

For  $x = 5$   $l'(p) = \frac{5}{p} > 0$  and  $l(p)$  is increasing

So for  $\Theta = [0, 1]$  The MLE is  $\frac{X}{5}$

One context in which the likelihood equations are especially nice is the exponential family - for  $X_1, X_2, \dots, X_n$

$$f(x|\eta) = c(\eta)^n \prod_{i=1}^n h(x_i) \exp \sum_{j=1}^k \eta_j \left( \sum_{i=1}^n T_j(x_i) \right)$$

$$\frac{\text{B+D Form}}{A(\eta) = -\log c(\eta)} \Leftrightarrow \exp(-nA(\eta)) \prod_{i=1}^n h(x_i) \exp \sum_{j=1}^k \eta_j \left( \sum_{i=1}^n T_j(x_i) \right)$$

$$\text{So } l(\eta) = -nA(\eta) + \sum_{i=1}^n \ln h(x_i) + \sum_{j=1}^k \eta_j \left( \sum_{i=1}^n T_j(x_i) \right)$$

$$\text{and } \frac{\partial}{\partial \eta_j} l(\eta) = -\frac{\partial}{\partial \eta_j} A(\eta) + \sum_{i=1}^n T_j(x_i)$$

By Corollary 1.6.1 B+D,  $A$  has nonempty interior

$$\begin{aligned} \text{or claim 9} \\ \text{on Exponential} \\ \text{Family handout} \end{aligned} \quad \frac{\partial}{\partial \eta_j} A(\eta) \left( = \frac{\partial}{\partial \eta_j} (-\log c(\eta)) \right)$$

$$= E_{\eta} T_j(X)$$

Thus, the  $j$ th of these equations is

$$-n E_{\eta} T_j(X) + \sum_{i=1}^n T_j(x_i) = 0$$

$$E_{\eta} T_j(X) = \frac{1}{n} \sum_{i=1}^n T_j(x_i) \quad \left. \vphantom{E_{\eta} T_j(X)} \right\} \begin{array}{l} \text{jth} \\ \text{likelihood} \\ \text{equation} \end{array}$$

So, the ML estimating equation in an exponential family with nonempty interior in  $\mathcal{E}$  is "set the theoretical mean of  $T$  equal to the empirical mean of  $T$ " — a natural exponential family is "MGM"

The question of whether there must be a solution to the exponential family likelihood equations and if any solution must be unique and maximize the likelihood is discussed in BTD Section 2.3 — one simple result from there is

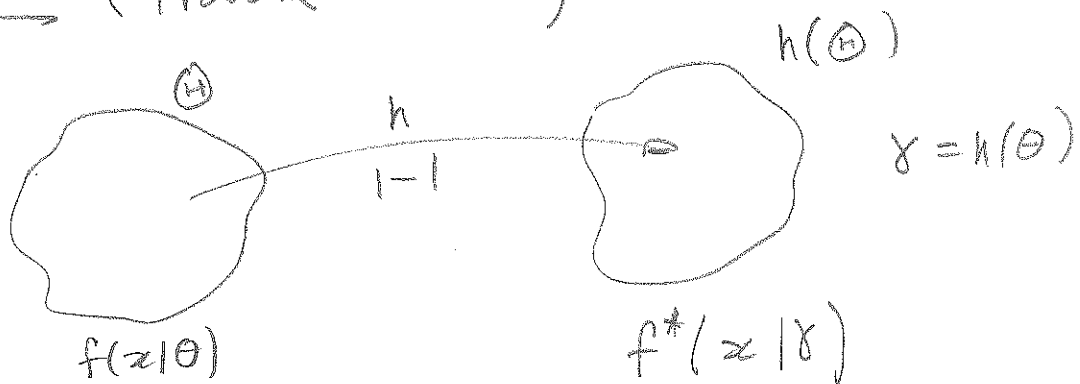
Corollary 2.3.2 If the equations

$$E_{\eta} T_j(X) = \frac{1}{n} \sum_{i=1}^n T_j(z_i) \quad j=1, 2, \dots, k$$

have a solution  $\eta$  in the interior of  $\mathcal{E}$  it is the unique MLE of  $\eta$

Both for application in exponential families where the "natural" parameterization isn't the "standard" or "usual" parameterization and in other contexts, there is the lemma

Lemma (Problem 2.2.16a)



clearly  $\hat{\theta}$  maximizes  $f(z|\theta)$  over  $\Theta$

$\Rightarrow \hat{\gamma} = h(\hat{\theta})$  maximizes  $f^*(z|\gamma)$  over  $h(\Theta)$

So, e.g.,  $\hat{\eta}$  MLE of natural parameter if  $\theta$  can be written as  $h(\eta) = \theta$  for a 1-1 function  $h$ , then

$$\hat{\theta} = h(\hat{\eta}) \text{ is the MLE of } \theta$$

Solving the likelihood equations

$$\nabla \ell(\theta) = 0$$

is "just" a numerical analysis problem — various iterative procedures (bisection, Newton-Raphson, etc.) might be applied to do the job — mostly, their discussion is not a topic for ST3 —

one procedure for the problem has some probabilistic/statistical content and can/should be touched here is the "EM algorithm" —

day 17

Basic idea: Sometimes an observable  $Y$  has a nasty loglikelihood  $\ell_Y(\theta)$  but could be thought of as distributionally equivalent to  $S(X)$  for some (potentially completely fictitious)  $X$  for which computation and optimization of  $\ell_X(\theta)$  is feasible and  $S(X)$  data is hand

$$E_{\theta_0} [\ell_X(\theta) \mid S(X) = y]$$

is feasible — In such cases I might

o. pick some starting value  $\theta^{(0)}$

1. Find

E step  $E_{\theta^{(0)}} [ l_X(\theta) \mid S(X) = y ]$

M step 2. Optimize this as a function of  $\theta$  to find  $\theta^{(1)}$

3. Replace  $\theta^{(0)}$  with  $\theta^{(1)}$

4. Iterate to convergence  $\frac{1}{\text{mean}}$

Example  $X \sim \text{exp}(\lambda)$

$$Y = X \mathbb{I}[1 < X < 2] + \mathbb{I}[X \leq 1] + 2 \mathbb{I}[2 \leq X]$$

(Y a censored version of X) - Suppose

$Y_1, Y_2, \dots, Y_n$  iid with this dsu

$$f(y|\lambda) = \begin{cases} 1 - e^{-\lambda} & y = 1 \\ \lambda e^{-\lambda y} & 1 < y < 2 \\ e^{-2\lambda} & y = 2 \end{cases}$$

If  $Y_i = 1$  then  $X_i$  has density  $\propto \lambda e^{-\lambda x}$  on  $(0, 1]$  and if  $Y_i = 2$   $X_i$  has density  $\propto \lambda e^{-\lambda x}$  on  $[2, \infty)$

$$l_X(\lambda) = n \log \lambda - \lambda \sum X_i$$

Then for a particular  $\lambda_0$  conditioned on  $Y = y$

$\sum X_i$  has mean

$$\sum_{y_i \in (1,2)} y_i + \# [y_i = 1] E_{\lambda_0} [X | X < 1] \\ + \# [y_i = 2] E_{\lambda_0} [X | X > 2]$$

$$E_{\lambda_0} [X | X < 1] = \frac{1}{1 - e^{-\lambda_0}} \int_0^1 x \lambda_0 e^{-x \lambda_0} dx \\ = \frac{1}{1 - e^{-\lambda_0}} \left[ -x e^{-\lambda_0 x} \Big|_0^1 + \int_0^1 e^{-x \lambda_0} dx \right] \\ = \frac{1}{1 - e^{-\lambda_0}} \left[ -e^{-\lambda_0} + \frac{1}{\lambda_0} (1 - e^{-\lambda_0}) \right] \\ = \frac{1}{\lambda_0} - \frac{e^{-\lambda_0}}{(1 - e^{-\lambda_0})}$$

$$E_{\lambda_0} [X | X > 2] = \frac{1}{e^{-2\lambda_0}} \int_2^{\infty} x \lambda_0 e^{-x \lambda_0} dx \\ = e^{2\lambda_0} \left[ -x e^{-\lambda_0 x} \Big|_2^{\infty} + \int_2^{\infty} e^{-x \lambda_0} dx \right] \\ = e^{2\lambda_0} \left[ 2e^{-2\lambda_0} + \frac{1}{\lambda_0} e^{-2\lambda_0} \right] \\ = 2 + \frac{1}{\lambda_0}$$

i.e. the  $\lambda_0$  conditional mean of  $l_X(\lambda)$  given  $Y=y$

is

$$n \log \lambda - \lambda \left[ \sum_{y \in (1,2)} y_i + \# [y_i = 1] \left( \frac{1}{\lambda_0} - \frac{e^{-\lambda_0}}{1 - e^{-\lambda_0}} \right) \right. \\ \left. + \# [y_i = 2] \left( 2 + \frac{1}{\lambda_0} \right) \right]$$

E-step

which is maximized at  $\lambda^{(1)} = \frac{n}{\text{above}}$   
M-step

So ultimately, an iterative E-M algorithm sets

$$\lambda^{(j+1)} = \frac{n}{\left[ \sum_{y_i \in (1,2)} y_i + \# [y_i=1] \left( \frac{1}{\lambda^{(j)}} - \frac{e^{-\lambda^{(j)}}}{1-e^{-\lambda^{(j)}}} \right) + \# [y_i=2] \left( 2 + \frac{1}{\lambda^{(j)}} \right) \right]}$$

and iterates to convergence as opposed to straight-up optimization of

$$l_Y(\lambda) = \# [y_i=1] \log(1-e^{-\lambda}) + \# [1 < y_i < 2] \log \lambda - \lambda \sum_{y_i \in (1,2)} y_i + \# [y_i=2] 2\lambda$$

day 18

Example  $k$  possible outcomes,  $n$  independent identical trials

$P_j$  = probability that any trial produces outcome  $j$

$$X_{ij} = I[\text{trial } i \text{ produces outcome } j]$$

$$n_j = \sum_{i=1}^n X_{ij} = \# \text{ of trials producing outcome } j$$

Suppose that  $p$  has each  $p_i \geq 0$  w  $\sum p_i = 1$

pdf for  $X$  is

$$f(x|p) = \prod_{j=1}^k P_j^{n_j}$$

(for  $x$  with  $x_{ij} = 0$  or  $1$   
 $\sum_j x_{ij} = 1$ ,  $n_j \geq 0$  integers  
 $\sum_{j=1}^k n_j = n$ )

MLE of  $p$  is (as it turns out)

$$\hat{p} = \left( \frac{n_1}{n}, \frac{n_2}{n}, \dots, \frac{n_k}{n} \right)$$

But now, what if, e.g.  $k=4$  and all we observe is this

Trial	Information Available	In Terms of $X$ 's
1	outcome is 1	$X_{11} = 1$
2	3	$X_{23} = 1$
3	2 or 4	$X_{32} + X_{34} = 1$
4	2	$X_{42} = 1$
5	3	$X_{53} = 1$
6	2 or 3	$X_{62} + X_{63} = 1$
7	1	$X_{71} = 1$
8	1 or 2	$X_{81} + X_{82} = 1$
9	2	$X_{92} = 1$
10	4	$X_{104} = 1$

That is, I don't get the 10  $X_i = (X_{i1}, X_{i2}, \dots, X_{i4})$   
but, e.g.  $Y_3 = (X_{31}, X_{33}, X_{32} + X_{34})$

The likelihood function based on  $Y_i$ 's (information available) not  $X_i$ 's is

$$L_Y(p) = p_1^2 p_2^2 p_3^2 (1 - p_1 - p_2 - p_3) \\ \times (1 - p_1 - p_3) \underset{\substack{\uparrow \\ p_2 + p_4}}{(p_2 + p_3)} (p_1 + p_2)$$

and I could probably optimize  $L_Y(p)$  or  $l_Y(p) = \log L_Y(p)$  via some numerical method

Another possibility is to use the EM algorithm

$$l_X(p) = n_1 \log p_1 + n_2 \log p_2 + n_3 \log p_3 + n_4 (1 - p_1 - p_2 - p_3)$$

For the data in hand this is

$$= (2 + X_{e1}) \log p_1 + (2 + X_{32} + X_{62} + X_{e2}) \log p_2 \\ + (2 + X_{63}) \log p_3 + (1 + X_{34}) \log (1 - p_1 - p_2 - p_3)$$

For any particular  $P_0 = (P_{01}, P_{02}, P_{03} + P_{04})$

$$E_{P_0} [X_{31} | Y = \text{data in hand}] = \frac{P_{01}}{P_{01} + P_{02}}$$

$$E_{P_0} [X_{32} | Y = \text{data in hand}] = \frac{P_{02}}{P_{02} + P_{04}}$$

$$E_{P_0} [X_{62} | Y = \text{data in hand}] = \frac{P_{02}}{P_{02} + P_{03}}$$



$$E_{p_0} [X_{82} | \quad] = \frac{p_{02}}{p_{01} + p_{02}}$$

$$E_{p_0} [X_{63} | \quad] = \frac{p_{03}}{p_{02} + p_{03}}$$

$$E_{p_0} [X_{34} | \quad] = \frac{p_{04}}{p_{02} + p_{04}}$$

Then

$$E_{p_0} [l_X(p) | Y = \text{data in hand}] \quad b(p_0)$$

$$\begin{aligned} a(p_0) &= \left(2 + \frac{p_{01}}{p_{01} + p_{02}}\right) \log p_1 + \left(2 + p_{02} \left(\frac{1}{p_{02} + p_{04}} + \frac{1}{p_{02} + p_{03}} + \frac{1}{p_{01} + p_{02}}\right)\right) \log p_2 \\ &+ \left(2 + \frac{p_{03}}{p_{02} + p_{03}}\right) \log p_2 + \left(1 + \frac{p_{04}}{p_{02} + p_{04}}\right) \log (1 - p_1 - p_2 - p_3) \\ &\quad c(p_0) \quad d(p_0) \end{aligned}$$

Note

$$a(p_0) + b(p_0) + c(p_0) + d(p_0) = 10 = n$$

$$p^{(j+1)} = \left( \frac{a(p^{(j)})}{10}, \frac{b(p^{(j)})}{10}, \frac{c(p^{(j)})}{10}, \frac{d(p^{(j)})}{10} \right)$$

iterate to a fixed point

BTD give some arguments as to why EM might work - see, in particular, Lemma 2.4.1 That says that at each step  $l_Y(\theta)$  never decreases - the standard complaint w/ EM is that its convergence is often very slow - you need a good starting point and if other methods are possible they may be faster -

One matter that should be raised regarding EM is that B+D phrase their version of it concerns not

$$E_{\theta_0} [l_x(\theta) | Y=y]$$

but rather

$$E_{\theta_0} [l_x(\theta) - l_x(\theta_0) | Y=y]$$

sometimes I'll be able to compute the 2nd when I couldn't compute the 1st - and optimization of the first is equivalent to optimization of the 2nd

B+D give some arguments why EM might work - see, in particular, Lemma 2.9.1 That says that at each step  $l_x(\theta)$  never decreases - the standard complaint about EM is that its convergence is often very slow - you need a good starting point and if other methods are possible, they may be faster -

### Bayes Estimators

we've said repeatedly that where  $g(\theta)$  specifies a prior  $dsn$

$$g(\theta | z) \propto L_z(\theta) g(\theta)$$

and that characteristics of the posterior serve as estimators of  $\gamma(\theta)$  - e.g.

under SEL, the Bayes estimator of  $Y(\theta)$  is

$$\delta(x) = E[Y(\theta) | X=x]$$

under AEL, the Bayes estimator of  $Y(\theta)$  is

$$\delta(x) = \text{median of the den of } Y(\theta) | X=x$$

under WEL with weight function  $w(\theta) \geq 0$ , the Bayes estimator of  $Y(\theta)$  is

$$\delta(x) = \frac{E[Y(\theta)w(\theta) | X=x]}{E[w(\theta) | X=x]}$$

In fact, this story can sometimes be generalized by not requiring  $g(\theta)$  to specify a probability den

Example

$$X \sim N(\theta, 1)$$

$$g(\theta) = 1$$

$$\int_{-\infty}^{\infty} 1 d\theta = \infty !!!$$

$$L_X(\theta)g(\theta) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(x-\theta)^2\right) \cdot 1$$

and we can use  $L_X(\theta)g(\theta)$  to specify a "posterior"  $(N(x, 1))$  and  $\therefore$

"generalized Bayes" estimator  $\delta(x) = x$  even though  $f(x|\theta)g(\theta)$  doesn't define a (joint) probability den for  $(X, \theta)$  —

note that

$$\begin{aligned} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(x-\theta)^2\right) dx d\theta \\ = \int_{-\infty}^{\infty} 1 d\theta = \infty \end{aligned}$$

What is (Bayes) optimal is in principle "clear"  
 The problem of actually computing characteristics  
 of a dsr with density proportional to

$$L(\theta)g(\theta)$$

Note that even finding an exact density for the  
 posterior requires finding

$$\int L(\theta)g(\theta) d\theta$$

the normalizer for the function  $L(\theta)g(\theta)$

Modern Bayes computation is based on simulation  
 to substitute for calculus - e.g. if I  
 am interested in

$$Q = \int q(\theta)g(\theta|z) d\theta$$

for some  $q: \mathbb{R}^k \rightarrow \mathbb{R}^1$  and can somehow  
 generate  $\theta_1^*, \theta_2^*, \dots$  iid  $G(\theta|z)$  I might  
 approximate  $Q$  by

$$\hat{Q}_n = \frac{1}{n} \sum q(\theta_i^*)$$

relying on the LLN to conclude that

$$\frac{1}{n} \sum_{i=1}^n q(\theta_i^*) \xrightarrow{P} Q$$

$G(\theta|z)$   
Probability

(e.g. I might for  $\theta = (\theta_1, \theta_2)$  have  $q(\theta) = \theta_1$  or  $q(\theta) = I[\theta_1 < \theta_2]$ )

The problem with this idea is that naive ways of  
 doing simulation would require that I know  
 $\int L(\theta)g(\theta) d\theta$  in order to simulate the  $\theta_i^*$  -

Happily there is the famous "rejection algorithm"