood"}$$

and total likelihood

$$p_{\vec{Z}}(\vec{z} | \mathcal{F} = f) = \int p_Z(z_i | \vec{X}=\vec{x}) p_X(\vec{x} | \mathcal{F} = f) d\vec{x} = \prod_{m=1}^M p_Z(z_m | \mathcal{F} = f)$$

discrete case: $P(Z=z | \mathcal{F} = f) = \sum_x P(z=z \wedge X=x | \mathcal{F} = f) = \sum_x P(z=z | X=x \wedge \mathcal{F} = f) P(X=x | \mathcal{F} = f)$

$$= \sum_x P(z=z | X=x) P(X=x | \mathcal{F} = f)$$

So, now we can perform MLE estimation (and even Bayesian inference, after specifying a "hyper-prior", an a priori distribution of all possible priors) for the hyper-parameter \mathcal{F} using the new likelihood model $\{p_Z(\cdot | \mathcal{F} = f) \mid f \in \mathcal{F}\}$

\uparrow set of considered priors

/assume

- Two cases:
- We know that the prior p_X has a certain parametric form, e.g. $p_X = N(\mu, \Sigma)$
 \rightsquigarrow hyper-parameters = $(\mu \in \mathbb{R}^d, \Sigma \in \mathbb{R}^{d \times d} \text{ spd.})$, finite-dim. problem
 $\uparrow \uparrow$ unknown
 - p_X is completely unknown \rightsquigarrow hyper-parameter = $\mathcal{F} \in \mathcal{U}_{\mathbb{R}^d}$, ∞ -dim. problem
 \uparrow set of all prob. dist. on \mathbb{R}^d

In both cases, the EM algorithm is the standard tool to maximize the (marginal) likelihood (PMLE & NPMLE).

(IV) Expectation - maximization (EM) algorithm

Iterative algorithm that (locally) maximizes the (marginal) log-likelihood

$$\mathcal{L}(f) = \log \int_{\vec{z}} p_{\vec{z}}(\vec{z} | F=f) = \log \prod_{m=1}^M p_z(z_m | F=f) = \sum_{m=1}^M \log \int_{\vec{z}} p_{\vec{z}}(z | X=x) \underbrace{p_X(x | F=f)}_{\text{often } f(x)} dx$$

Starting with f_0 iterate:

E-step: Formulate the "complete data likelihood fct."
 $(\vec{X} = \vec{x}, Z = \vec{z})$

$$L_c(\vec{x}, \vec{z} | F=f) = \prod_{m=1}^M p_z(z_m | X=x_m) \underbrace{p_X(x_m | F=f)}_{\text{often } f(x)}$$

and the expectation value in \vec{X} given $F=f_n$, $\vec{Z}=\vec{z}$

$$\begin{aligned} Q(f | f_n) &:= \mathbb{E} \left[\log L_c((\vec{X} | \vec{Z}=\vec{z}, F=f_n), \vec{z} | F=f) \right] \\ &= \mathbb{E} \left[\log L_c(\vec{X}, \vec{z} | F=f) \mid \vec{Z}=\vec{z}, F=f_n \right] \\ &= \int \log L_c(\vec{x}, \vec{z} | F=f) \underbrace{p_{\vec{X}}(\vec{x} | \vec{Z}=\vec{z}, F=f_n)}_{\substack{\text{density of } (\vec{X} | \vec{Z}=\vec{z}, F=f_n) \\ = \text{posterior, if } f_n \text{ was the prior}}} d\vec{x} \end{aligned}$$

M-step: maximize $Q(\cdot | f_n)$: $f_{n+1} = \underset{f}{\operatorname{argmax}} Q(f | f_n)$

(IV.1) Some theory:

$$\begin{aligned} 1) \mathcal{L}(f) &= \log \int_{\vec{z}} p_{\vec{z}}(\vec{z} | F=f) = \log \int_{\vec{z}} \int_{\vec{X}} p_{\vec{z}}(\vec{z} | \vec{X}=\vec{x}) p_{\vec{X}}(\vec{x} | F=f) \cdot \underbrace{\frac{p_{\vec{X}}(\vec{x} | \vec{Z}=\vec{z}, F=f_n)}{p_{\vec{X}}(\vec{x} | \vec{Z}=\vec{z}, F=f_n)}}_{\substack{\text{density in } \vec{X} \\ d\vec{x}}} \\ &\stackrel{\text{Jensen}}{\geq} \int_{\vec{X}} p_{\vec{X}}(\vec{x} | \vec{Z}=\vec{z}, F=f_n) \log \left(\frac{p_{\vec{z}}(\vec{z} | \vec{X}=\vec{x}) p_{\vec{X}}(\vec{x} | F=f)}{p_{\vec{X}}(\vec{x} | \vec{Z}=\vec{z}, F=f_n) p_{\vec{z}}(\vec{z} | F=f_n)} \right) d\vec{x} \\ &= \mathcal{L}(f_n) + \underbrace{\int_{\vec{X}} p_{\vec{X}}(\vec{x} | \vec{Z}=\vec{z}, F=f_n) \log \left(\frac{p_{\vec{z}}(\vec{z} | \vec{X}=\vec{x}) p_{\vec{X}}(\vec{x} | F=f)}{p_{\vec{X}}(\vec{x} | \vec{Z}=\vec{z}, F=f_n) p_{\vec{z}}(\vec{z} | F=f_n)} \right) d\vec{x}}_{=: l(f | f_n)} \end{aligned}$$

$$2) l(f | f_n) = 0 \text{ for } f=f_n, \text{ because then denominator} \xrightarrow{\text{Bayes}} \frac{p_{\vec{X}}(\vec{x} | F=f_n) p_{\vec{z}}(\vec{z} | \vec{X}=\vec{x}, F=f_n)}{p_{\vec{z}}(\vec{z} | F=f_n)} = \frac{p_{\vec{z}}(\vec{z} | F=f_n)}{p_{\vec{z}}(\vec{z} | F=f_n)} = \text{numerate}$$

$$3) \underset{f}{\operatorname{argmax}} \mathcal{L}(f | f_n) = \underset{f}{\operatorname{argmax}} \int_{\vec{X}} p_{\vec{X}}(\vec{x} | \vec{Z}=\vec{z}, F=f_n) \log \left(\frac{p_{\vec{z}}(\vec{z} | \vec{X}=\vec{x}) p_{\vec{X}}(\vec{x} | F=f)}{p_{\vec{z}}(\vec{z} | F=f_n) p_{\vec{X}}(\vec{x} | F=f_n)} \right) d\vec{x} = \underset{f}{\operatorname{argmax}} Q(f | f_n)$$

Therefore:

$$f_{n+1} = \underset{f}{\operatorname{argmax}} Q(f | f_n) \stackrel{(3)}{=} \underset{f}{\operatorname{argmax}} l(f | f_n) \stackrel{(2)}{\Rightarrow} l(f_{n+1} | f_n) \geq 0 \stackrel{(1)}{\Rightarrow} \boxed{\mathcal{L}(f_{n+1}) \geq \mathcal{L}(f_n)}$$

Picture?

(IV.2) Further applications of EM

The EM algorithm has further applications in statistics and data analysis, e.g.
in incomplete data analysis
grouped / censored / truncated

Example: You model the age of death (AoD) of patients undergoing some treatment by an exponential distribution $X_m = \text{AoD}_m \sim \text{Exp}(\mu)$ with unknown mean $\mu = \mu$

$$p_X(x | \mu = \mu) = \mu^{-1} \exp(-x/\mu) \cdot \mathbb{1}_{(0, \infty)}$$

However, when evaluating the data some patients (w.l.o.g. $m=1, \dots, r$) are still alive, so your data is incomplete.

$$\vec{z}_m = \begin{pmatrix} a_m & d_m \end{pmatrix}$$

↑
age reached ↗
 $d_m = \begin{cases} 1 & \text{if uncensored (person died at age } a_m) \\ 0 & \text{if censored (AoD} \geq a_m\text{)} \end{cases}$

Since $\text{Exp}(\mu)$ is memoryless we have for $m=1, \dots, r$

$$(X_m | \vec{z}_m = \vec{z}_m) = a_m + \tilde{X}_m, \text{ where } \tilde{X}_m \sim \text{Exp}(\mu)$$

E-step: $L_c(\vec{x}, \vec{z} | \mu = \mu) = \prod_{m=1}^M \mu^{-1} \exp(-x_m/\mu) = \prod_{m=1}^M \mu^{-1} \exp(-a_m/\mu) \cdot \prod_{m=1}^r \exp(-\tilde{x}_m/\mu)$

$$\log L_c(\vec{x}, \vec{z} | \mu = \mu) = -M \log \mu - \sum_{m=1}^M a_m/\mu - \sum_{m=1}^r \tilde{x}_m/\mu$$

$$Q(\mu | \mu_n) = \mathbb{E} \left[\log L_c((\vec{x} | \vec{z} = \vec{z}, \mu = \mu_n), \vec{z} | \mu = \mu) \right] = -M \log \mu - \sum_{m=1}^M a_m/\mu - \sum_{m=1}^r \frac{\mu_n}{\mu}$$

M-step: $Q(\mu | \mu_n) = -M \log \mu - \mu^{-1} (r \cdot \mu_n + \sum_{m=1}^M a_m)$

$$0 \stackrel{!}{=} \partial_\mu Q(\mu | \mu_n) = -\frac{M}{\mu} + \frac{1}{\mu^2} (r \cdot \mu_n + \sum_{m=1}^M a_m)$$

$$\Rightarrow \mu_{n+1} = \frac{1}{M} (r \cdot \mu_n + \sum_{m=1}^M a_m)$$

$\boxed{\mu_n \rightarrow \mu_* \text{ with } \mu_* = \frac{1}{M} (r \cdot \mu_n + \sum_{m=1}^M a_m)}, \text{ i.e. } \mu_* = \frac{1}{M-r} \sum_{m=1}^M a_m$

μ_* truly maximizes the likelihood

$$p_{\vec{z}}(\vec{z} | \mu = \mu) = \prod_{m=1}^r \underbrace{\frac{P(X \geq a_m | \mu = \mu)}{\exp(-a_m/\mu)}}_{\mu^{-1} \exp(-a_m/\mu)} \cdot \prod_{m=r+1}^M \underbrace{p_X(a_m | \mu = \mu)}_{\mu^{-1} \exp(-a_m/\mu)}$$

as can be seen by differentiating its logarithm