

## Contents

<b>1</b>	<b>Foundation</b>	<b>2</b>
1.1	Review on MLE . . . . .	2
1.2	Review on GLM . . . . .	4
1.2.1	Newly introduced in 711 . . . . .	5
1.3	Families of distributions, moments / cumulants, and quantiles / percentiles . . . . .	6
1.4	Transformation of random variables . . . . .	8
1.5	Others . . . . .	8
<b>2</b>	<b>Likelihood construction and estimation</b>	<b>8</b>
2.1	Introduction . . . . .	8
2.2	Likelihood construction . . . . .	9
2.3	Proportional likelihoods . . . . .	10
2.4	Empirical distribution function as an MLE . . . . .	11
2.5	Likelihoods for censored / truncated data . . . . .	12
2.6	Likelihoods for regression models . . . . .	12
2.7	Marginal and conditional likelihoods . . . . .	15
2.8	MLE and information matrix . . . . .	16
2.8.1	Transformed and modeled parameters . . . . .	19
2.9	Methods for maximizing the likelihood . . . . .	19
2.9.1	Why does EM work? . . . . .	22
2.9.2	Calculating observed info matrix after EM . . . . .	23
2.10	Uniqueness of MLE . . . . .	24
<b>3</b>	<b>Likelihood-based tests and confidence regions</b>	<b>26</b>
3.1	Simple null hypothesis . . . . .	27
3.2	Composite null hypothesis . . . . .	28
3.3	Confidence interval . . . . .	31
3.4	Nonstandard hypothesis testing problems . . . . .	32
<b>4</b>	<b>Bayesian methods</b>	<b>33</b>
4.1	Introduction . . . . .	33
4.2	Bayesian estimator from decision theory perspective . . . . .	34
4.3	Credible intervals . . . . .	35
4.4	Conjugate prior . . . . .	35
4.5	Noninformative prior . . . . .	36
4.6	Normal data examples . . . . .	37
4.6.1	One sample with unknown mean and variance . . . . .	37
4.6.2	Two samples . . . . .	38
4.6.3	Normal linear model . . . . .	38
4.7	Hierarchical Bayes and empirical Bayes . . . . .	39
4.7.1	James–Stein estimation . . . . .	39
4.7.2	Meta-analysis applications of hierarchical and empirical Bayes . . . . .	40
4.8	Monte Carlo estimation of a posterior . . . . .	41
4.8.1	Noniterative Monte Carlo methods . . . . .	41
4.9	MCMC methods . . . . .	42
4.9.1	Substitution sampling . . . . .	42
4.9.2	Gibbs sampling . . . . .	43
4.9.3	Metropolis-Hastings algorithm . . . . .	44
4.9.4	Hybrid forms . . . . .	45
<b>5</b>	<b>Large sample theory</b>	<b>45</b>

<b>6 M-Estimation (Estimating Equations)</b>	<b>45</b>
6.1 Introduction . . . . .	45
6.2 Basic approach . . . . .	46
6.2.1 Estimation for A, B, and V . . . . .	46
6.3 Delta method via M-estimation . . . . .	47

# 1 Foundation

## 1.1 Review on MLE

**Central limit theorem:** Suppose  $\{X_1, X_2, \dots, X_n\}$  is a sequence of i.i.d. random variables with  $E(X_i) = \mu$  and  $Var(X_i) = \sigma^2 < \infty$  then as  $n \rightarrow \infty$ ,  $\bar{X} \rightarrow N(\mu, \frac{\sigma^2}{n})$

**Score:** The partial derivative with respect to  $\theta$  of the natural logarithm of the likelihood function is called the score

$$Z = l' = \frac{\partial}{\partial \theta} \log f(X; \theta)$$

$$E(Z) = 0 \text{ and } Z \xrightarrow{d} N(0, I(\theta_0))$$

under  $\theta_0$

**Fisher information:** The variance of the score is defined to be the Fisher information

$$\mathcal{I}(\theta) = E \left[ \left( \frac{\partial}{\partial \theta} \log f(X; \theta) \right)^2 \mid \theta \right] = -E \left[ \frac{\partial^2}{\partial \theta^2} \log f(X; \theta) \mid \theta \right]$$

Total information and observed total information

$$\mathbf{I}_T(\boldsymbol{\theta}) = -E \left\{ \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} \log L(\boldsymbol{\theta} \mid \mathbf{Y}) \right\}$$

$$\mathbf{I}_T(\mathbf{Y}, \boldsymbol{\theta}) = -\frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} \log L(\boldsymbol{\theta} \mid \mathbf{Y})$$

Information estimation

- Version 1

$$\begin{aligned} \bar{\mathbf{I}}(\mathbf{Y}, \boldsymbol{\theta}) &= n^{-1} \sum_{i=1}^n \left\{ -\frac{\partial}{\partial \boldsymbol{\theta}} \mathbf{s}(Y_i, \boldsymbol{\theta}) \right\} \\ &= n^{-1} \sum_{i=1}^n \left\{ -\frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} \log f(Y_i; \boldsymbol{\theta}) \right\} \end{aligned}$$

- Version 2

$$\begin{aligned} \bar{\mathbf{I}}^*(\mathbf{Y}, \boldsymbol{\theta}) &= n^{-1} \sum_{i=1}^n \mathbf{s}(Y_i, \boldsymbol{\theta})^{\otimes 2} \\ &= n^{-1} \sum_{i=1}^n \left\{ \frac{\partial}{\partial \boldsymbol{\theta}^T} \log f(Y_i; \boldsymbol{\theta}) \right\}^{\otimes 2} \end{aligned}$$

Property 1. If  $\hat{\theta}$  is the MLE estimate of  $\theta_0$ , then it has the following property:

$$\sqrt{n}(\hat{\theta} - \theta_0) \xrightarrow{d} N\left(0, \frac{1}{I(\theta_0)}\right)$$

Ways for estimating MLE:

- Newton-Raphson algorithm

$$\begin{aligned} \mathbf{0} &= \mathbf{S}(\boldsymbol{\theta}) \approx \mathbf{S}(\boldsymbol{\theta}^{(\nu)}) + \left\{ \frac{\partial}{\partial \boldsymbol{\theta}} \mathbf{S}(\boldsymbol{\theta}) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^{(\nu)}} \right\} (\boldsymbol{\theta} - \boldsymbol{\theta}^{(\nu)}) \\ &= \mathbf{S}(\boldsymbol{\theta}^{(\nu)}) - \mathbf{I}_T(\mathbf{Y}, \boldsymbol{\theta}^{(\nu)}) (\boldsymbol{\theta} - \boldsymbol{\theta}^{(\nu)}) \\ \boldsymbol{\theta}^{(\nu+1)} &= \boldsymbol{\theta}^{(\nu)} + \mathbf{I}_T(\mathbf{Y}, \boldsymbol{\theta}^{(\nu)})^{-1} \mathbf{S}(\boldsymbol{\theta}^{(\nu)}) \end{aligned}$$

1. start with initial  $\theta^{(0)}$

2. update from current  $\theta^{(\nu)}$  to obtain  $\theta^{(\nu+1)}$
3. stop if  $\left\| \mathbf{S} \left( \theta^{(\nu+1)} \right) \right\|$  or  $\left\| \theta^{(\nu+1)} - \theta^{(\nu)} \right\|$  is sufficiently small

Note: Newton-Raphson estimator has local quadratic convergence property

$$\left\| \theta^{(\nu+1)} - \hat{\theta}_{\text{MLE}} \right\| \leq c \left\| \theta^{(\nu)} - \hat{\theta}_{\text{MLE}} \right\|^2 \text{ for some } c > 0$$

This property exists for the situation where the estimator is very close to the true parameter. Thus, at least near the solution, convergence is fast for a Newton method.

More specifically, the local quadratic convergence holds under the following conditions

1.  $I_T(\mathbf{Y}, \theta) \neq 0$  in a neighborhood of  $\hat{\theta}_{\text{MLE}}$
2.  $S''(\theta)$  is bounded
3.  $\theta^{(\nu)}$  is sufficiently close to  $\hat{\theta}_{\text{MLE}}$

Proof draft:

$$\begin{aligned} 0 &= S \left( \hat{\theta}_{\text{MLE}} \right) = S \left( \theta^{(\nu)} \right) - I_T \left( \mathbf{Y}, \theta^{(\nu)} \right) \left( \hat{\theta}_{\text{MLE}} - \theta^{(\nu)} \right) \\ &\quad + \frac{1}{2} S'' \left( \vartheta^{(\nu)} \right) \left( \hat{\theta}_{\text{MLE}} - \theta^{(\nu)} \right)^2 \\ \hat{\theta}_{\text{MLE}} - \theta^{(\nu+1)} &= \frac{1}{2} \frac{S'' \left( \vartheta^{(\nu)} \right)}{I_T \left( \mathbf{Y}, \theta^{(\nu)} \right)} \left( \hat{\theta}_{\text{MLE}} - \theta^{(\nu)} \right)^2 \end{aligned}$$

where  $\vartheta^{(\nu)} = \theta^{(\nu)} + k * (\hat{\theta}_{\text{MLE}} - \theta^{(\nu)})$ ,  $k \in [0, 1]$ . This uses the Taylor expansion with Lagrange remainder. *One-step estimator*: typically,  $\hat{\theta}_{\text{MLE}} - \theta = O_p(n^{-1/2})$ . If one starts with  $\theta^{(0)}$  such that  $\theta^{(0)} - \theta = O_p(n^{-1/2})$ , then

$$\theta^{(1)} - \hat{\theta}_{\text{MLE}} = O_p(n^{-1})$$

under regularity conditions. That is  $\theta^{(0)}$  and  $\hat{\theta}_{\text{MLE}}$  is asymptotically equivalent.

- Fisher scoring algorithm

$I_T \left( \mathbf{Y}, \theta^{(\nu)} \right)$  replaced by its expectation  $\mathbf{I}_T \left( \theta^{(\nu)} \right)$ , which is

$$\theta^{(\nu+1)} = \theta^{(\nu)} + \mathbf{I}_T \left( \theta^{(\nu)} \right) \mathbf{S} \left( \theta^{(\nu)} \right)$$

- EM algorithm

The basic idea of the EM Algorithm is to view the observed data  $\mathbf{Y}$  as incomplete, that somehow there is missing data  $\mathbf{Z}$  that would make the problem simpler if we had it. In some cases  $\mathbf{Z}$  could truly be missing data, but in others it is just additional data that we wish we had.

1. The first step is to write down the joint likelihood of the “complete” data  $(\mathbf{Y}, \mathbf{Z})$ , call it  $L_C(\theta | \mathbf{Y}, \mathbf{Z})$ . The “E” step of the EM Algorithm is to compute the conditional expectation of  $\log L_C(\theta | \mathbf{Y}, \mathbf{Z})$  given  $\mathbf{Y}$  assuming the true parameter value is  $\theta^{(\nu)}$

$$\begin{aligned} Q \left( \theta, \theta^{(\nu)}, \mathbf{Y} \right) &= E_{\theta^{(\nu)}} \{ \log L_C(\theta | \mathbf{Y}, \mathbf{Z}) | \mathbf{Y} \} \\ &= \int \log L_C(\theta | \mathbf{Y}, \mathbf{z}) f_{\mathbf{Z}|\mathbf{Y}} \left( \mathbf{z} | \mathbf{Y}, \theta^{(\nu)} \right) dz \end{aligned}$$

2. M step: calculate  $\theta^{(\nu+1)}$  that maximizes  $Q \left( \theta, \theta^{(\nu)}, \mathbf{Y} \right)$  wrt  $\theta$

*The plausibility of EM*

**Lemma 1.1.** For  $\mathbf{Y} \sim f(y; \theta_0)$ ,  $E_{\theta_0} \{ \log f(\mathbf{Y}; \theta) \}$  is maximized at  $\theta = \theta_0$

*Proof:* Note that  $\psi(x) = -\log(x)$  is a convex function for  $x \in (0, \infty)$ , since  $\psi''(x) = x^{-2} > 0$ . Then by Jensen’s inequality we have

$$-\log \left[ E_{\theta_0} \left\{ \frac{f(\mathbf{Y}; \theta)}{f(\mathbf{Y}; \theta_0)} \right\} \right] \leq -E_{\theta_0} \left[ \log \left\{ \frac{f(\mathbf{Y}; \theta)}{f(\mathbf{Y}; \theta_0)} \right\} \right]$$

and since  $E_{\theta_0} \left\{ \frac{f(\mathbf{Y}; \theta)}{f(\mathbf{Y}; \theta_0)} \right\} = \int \frac{f(y; \theta)}{f(y; \theta_0)} f(y; \theta_0) dy = 1$  we get  $E_{\theta_0} \left[ \log \left\{ \frac{f(\mathbf{Y}; \theta)}{f(\mathbf{Y}; \theta_0)} \right\} \right] \leq 0$

## 1.2 Review on GLM

some materials: **on two forms, canonical forms, GLM and exponential family, Exponential family**: Y with distribution with the density of the form

$$Y \sim f(y; \theta, \phi) = \exp \left\{ \frac{y\theta - b(\theta)}{a(\phi)} + c(y, \phi) \right\}$$

where

- $\theta$  is the canonical (natural) parameter – often the parameter of interest
- $\Phi$  is dispersion (scale) parameter – often nuisance parameter

Exponential family also has **another** commonly used form **A good material here**

$$p(x | \eta) = h(x) \exp \{ \eta^T T(x) - A(\eta) \}$$

parameter vector  $\eta$ , often referred to as the *canonical parameter* for given functions  $T$  and  $h$ . The statistic  $T(X)$  is referred to as a *sufficient statistic*. The function  $A(\eta)$  is known as the *cumulant function* and

$$A(\eta) = \log \int h(x) \exp \{ \eta^T T(x) \} \nu(dx)$$

with respect to the measure  $\nu$

For example, for normal distribution,

$$\begin{aligned} f(y; \mu, \sigma^2) &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{(y - \mu)^2}{2\sigma^2} \right\} \\ &= \exp \left\{ \frac{y\mu - \mu^2/2}{\sigma^2} - \frac{y^2}{2\sigma^2} - \frac{1}{2} \log(2\pi\sigma^2) \right\} \end{aligned}$$

then

- $\theta = \mu$
- $\phi = \sigma^2$
- $a(\phi) = \phi$
- $b(\theta) = \frac{\theta^2}{2}$
- $c(y, \phi) = -\frac{y^2}{2\phi} - \frac{1}{2} \log(2\pi\phi)$

or

$$f(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma}} \exp \left\{ -\frac{1}{2\sigma^2}(x - \mu)^2 \right\} = \frac{1}{\sqrt{2\pi}} \exp \left\{ \frac{\mu}{\sigma^2}x - \frac{1}{2\sigma^2}x^2 - \frac{1}{2\sigma^2}\mu^2 - \log \sigma \right\}$$

then

$$\begin{aligned} \eta &= \begin{bmatrix} \mu/\sigma^2 \\ -1/2\sigma^2 \end{bmatrix} \\ T(x) &= \begin{bmatrix} x \\ x^2 \end{bmatrix} \\ A(\eta) &= \frac{\mu^2}{2\sigma^2} + \log \sigma = -\frac{\eta_1^2}{4\eta_2} - \frac{1}{2} \log(-2\eta_2) \\ h(x) &= \frac{1}{\sqrt{2\pi}} \end{aligned}$$

We also have the following property

Property 1.

$$\begin{aligned} E(Y) &= b'(\theta) \\ \text{Var}(Y) &= b''(\theta)a(\phi) \end{aligned}$$

**Proof**

$$\ell(\theta) = \log f(Y; \theta, \phi) = \frac{Y\theta - b(\theta)}{a(\phi)} + c(Y, \phi)$$

Then the score function for  $\theta$  is

$$U(\theta) = \frac{\partial \ell(\theta)}{\partial \theta} = \frac{Y - b'(\theta)}{a(\phi)}$$

Since  $E(U(\theta)) = 0$ , we have

$$E\left(\frac{Y - b'(\theta)}{a(\phi)}\right) = 0 \implies E(Y) = b'(\theta)$$

Since  $E(U(\theta)) = 0$ , we have

$$\text{Var}(U(\theta)) = E\left(\frac{\partial \ell(\theta)}{\partial \theta}\right)^2 = -E\left(\frac{\partial^2 \ell(\theta)}{\partial \theta^2}\right) \implies \frac{\text{Var}(Y)}{a^2(\phi)} = \frac{b''(\theta)}{a(\phi)} \implies \text{Var}(Y) = b''(\theta)a(\phi)$$

### 1.2.1 Newly introduced in 711

Exponential family distributions  $\{\mathcal{F}_\theta, \theta \in \Omega\}$  have densities of the general form

$$f(x; \theta) = h(x) \exp\left\{\sum_{i=1}^s g_i(\theta)T_i(x) - B(\theta)\right\}$$

this representation is not unique. A version with the smallest  $s$  is called *minimal exponential family* under this case both  $T_i$ 's and  $g_i$ 's are free of linear constraints, which means

$$\sum_{i=1}^s c_i T_i(x) = c_{s+1} \text{ a.e. } \nu \implies c_i = 0, i = 1, \dots, s+1$$

With minimal exponential family  $(T_1, \dots, T_s)$  is minimal sufficient. *Natural parameter space*

$$\Omega_0 = \left\{ \theta : \int h(x) \exp\left\{\sum_{i=1}^s g_i(\theta)T_i(x)\right\} d\nu(x) < \infty \right\}$$

*Canonical representation*

$$f(x; \eta) = h(x) \exp\left\{\sum_{i=1}^s \eta_i T_i(x) - A(\eta)\right\}$$

by integration equal to 1 we have

$$\int h(x) \exp\left\{\sum_{i=1}^s \eta_i T_i(x)\right\} d\nu(x) = \exp[A(\eta)]$$

with the canonical parameter  $\eta = g(\theta)$  and accordingly

$$\Omega_\eta = \left\{ \eta : \int h(x) \exp\left\{\sum_{i=1}^s \eta_i T_i(x)\right\} d\nu(x) < \infty \right\}$$

is another form of the natural parameter space.

If  $\eta$  is in the interior of  $\Omega_\eta$  then the m.g.f of  $T$  exists in a neighborhood of  $u = 0$

$$\begin{aligned} M_{\mathbf{T}}(\mathbf{u}) &= Ee^{u_1 T_1 + \dots + u_s T_s} \\ &= \int \exp\left\{\sum u_i T_i\right\} \cdot h(x) \exp\left\{\sum \eta_i T_i - A(\eta)\right\} d\nu(x) \\ &= \int h(x) \exp\left\{\sum (\eta_i + u_i) T_i - A(\eta)\right\} d\nu(x) \\ &= \int h(x) \frac{\exp\left\{\sum (\eta_i + u_i) T_i\right\}}{\exp\{A(\eta)\}} d\nu(x) \\ &= \frac{\exp\{A(\eta + \mathbf{u})\}}{\exp\{A(\eta)\}} = e^{A(\eta + \mathbf{u}) - A(\eta)} \end{aligned}$$

$$K_{\mathbf{T}}(\mathbf{u}) = \log \{M_{\mathbf{T}}(\mathbf{u})\} = A(\boldsymbol{\eta} + \mathbf{u}) - A(\boldsymbol{\eta})$$

therefore

$$\frac{\partial A(\boldsymbol{\eta})}{\partial \boldsymbol{\eta}^T} = E\{\mathbf{T}(X)\}, \quad \frac{\partial^2 A(\boldsymbol{\eta})}{\partial \boldsymbol{\eta} \partial \boldsymbol{\eta}^T} = \text{Var}\{\mathbf{T}(X)\}$$

also

$$\mathbf{I}(\boldsymbol{\eta}) = \frac{\partial^2 A(\boldsymbol{\eta})}{\partial \boldsymbol{\eta} \partial \boldsymbol{\eta}^T} = \text{Var}\{\mathbf{T}(X)\}$$

*curved exponential family:*  $\dim(\boldsymbol{\eta}) > \dim(\theta)$

*Completeness:* a set of statistics  $\mathbf{T} = (T_1, \dots, T_s)$  is complete with respect to the family of their induced distributions indexed by  $\theta$  if there are no functions  $\Phi(T)$  other than  $\Phi = 0$  such that  $E_{\theta} \Phi(\mathbf{T}) = 0$  for all  $\theta \in \Omega$

complete sufficient statistics  $\xrightarrow{\text{if}} \leftarrow$  minimal sufficient statistics

*Lehmann–Scheffe theorem:* If a statistic is unbiased, complete and sufficient, then it is the unique best unbiased estimator.

Sufficient statistic  $(T_1, \dots, T_s)$  of a canonical exponential family is complete sufficient provided that the family is minimal and that the parameter space contains an  $s$ -dimensional rectangle (which ruled out the curve exponential family).

### Distribution of sufficient statistics

$$X \sim \text{exponential family with } \boldsymbol{\eta} = (\boldsymbol{\zeta}, \boldsymbol{\psi}) : \\ f(x; \boldsymbol{\zeta}, \boldsymbol{\psi}) = h(x) \exp \left\{ \sum_{i=1}^r \zeta_i U_i(x) + \sum_{j=1}^s \psi_j T_j(x) - A(\boldsymbol{\zeta}, \boldsymbol{\psi}) \right\}$$

then

1. the distribution of  $(U, T)$  is an exponential family
2. marginally  $T$  is an exponential family

$$f(\mathbf{t}; \boldsymbol{\zeta}, \boldsymbol{\psi}) = q(\mathbf{t}) C(\boldsymbol{\zeta}) \exp \left\{ \sum_{j=1}^s \psi_j t_j - A(\boldsymbol{\zeta}, \boldsymbol{\psi}) \right\}$$

3. conditional distribution of  $U$  given  $T = t$  is an exponential family

$$f(\mathbf{u} \mid \mathbf{t}; \boldsymbol{\zeta}) = q_{\mathbf{t}}(\mathbf{u}) \exp \left\{ \sum_{i=1}^r \zeta_i u_i - A_{\mathbf{t}}(\boldsymbol{\zeta}) \right\}$$

## 1.3 Families of distributions, moments / cumulants, and quantiles / percentiles

### Families of distributions

- **Parametric** model or parametric family: a set of distributions indexed by a finite dimensional parameter

$$\boldsymbol{\theta} = (\theta_1, \dots, \theta_b)^T, \quad \{F(y; \boldsymbol{\theta}) : \boldsymbol{\theta} \in \Theta\}$$

- **Semiparametric** model: a set of distributions indexed by a finite dimensional  $\theta$  and some infinite dimensional parameters

for example

$$\{f(y; \mu, f_0) = f_0(y - \mu) : -\infty < \mu < \infty, f_0 \in C, \text{ where } C \text{ is a class of continuous unimodal densities}\}$$

where  $\mu$  belongs to the parametric part and the  $f_0$  belongs to the nonparametric part

- **Nonparametric** model: a set of distributions indexed by infinite dimensional parameters

for example

$$Y_i = g(X_i) + e_i, \quad \text{density of } e \in C$$

where  $g(\cdot)$  belongs to a class of functions with certain smoothness

- **$j$ th population moment** of random variable  $Y$

$$\mu'_j = E(Y^j), j = 1, 2, \dots$$

with its sample counterparts as

$$m'_j = n^{-1} \sum_{i=1}^n Y_i^j$$

- **jth population central moment** of rv  $Y$

$$\mu_j = E[\{Y - E(Y)\}^j]$$

with its sample counterparts as

$$m_j = n^{-1} \sum_{i=1}^n (Y_i - \bar{Y})^j$$

- **moment generating function (MGF)**

$$m(t) = E\{\exp(tY)\}$$

a useful trick is to write the MGF as

$$\begin{aligned} M(t) &= E(e^{tX}) = E(1 + tX + \dots + t^r X^r / r! + \dots) \\ &= \sum_{r=0}^{\infty} \mu_r t^r / r! \end{aligned}$$

which implies that

$$\mu'_r = M^{(r)}(0)$$

- **cumulant generating function (CGF)**

$$k(t) = \log\{m(t)\}$$

it's useful to remember that

$$K(t) = \log M(t) = \sum_r \kappa_r t^r / r!$$

which is just from the Taylor expansion for  $K(t)$  (note: the Taylor expansion around 0 is  $f(x) = \sum_{n=0}^{\infty} \frac{f^{(n)}(0)}{n!} x^n$ ) **a nice material on this topic**

- **jth population cumulant**

$$\kappa_j = k^{(j)}(0)$$

for example

$$\kappa_1 = \mu'_1 \text{ ( mean )}$$

$$\kappa_2 = \mu'_2 - \kappa_1^2 \text{ ( variance )}$$

*proof:*

$$\begin{aligned} \kappa_1 &= \kappa'(0) = \left. \frac{m'(t)}{m(t)} \right|_{t=0} \\ &= \left. \frac{E[\exp(tY)Y]}{E[\exp(tY)]} \right|_{t=0} \\ &= E[Y] = \mu'_1 \end{aligned}$$

and

$$\begin{aligned} \kappa''(t) &= \left. \frac{m''(t)m(t) - [m'(t)]^2}{m^2(t)} \right|_{t=0} \\ &= E(Y^2) - E(Y)^2 = \mu'_2 - \kappa_1^2 \end{aligned}$$

- **skewness**

$$\text{skew} = E \left[ \left( \frac{X - \mu}{\sigma} \right)^3 \right] = \frac{\mu_3}{\sigma^3} = \frac{E[(X - \mu)^3]}{(E[(X - \mu)^2])^{3/2}} = \mu_3 / \mu_2^{3/2}$$

which measures symmetry. Note that we are using the central moment in the definition

- **kurtosis**

$$\text{kurt} = E \left[ \left( \frac{X - \mu}{\sigma} \right)^4 \right] = \frac{E[(X - \mu)^4]}{(E[(X - \mu)^2])^2} = \mu_4 / \mu_2^2$$

which measures tail mass. Note that for Normal distribution skew = 0, kurt = 3

- **coefficient of variation**

$$CV = \frac{\sigma}{\mu} = \mu_2^{1/2} / \mu$$

for non-negative random variable, providing a relative standard deviation

- **quantile**

$$\eta_p \equiv \inf\{x : F(x) \geq p\} = F^{-1}(p)$$

the last equality holds only when  $F$  is a strictly monotone function and  $X$  is continuous

## 1.4 Transformation of random variables

**Jacobian method:** Let  $f_{X_1 X_2}(x_1, x_2)$  be the value of the joint probability density of the continuous random variables  $X_1$  and  $X_2$  at  $(x_1, x_2)$ . If the functions given by  $y_1 = u_1(x_1, x_2)$  and  $y_2 = u_2(x_1, x_2)$  are partially differentiable with respect to  $x_1$  and  $x_2$  and represent a one-to-one transformation for all values within the range of  $X_1$  and  $X_2$  for which  $f_{X_1 X_2}(x_1, x_2) \neq 0$ , then, for these values of  $x_1$  and  $x_2$ , the equations  $y_1 = u_1(x_1, x_2)$  and  $y_2 = u_2(x_1, x_2)$  can be uniquely solved for  $x_1$  and  $x_2$  to give  $x_1 = w_1(y_1, y_2)$  and  $x_2 = w_2(y_1, y_2)$  and for corresponding values of  $y_1$  and  $y_2$ , the joint probability density of  $Y_1 = u_1(X_1, X_2)$  and  $Y_2 = u_2(X_1, X_2)$  is given by

$$f_{Y_1 Y_2}(y_1, y_2) = f_{X_1 X_2}[w_1(y_1, y_2), w_2(y_1, y_2)] \cdot |J|$$

where  $J$  is the Jacobian of the transformation

$$J = \begin{vmatrix} \frac{\partial x_1}{\partial y_1} & \frac{\partial x_1}{\partial y_2} \\ \frac{\partial x_2}{\partial y_1} & \frac{\partial x_2}{\partial y_2} \end{vmatrix}$$

all other points of  $f_{Y_1 Y_2}(y_1, y_2) = 0$

**Distribution function method:** suppose we have the distribution function for  $Y$  as  $F_Y(y; \theta) = P(Y \leq y)$ . But we want the distribution function of  $X = g(Y)$ . Then

$$F_X(x; \theta) = P\{X \leq x\} = P\{g(Y) \leq x\}$$

if  $g$  is a strictly increasing function which means  $g^{-1}$  exists then

$$F_X(x; \theta) = P\{g(Y) \leq x\} = P\{Y \leq g^{-1}(x)\} = F_Y\{g^{-1}(x); \theta\}$$

if functions are differentiable then

$$f_X(x; \theta) = f_Y\{g^{-1}(x); \theta\} \frac{dg^{-1}(x)}{dx}$$

other transformation methods

## 1.5 Others

**Theorem 1.2. Singular value decomposition**

Suppose  $A$  is an  $n \times p$  matrix of rank  $r$ , where  $r \leq \min(n, p)$ . There exists orthogonal matrices  $U_{p \times p}$  and  $V_{n \times n}$  such that

$$V^T A U = \begin{pmatrix} \Delta & 0 \\ 0 & 0 \end{pmatrix} \rightarrow A = V D U^T, \text{ where } D = \begin{pmatrix} \Delta & 0 \\ 0 & 0 \end{pmatrix}$$

where  $\Delta = \text{diag}(\delta_1, \dots, \delta_r)$  is an  $r \times r$  diggonal matrix with  $\delta_1 \geq \delta_2 \dots \geq \delta_r > 0$ . The  $\delta_i$  are called the singular values of  $A$ .

# 2 Likelihood construction and estimation

## 2.1 Introduction

**Likelihood** is the joint density of the observed data.

For example

- in models for censored or missing data, the likelihood is not the density of the so-called “complete” data that includes the censored or missing values. Rather it is the density of only those components of the data that are observed and used in the statistical analysis
- consider an iid sample  $Y_1, \dots, Y_n$ , and a parametric transformation  $h(y, \alpha)$  strictly increasing in  $y$  for each  $\alpha$ . The model assumption is  $h(Y_1, \alpha), \dots, h(Y_n, \alpha)$  are iid with common density  $f(y; \theta)$  and distribution function  $F(y; \theta)$ , where  $\alpha$  and  $\theta$  are both parameters in the model. To construct the likelihood, we need to find the likelihood of the observed data  $Y$  rather than  $h(Y)$ . The distribution of  $Y_i$  is (we can also use Jacobian method here)

$$P\{Y_i \leq y\} = P\{h(Y_i, \alpha) \leq h(y, \alpha)\} = F\{h(y, \alpha); \theta\}$$

taking derivative wrt  $y$

$$f_Y(y; \theta, \alpha) = f\{h(y, \alpha); \theta\} \frac{\partial h(y, \alpha)}{\partial y}$$



thus the likelihood is

$$L(\boldsymbol{\theta}, \alpha; \mathbf{Y}) = \prod_{i=1}^n f\{h(Y_i, \alpha); \boldsymbol{\theta}\} \left\{ \frac{\partial h(y, \alpha)}{\partial y} \Big|_{y=Y_i} \right\}$$

a common mistake is missing the  $\left\{ \frac{\partial h(y, \alpha)}{\partial y} \Big|_{y=Y_i} \right\}$  part here

$\hat{\boldsymbol{\theta}}_{\text{MLE}}$  of size  $b \times 1$  is derived by maximizing

$$\ell(\boldsymbol{\theta}) = \log\{L(\boldsymbol{\theta} | \mathbf{Y})\}$$

under differentiable assumptions we can also solve for

$$\mathbf{S}(\boldsymbol{\theta}) = \{\ell'(\boldsymbol{\theta})\}^T$$

in vector calculus notation

$$\ell'(\boldsymbol{\theta}) = \frac{\partial \ell}{\partial \boldsymbol{\theta}} = \left( \frac{\partial \ell}{\partial \theta_1}, \dots, \frac{\partial \ell}{\partial \theta_b} \right)$$

$$\ell'(\boldsymbol{\theta})^T = \left( \frac{\partial \ell}{\partial \boldsymbol{\theta}} \right)^T = \frac{\partial \ell}{\partial \boldsymbol{\theta}^T} = \begin{pmatrix} \frac{\partial \ell}{\partial \theta_1} \\ \vdots \\ \frac{\partial \ell}{\partial \theta_b} \end{pmatrix}$$

and

$$\ell''(\boldsymbol{\theta}) = \frac{\partial}{\partial \boldsymbol{\theta}} \left( \frac{\partial \ell}{\partial \boldsymbol{\theta}^T} \right) = \frac{\partial^2 \ell}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} = \begin{pmatrix} \frac{\partial^2 \ell}{\partial \theta_1^2} & \cdots & \frac{\partial^2 \ell}{\partial \theta_b \partial \theta_1} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 \ell}{\partial \theta_1 \partial \theta_b} & \cdots & \frac{\partial^2 \ell}{\partial \theta_b^2} \end{pmatrix} = \left( \frac{\partial^2 \ell}{\partial \theta_j \partial \theta_i} \right)_{i=1, \dots, b, j=1, \dots, b}$$

## 2.2 Likelihood construction

- For discrete IID case, multinomial distribution is widely used cause
  1. any discrete distribution with a finite support  $\implies$  multinomial model
  2. any discrete distribution with an infinite support (e.g. Poisson)  $\xrightarrow{\text{grouped}}$  multinomial model

$$N \sim \text{binomial}(n, p), f(N; p) = \binom{n}{N} p^N (1-p)^{n-N}$$

more generally

$$(N_1, \dots, N_k) \sim \text{multinomial}(n; p_1, \dots, p_k), k \geq 2, f(N_1, \dots, N_k; p_1, \dots, p_k) = \frac{n!}{N_1! \dots N_k!} p_1^{N_1} \dots p_k^{N_k}$$

where  $\sum_{i=1}^k p_i = 1, \sum_{i=1}^k N_i = n$ , and the model can be interpreted as tossing  $n$  balls into  $k$  urns. Also

$$E(N_i) = np_i$$

$$\text{var}(N_i) = np_i(1-p_i)$$

$$\text{cov}(N_i, N_j) = -np_i p_j, \quad i \neq j$$

*hint:* to prove the covariance, we can write  $N_i = \sum_k M_{ki}$  where  $M_{ki}$  means toss  $k$ -th ball into  $i$ -th urn.

- Continuous IID case: skipped (nothing very important here)
- Connection Between Discrete and Continuous Likelihoods: **2h-method**

1. basic formula for continuous variable:

$$f(y) = \lim_{h \rightarrow 0^+} \frac{F(y+h) - F(y-h)}{2h} = \lim_{h \rightarrow 0^+} \frac{P(Y \in (y-h, y+h])}{2h}$$

2. bivariate data with both  $X$  and  $Y$  being continuous

$$f_{X,Y}(x, y) = \lim_{h \rightarrow 0^+} (2h)^{-2} \{F_{X,Y}(x+h, y+h) - F_{X,Y}(x-h, y+h)$$

$$- F_{X,Y}(x+h, y-h) + F_{X,Y}(x-h, y-h)$$

$$= \lim_{h \rightarrow 0^+} \frac{P(X \in (x-h, x+h], Y \in (y-h, y+h])}{(2h)^2}$$

### 3. bivariate data with $X$ discrete and $Y$ continuous

$$f_{X,Y}(x,y) = \lim_{h \rightarrow 0^+} \frac{P(X \in (x-h, x+h], Y \in (y-h, y+h])}{2h} = \lim_{h \rightarrow 0^+} \frac{P(X = x, Y \in (y-h, y+h])}{2h}$$

note: it is  $2h$  in the denominator

The likelihood based on the  $2h$ -method can be summarized as follows

(a) for continuous case

$$\begin{aligned} L(\boldsymbol{\theta} | \mathbf{Y}) &= \prod_{i=1}^n f(Y_i; \boldsymbol{\theta}) \\ &= \prod_{i=1}^n \lim_{h \rightarrow 0^+} \frac{F(Y_i + h; \boldsymbol{\theta}) - F(Y_i - h; \boldsymbol{\theta})}{2h} \\ &= \lim_{h \rightarrow 0^+} \left( \frac{1}{2h} \right)^n \prod_{i=1}^n \{F(Y_i + h; \boldsymbol{\theta}) - F(Y_i - h; \boldsymbol{\theta})\} \end{aligned}$$

(b) for discrete case

$$\begin{aligned} \lim_{h \rightarrow 0^+} \prod_{i=1}^n \{F(Y_i + h; \boldsymbol{\theta}) - F(Y_i - h; \boldsymbol{\theta})\} &= \prod_{i=1}^n \{F(Y_i^+; \boldsymbol{\theta}) - F(Y_i^-; \boldsymbol{\theta})\} \\ &= \prod_{i=1}^n f(Y_i; \boldsymbol{\theta}) \\ &= L(\boldsymbol{\theta} | \mathbf{Y}) \end{aligned}$$

(c) for combination of continuous and discrete random variables

$$L(\boldsymbol{\theta} | \mathbf{Y}) = \lim_{h \rightarrow 0^+} \left( \frac{1}{2h} \right)^m \prod_{i=1}^n \{F_i(Y_i + h; \boldsymbol{\theta}) - F_i(Y_i - h; \boldsymbol{\theta})\}$$

where  $1 \leq m \leq n$  depends on the number of continuous components in the data.

## 2.3 Proportional likelihoods

Likelihoods are equivalent if they are proportional and the constant of proportionality does not depend on unknown parameters.

- transformation of variables

$$Y_1, \dots, Y_n \sim \text{iid } f_Y(y; \boldsymbol{\theta})$$

we'd like to focus on transformed data  $X_i = g(Y_i)$ ,  $g(\cdot)$  known, increasing, and continuously differentiable

$$\begin{aligned} f_X(x; \boldsymbol{\theta}) &= f_Y\{h(x); \boldsymbol{\theta}\} h'(x), \quad h(\cdot) = g^{-1}(\cdot) \\ L(\boldsymbol{\theta} | \mathbf{X}) &= L(\boldsymbol{\theta} | \mathbf{Y}) \prod_{i=1}^n h'(X_i) = L(\boldsymbol{\theta} | \mathbf{Y}) \prod_{i=1}^n \frac{1}{g'(Y_i)} \end{aligned}$$

note that  $\prod_{i=1}^n \frac{1}{g'(Y_i)}$  is constant, thus the estimate of  $\boldsymbol{\theta}$  will be the same

- sufficient statistic e.g.,  $Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} \text{Bernoulli}(p)$

$$L(p | \mathbf{Y}) = \prod_{i=1}^n p^{Y_i} (1-p)^{1-Y_i} = p^S (1-p)^{n-S}$$

for  $S = \sum_{i=1}^n Y_i$ , which is a sufficient statistic

$$L(p | S) = \binom{n}{S} p^S (1-p)^{n-S} \propto L(p | \mathbf{Y})$$

- different sampling plans

the above two examples give the same sampling plan and data for the two likelihoods. Here we give an example that different sampling plan leads to proportional likelihoods.

1. as in the sufficient statistic example

$$Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} \text{Bernoulli}(p)$$

which leads to the first likelihood

$$L_1(p | S) = \binom{12}{S} p^S (1-p)^{12-S}$$

2. negative binomial:  $Y_i$ 's are observed until 3 0's appear, which leads to likelihood

$$L_2(p | S) = \binom{S+2}{S} p^S (1-p)^3$$

(it's  $S+2$  cause the last observation must be 0)

the ratio is

$$\frac{\binom{12}{S}}{\binom{S+2}{S}} (1-p)^{9-S}$$

which depends on  $p$  except for  $S=9$ . When  $S=9$ , the MLE estimate for  $p$  will be the same but the inference like hypothesis testing (p-value) is **different** cause null distribution in the two plans are different

## 2.4 Empirical distribution function as an MLE

$$Y_1, \dots, Y_n \sim \text{iid } F(y)$$

$$\begin{aligned} L_h(F | \mathbf{Y}) &= \prod_{i=1}^n \{F(Y_i + h) - F(Y_i - h)\} \\ &= \prod_{i=1}^n p_{i,h} \text{ assuming no ties} \end{aligned}$$

for the likelihood here we ignored the  $(2h)^{-m}$  factor and the  $p_{i,h}$  here needs to satisfy that  $p_{i,h} \geq 0$ ,  $\sum_{i=1}^n p_{i,h} \leq 1$ . Obviously, the likelihood increases as  $p_{i,h}$  increases until it reaches that

$$\sum_{i=1}^n p_{i,h} = 1$$

Thus use Lagrange multipliers, we solve for

$$\log L_h(F | \mathbf{Y}) + \lambda \left( \sum_{i=1}^n p_{i,h} - 1 \right) = \sum_{i=1}^n \log p_{i,h} + \lambda \left( \sum_{i=1}^n p_{i,h} - 1 \right)$$

then

$$\begin{aligned} \frac{\partial g}{\partial p_{i,h}} &= \frac{1}{p_{i,h}} + \lambda = 0, \quad i = 1, \dots, n \\ \frac{\partial g}{\partial \lambda} &= \sum_{i=1}^n p_{i,h} - 1 = 0 \end{aligned}$$

we can show that

$$\hat{F}_h \xrightarrow{d} \hat{F}_{\text{EMP}}(y) = n^{-1} \sum_{i=1}^n I(Y_i \leq y)$$

therefore we take the MLE of  $F(y)$  as the empirical distribution function

$$\hat{F}_{\text{MLE}}(y) = F_{\text{EMP}}(y) = \frac{1}{n} \sum_{i=1}^n I(Y_i \leq t)$$

## 2.5 Likelihoods for censored / truncated data

- fixed censoring (for censoring, we know exactly which individuals are censored) e.g. suppose we have underlying random variable  $X \sim N(\mu, \sigma^2)$  and we generate  $Y$  as follows

$$Y = \begin{cases} 0 & X \leq 0 \\ X & X > 0 \end{cases}$$

$$F(y) = \begin{cases} 0 & y < 0 \\ \Phi\{(y - \mu)/\sigma\} & y \geq 0 \end{cases}$$

then the likelihood is

$$Y_1, \dots, Y_n \sim \text{iid}$$

$$L(\mu, \sigma \mid \mathbf{Y}) = \left\{ \prod_{i:Y_i=0} \Phi(-\mu/\sigma) \right\} \left[ \prod_{i:Y_i>0} \sigma^{-1} \phi\{(Y_i - \mu)/\sigma\} \right]$$

- truncation (unlike censoring, here we are unaware of the truncated individuals) e.g. suppose income  $X \sim F_X(x; \theta)$  and sample  $Y_1, \dots, Y_n$  comes from incomes above  $L_0$  (left truncated). Then

$$\begin{aligned} P(Y \leq y) &= P(X \leq y \mid X > L_0) I(y > L_0) \\ &= \frac{P(X \leq y, X > L_0)}{P(X > L_0)} I(y > L_0) \\ &= \frac{F_X(y; \theta) - F_X(L_0; \theta)}{1 - F_X(L_0; \theta)} I(y > L_0) \\ f_Y(y) &= \frac{f_X(y; \theta)}{1 - F_X(L_0; \theta)} I(y > L_0) \end{aligned}$$

- random censoring: we have underlying random variables  $X$  and  $R$  (e.g.  $X$ : survival time,  $R$ : censoring time) and  $X \perp R$ . What we observed is  $Y = \min(X, R)$  and  $\delta = I(X \leq R)$  then

$$\begin{aligned} f_{Y,\delta}(y, \delta = 1) &= f_X(y; \theta) \{1 - F_R(y)\} \\ f_{Y,\delta}(y, \delta = 0) &= \{1 - F_X(y; \theta)\} f_R(y) \\ L(\theta \mid \mathbf{Y}, \delta) &= \prod_{i=1}^n f_X(Y_i; \theta)^{\delta_i} \{1 - F_R(Y_i)\}^{\delta_i} \times \{1 - F_X(Y_i; \theta)\}^{1-\delta_i} f_R(Y_i)^{1-\delta_i} \\ &\propto \prod_{i=1}^n f_X(Y_i; \theta)^{\delta_i} \{1 - F_X(Y_i; \theta)\}^{1-\delta_i} \\ &\text{if } F_R(\cdot) \text{ is noninformative of } \theta \end{aligned}$$

proof: we can use the 2h-method to get the density in this case.

$$\begin{aligned} \lim_{h \rightarrow 0} \frac{P(Y \in (y - h, y + h], \delta = 1)}{2h} &= \lim_{h \rightarrow 0} \frac{P(x \in (y - h, y + h), R > x)}{2h} \\ &= \lim_{h \rightarrow 0} \frac{\iint I(x \in (y - h, y + h], r > x) f_X(x) f_R(r) dx dr}{2h} \\ &= \lim_{h \rightarrow 0} \frac{\iint [I(r > x) f_R(r)] dr I(x \in (y - h, y + h]) f_X(x) dx}{2h} \\ &= \lim_{h \rightarrow 0} \frac{\iint [1 - F_R(x)] I(x \in (y - h, y + h]) f_X(x) dx}{2h} \\ &= \lim_{h \rightarrow 0} \frac{\int_{y-h}^{y+h} \{1 - F_R(x)\} f_X(x) dx}{2h} = [1 - F_R(y)] f_X(y) \end{aligned}$$

## 2.6 Likelihoods for regression models

- normal linear model

$$Y_i = \mathbf{x}_i^T \boldsymbol{\beta} + e_i, \quad i = 1, \dots, n, \quad e_1, \dots, e_n \stackrel{\text{iid}}{\sim} N(0, \sigma^2) \quad (1)$$

- fixed design:  $x_1, x_2, \dots, x_n$  are known nonrandom p-vectors

$$L(\boldsymbol{\beta}, \sigma \mid \{Y_i, \mathbf{x}_i\}_{i=1}^n) = \left( \frac{1}{\sqrt{2\pi}\sigma} \right)^n \exp \left\{ - \sum_{i=1}^n \frac{(Y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2}{2\sigma^2} \right\}$$

$$\hat{\beta}_{\text{MLE}} = \hat{\beta}_{\text{LS}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

where  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^T$ ,  $\mathbf{Y} = (Y_1, \dots, Y_n)^T$  and

$$\hat{\sigma}_{\text{MLE}}^2 = n^{-1} \sum_{i=1}^n \hat{e}_i^2 = n^{-1} \sum_{i=1}^n \left( Y_i - \mathbf{x}_i^T \hat{\beta}_{\text{MLE}} \right)^2$$

whereas

$$\hat{\sigma}_{\text{unbiased}} = (n - p)^{-1} \sum_{i=1}^n \hat{e}_i^2$$

2. random design:  $\{Y_i, \mathbf{x}_i\}_{i=1}^n \stackrel{\text{iid}}{\sim} (Y, \mathbf{x})$  and marginally  $\mathbf{x} \sim f_{\mathbf{x}}(\mathbf{x}; \boldsymbol{\tau})$  then

$$L(\boldsymbol{\beta}, \sigma, \boldsymbol{\tau} \mid \{Y_i, \mathbf{x}_i\}_{i=1}^n) = \text{above} \times \prod_{i=1}^n f_{\mathbf{x}}(\mathbf{x}_i; \boldsymbol{\tau}) = \left( \frac{1}{\sqrt{2\pi}\sigma} \right)^n \exp \left\{ - \sum_{i=1}^n \frac{(Y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2}{2\sigma^2} \right\} \times \prod_{i=1}^n f_{\mathbf{x}}(\mathbf{x}_i; \boldsymbol{\tau})$$

- additive errors nonlinear model: it is very similar to 1 but with the following the modification

$$Y_i = g(\mathbf{x}_i, \boldsymbol{\beta}) + e_i$$

with function  $g$  known. Then

$$\hat{\sigma}_{\text{MLE}}^2 = n^{-1} \sum_{i=1}^n \hat{e}_i^2 = n^{-1} \sum_{i=1}^n \left( Y_i - \mathbf{x}_i^T \hat{\beta}_{\text{MLE}} \right)^2$$

but the estimator for  $\boldsymbol{\beta}$  tend to have no closed form and need to solve numerically.

- generalized linear model

**Exponential family** is of the following form:

$$f(y|\theta, \phi) = \exp \left\{ \frac{y\theta - b(\theta)}{a(\phi)} + c(y, \phi) \right\}$$

where  $\theta$  is the *canonical parameter*,  $\phi$  is the *dispersion parameter*,  $b(\theta)$  is the *cumulative function* (different from cumulative generating function) and  $a(\cdot), b(\cdot), c(\cdot)$  are known functions.

To fit into the regression framework, we have

$$\log f(y_i; \theta_i, \phi) = \frac{y_i \theta_i - b(\theta_i)}{a_i(\phi)} + c(y_i, \phi)$$

link function is  $g(\cdot)$  such that

$$g(\mu_i) = \mathbf{x}_i^T \boldsymbol{\beta} \implies \mu_i = g^{-1}(\mathbf{x}_i^T \boldsymbol{\beta}) \implies \theta_i = b'^{-1} \{g^{-1}(\mathbf{x}_i^T \boldsymbol{\beta})\}$$

a special case of the link function is the *canonical link*

$$g(\mu_i) = \theta_i = \mathbf{x}_i^T \boldsymbol{\beta}$$

in this setting the log-likelihood becomes

$$\log L(\boldsymbol{\beta}, \phi \mid \{Y_i, \mathbf{x}_i\}_{i=1}^n) = \sum_{i=1}^n \left\{ \frac{y_i \mathbf{x}_i^T \boldsymbol{\beta} - b(\mathbf{x}_i^T \boldsymbol{\beta})}{a_i(\phi)} + c(y_i, \phi) \right\}$$

If  $Y \sim f(y; \theta, \phi)$  and  $\phi$  is fixed, then

Property 1.  $E(Y) = b'(\theta)$

Property 2.  $\text{Var}(Y) = b''(\theta)a(\phi)$

**Proof**

We first prove that

$$E[l'(\theta)] = 0$$

and

$$E[l''(\theta)] = -E\{[l'(\theta)]^2\}$$

$$\begin{aligned}\ell(\theta) &= \log f(y; \theta, \phi) = \frac{y\theta - b(\theta)}{a(\phi)} + c(y, \phi) \\ E \left[ \frac{\partial}{\partial \theta} \log f(y; \theta, \phi) \right] \\ &= \int \left\{ \frac{\partial}{\partial \theta} \log f(y; \theta, \phi) \right\} f(y; \theta, \phi) dy \\ &= \int f'(y; \theta, \phi) dy \\ &\stackrel{\text{under some regularity conditions}}{=} \frac{\partial}{\partial \theta} \int f(y; \theta, \phi) dy = 0\end{aligned}$$

and

$$\begin{aligned}E \left[ -\frac{\partial^2}{\partial \theta^2} \log f(y; \theta, \phi) \right] \\ &= - \int \left\{ \frac{f''(y; \theta, \phi) \cdot f(y; \theta, \phi) - [f'(y; \theta, \phi)]^2}{f(y; \theta, \phi)^2} \right\} \cdot f(y; \theta, \phi) dy \\ &= - \int f''(y; \theta, \phi) dy + \int \left[ \frac{f'(y; \theta, \phi)}{f(y; \theta, \phi)} \right]^2 dy \\ &= - \frac{\partial^2}{\partial \theta^2} \int f(y; \theta, \phi) dy + E \left[ [(\log f(y; \theta, \phi))']^2 \right] \\ &= E \left\{ [(\log f(y; \theta, \phi))']^2 \right\}\end{aligned}$$

Then the score function for  $\theta$  is

$$U(\theta) = \frac{\partial \ell(\theta)}{\partial \theta} = \frac{Y - b'(\theta)}{a(\phi)}$$

Since  $E(U(\theta)) = 0$ , we have

$$E \left( \frac{Y - b'(\theta)}{a(\phi)} \right) = 0 \implies E(Y) = b'(\theta)$$

Since  $E(U(\theta)) = 0$ , we have

$$\text{Var}(U(\theta)) = E \left( \frac{\partial \ell(\theta)}{\partial \theta} \right)^2 = -E \left( \frac{\partial^2 \ell(\theta)}{\partial \theta^2} \right) \implies \frac{\text{Var}(Y)}{a^2(\phi)} = \frac{b''(\theta)}{a(\phi)} \implies \text{Var}(Y) = b''(\theta)a(\phi)$$

- generalized linear mixed model (GLMM)

in this model we extend the canonical-link GLM to accommodate random effects

$$\theta_i = \mathbf{x}_i^T \boldsymbol{\beta} + \mathbf{z}_i^T \mathbf{U}, \quad \mathbf{U} \sim f_{\mathbf{U}}(\mathbf{u}; \boldsymbol{\nu})$$

then the likelihood is

$$L(\boldsymbol{\beta}, \phi, \boldsymbol{\nu} \mid \{Y_i, \mathbf{x}_i, \mathbf{z}_i\}_{i=1}^n) = \prod_{i=1}^n \int f_{Y_i|\mathbf{U}}(Y_i \mid \mathbf{u}, \mathbf{x}_i, \mathbf{z}_i, \boldsymbol{\beta}, \phi) f_{\mathbf{U}}(\mathbf{u}; \boldsymbol{\nu}) d\mathbf{u}$$

- accelerated failure time model

$$\log T = \mathbf{x}^T \boldsymbol{\beta} + \sigma e, \quad T \perp R \mid \mathbf{x}$$

usually  $e \sim N(0, 1)$

the observed random variable is  $Y = \min(T, R)$ ,  $\delta = I(T \leq R)$  then the likelihood is

$$L(\boldsymbol{\beta}, \sigma \mid \{Y_i, \delta_i, \mathbf{x}_i\}_{i=1}^n) \propto \prod_{i=1}^n \left\{ \frac{1}{\sigma} f_e \left( \frac{\log Y_i - \mathbf{x}_i^T \boldsymbol{\beta}}{\sigma} \right) \right\}^{\delta_i} \left\{ 1 - F_e \left( \frac{\log Y_i - \mathbf{x}_i^T \boldsymbol{\beta}}{\sigma} \right) \right\}^{1-\delta_i}$$

## 2.7 Marginal and conditional likelihoods

Suppose  $\mathbf{Y} \iff (\mathbf{W}, \mathbf{V})$  and we have

$$\begin{cases} \theta_1: \text{parameter of interest} \\ \theta_2: \text{nuisance parameter} \end{cases} \quad (2)$$

$$\begin{aligned} f_{\mathbf{Y}}(\mathbf{y}; \boldsymbol{\theta}_1, \boldsymbol{\theta}_2) &= f_{\mathbf{W}, \mathbf{V}}(\mathbf{w}, \mathbf{v}; \boldsymbol{\theta}_1, \boldsymbol{\theta}_2) \\ &= f_{\mathbf{W}|\mathbf{V}}(\mathbf{w} | \mathbf{V}; \boldsymbol{\theta}_1, \boldsymbol{\theta}_2) f_{\mathbf{V}}(\mathbf{v}; \boldsymbol{\theta}_1, \boldsymbol{\theta}_2) \end{aligned} \quad (3)$$

1. if  $f_{\mathbf{V}}(\mathbf{v}; \boldsymbol{\theta}_1, \boldsymbol{\theta}_2) = f_{\mathbf{V}}(\mathbf{v}; \boldsymbol{\theta}_1)$ , it is a marginal likelihood
2. if  $f_{\mathbf{W}|\mathbf{V}}(\mathbf{w} | \mathbf{v}; \boldsymbol{\theta}_1, \boldsymbol{\theta}_2) = f_{\mathbf{W}|\mathbf{V}}(\mathbf{w} | \mathbf{v}; \boldsymbol{\theta}_1)$ , it is a conditional likelihood

Let's illustrate with an example

*Neyman–Scott problem*

$$\begin{aligned} Y_{ij} &\stackrel{i.i.d.}{\sim} N(\mu_i, \sigma^2), i = 1, \dots, n, j = 1, 2 \\ \boldsymbol{\theta} &= (\sigma^2, \mu_1, \dots, \mu_n)^T \text{ dim } n + 1 \text{ (increases with sample size)} \end{aligned}$$

then we have

$$\begin{aligned} \hat{\sigma}_{\text{MLE}}^2 &= (2n)^{-1} \sum_{i=1}^n \sum_{j=1}^2 (Y_{ij} - \hat{\mu}_{i, \text{MLE}})^2 \\ &= \frac{n^{-1} \sum_{i=1}^n (Y_{i1} - Y_{i2})^2}{4} \\ E(\hat{\sigma}_{\text{MLE}}^2) &= \frac{\sigma^2}{2} \implies \hat{\sigma}_{\text{MLE}}^2 \xrightarrow{p} \frac{\sigma^2}{2} \text{ by WLLN} \end{aligned}$$

as we can see, consistency no longer exist for MLE in this case cause the number of parameters increases with sample size. Then how to achieve consistency?

$$\begin{aligned} V_i &= \frac{(Y_{i1} - Y_{i2})}{\sqrt{2}} \sim N(0, \sigma^2) \\ W_i &= \frac{(Y_{i1} + Y_{i2})}{\sqrt{2}} \sim N(\sqrt{2}\mu_i, \sigma^2) \end{aligned}$$

1. marginal likelihood

$$\begin{aligned} L(\sigma | \mathbf{V}) &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{V_i^2}{2\sigma^2}\right) \\ \hat{\sigma}_{\text{MMLE}}^2 &= n^{-1} \sum_{i=1}^n V_i^2 = \frac{n^{-1} \sum_{i=1}^n (Y_{i1} - Y_{i2})^2}{2} \end{aligned}$$

in this case, only  $V_i$  needs to be identified. However,  $W_i$  may provide a sense of information loss

2. conditional likelihood

the key for this approach is to identify sufficient statistic for nuisance parameter under the assumption that the parameter of interest is known. In our case,  $T_i = Y_{i1} + Y_{i2}$  is sufficient for  $\mu_i \implies (Y_{i1}, Y_{i2}) | T_i$  doesn't depend on  $\mu_i$

$$\begin{aligned} Y_{i1} | T_i &\sim N\left(\frac{T_i}{2}, \frac{\sigma^2}{2}\right) \\ \hat{\sigma}_{\text{CMLE}}^2 &= 2n^{-1} \sum_{i=1}^n \left(Y_{i1} - \frac{T_i}{2}\right)^2 \\ &= \frac{n^{-1} \sum_{i=1}^n (Y_{i1} - Y_{i2})^2}{2} \end{aligned}$$

*Example1: Logistic regression measurement error model*

$$\begin{aligned} P(Y = 1 | X) &= \frac{\exp(\alpha + \beta X)}{1 + \exp(\alpha + \beta X)} \\ W &= X + U, \quad U \sim N(0, \sigma_U^2) \text{ with } \sigma_U^2 \text{ known} \end{aligned}$$

with  $U \perp Y$  and  $U \perp X$ ,  $\sigma_U^2$  is known

- data:  $(Y_i, W_i), i = 1, \dots, n$
- parameters of interest:  $\alpha, \beta$

then

$$f(Y_i, W_i | \alpha, \beta, X_i) = \frac{\exp\{Y_i(\alpha + \beta X_i)\}}{1 + \exp(\alpha + \beta X_i)} \frac{1}{\sqrt{2\pi}\sigma_U} \exp\left\{-\frac{(W_i - X_i)^2}{2\sigma_U^2}\right\}$$

$$= \frac{1}{\sqrt{2\pi}\sigma_U} \frac{1}{1 + \exp(\alpha + \beta X_i)} \times \exp\left\{\frac{(W_i + Y_i\sigma_U^2\beta)X_i}{\sigma_U^2} - \frac{X_i^2}{2\sigma_U^2} + \alpha Y_i - \frac{W_i^2}{2\sigma_U^2}\right\}$$

$W_i + Y_i\sigma_U^2\beta$  or equivalently  $T_i = W_i + (Y_i - 1/2)\sigma_U^2\beta$  is sufficient for  $X_i$ , assuming  $\alpha$  and  $\beta$  are known (Sufficient statistic depends on the parameter of interest). Then

$$P(Y_i = 1 | T_i) = \frac{\exp(\alpha + \beta T_i)}{1 + \exp(\alpha + \beta T_i)}$$

and the conditional likelihood score is based on

$$\frac{\partial}{\partial(\alpha, \beta)} \log P(Y_i | T_i) = \left\{ Y_i - \frac{\exp(\alpha + \beta T_i)}{1 + \exp(\alpha + \beta T_i)} \right\} \begin{pmatrix} 1 \\ T_i \end{pmatrix}$$

and the score is

$$\sum_{i=1}^n \left\{ Y_i - \frac{\exp(\alpha + \beta T_i)}{1 + \exp(\alpha + \beta T_i)} \right\} \begin{pmatrix} 1 \\ T_i \end{pmatrix} \Bigg|_{T_i = W_i + (Y_i - 1/2)\sigma_U^2\beta} = \mathbf{0}$$

Note: when taking the differentiation,  $T_i$  is treated as given but not as a function of  $\beta$

*Example2: exponential families with canonical parameter*

$$f(\mathbf{y}; \boldsymbol{\theta}_1, \boldsymbol{\theta}_2) = h(\mathbf{y}) \exp\left\{ \sum_i \theta_{1i} W_i + \sum_j \theta_{2j} V_j - A(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) \right\} \implies \mathbf{W} | \mathbf{V} \sim \text{exponential family indexed by } \boldsymbol{\theta}_1 \text{ only}$$

*Example3: conditional logistic regression*

$$\text{logit}\{P(Y_i = 1)\} = \mathbf{x}_i^T \boldsymbol{\beta}$$

$$P(Y_i = 1) = \frac{\exp(\mathbf{x}_i^T \boldsymbol{\beta})}{1 + \exp(\mathbf{x}_i^T \boldsymbol{\beta})}$$

$$L(\boldsymbol{\beta} | \mathbf{Y}, \mathbf{X}) = \prod_{i=1}^n \frac{\exp(\mathbf{x}_i^T \boldsymbol{\beta})^{Y_i}}{1 + \exp(\mathbf{x}_i^T \boldsymbol{\beta})} = \frac{\exp(\sum_{i=1}^n Y_i \mathbf{x}_i^T \boldsymbol{\beta})}{\prod_{i=1}^n \{1 + \exp(\mathbf{x}_i^T \boldsymbol{\beta})\}}$$

$T_j = \sum_{i=1}^n x_{ij} Y_i$  is sufficient for  $\beta_j, j = 1, \dots, p$ . If  $\beta_k$  is the parameter of interest then let  $W = T_k$  and  $\mathbf{V} = (T_1, \dots, T_{k-1}, T_{k+1}, \dots, T_p)^T$

$$P(W | \mathbf{V}) = \frac{c(T_1, \dots, T_p) \exp(\beta_k T_k)}{\sum_u c(T_1, \dots, T_{k-1}, u, T_{k+1}, \dots, T_p) \exp(\beta_k u)}$$

where  $P(W, V) = c(T_1, \dots, T_p) \exp(\beta_k T_k)$  and  $P(V) = \sum_u P(W, V) = \sum_u c(T_1, \dots, T_{k-1}, u, T_{k+1}, \dots, T_p) \exp(\beta_k u)$

## 2.8 MLE and information matrix

Recall that

$$\hat{\boldsymbol{\theta}}_{\text{MLE}} = \arg \max_{\boldsymbol{\theta} \in \boldsymbol{\Theta}} L(\boldsymbol{\theta} | \mathbf{Y})$$

$$\ell(\boldsymbol{\theta}) = \log L(\boldsymbol{\theta} | \mathbf{Y})$$

if  $\ell(\boldsymbol{\theta})$  is continuously differentiable, then likelihood score

$$\mathbf{S}(\boldsymbol{\theta}) = \mathbf{S}(\mathbf{Y}, \boldsymbol{\theta}) = \frac{\partial \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}^T} = \{\ell'(\boldsymbol{\theta})\}^T$$

exists and  $\hat{\boldsymbol{\theta}}_{\text{MLE}}$  satisfies  $\mathbf{S}(\hat{\boldsymbol{\theta}}_{\text{MLE}}) = \mathbf{0}$ . Under most of the cases,  $\hat{\boldsymbol{\theta}}_{\text{MLE}}$  is unique, or at least there is a principled strategy for choosing a single solution from among the possibly multiple values. But there are also special cases, for example

$$Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} f(y; \mu) = \begin{cases} e^{-(y-\mu)} & y > \mu \\ 0 & \text{otherwise} \end{cases}$$

any  $\mu < \min(Y_i)$  leads to likelihood to be 0.



Property 1. invariant property:

$$\hat{\boldsymbol{\theta}}_{\text{MLE}} \text{ for } \boldsymbol{\theta} \implies \hat{\boldsymbol{\tau}}_{\text{MLE}} = g(\hat{\boldsymbol{\theta}}_{\text{MLE}}) \text{ for } \boldsymbol{\tau} = g(\boldsymbol{\theta})$$

Property 2. asymptotic normality: under regularity conditions

$$\sqrt{n}(\hat{\boldsymbol{\theta}}_{\text{MLE}} - \boldsymbol{\theta}) \xrightarrow{d} N\{\mathbf{0}, \mathbf{I}(\boldsymbol{\theta})^{-1}\}, \quad \text{as } n \rightarrow \infty$$

where

$$\begin{aligned} \mathbf{I}(\boldsymbol{\theta}) &= E \left[ \left\{ \frac{\partial}{\partial \boldsymbol{\theta}^T} \log f(Y; \boldsymbol{\theta}) \right\}^{\otimes 2} \right] \\ &= \text{var} \left\{ \frac{\partial}{\partial \boldsymbol{\theta}^T} \log f(Y; \boldsymbol{\theta}) \right\} \\ &= -E \left\{ \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} \log f(Y; \boldsymbol{\theta}) \right\} \end{aligned}$$

How to estimate  $I(\theta)$ ?

$I(\theta)$  has 2 components:  $I(\cdot)$  and  $\theta$ . It is a natural choice for using  $\hat{\theta}_{\text{MLE}}$  to estimate  $\theta$ . But what about  $I(\cdot)$ ?

1. theoretical  $I(\cdot) \implies \mathbf{I}(\hat{\boldsymbol{\theta}}_{\text{MLE}})$  is the MLE for  $I(\theta)$  (invariant property)

2. estimated  $I(\cdot)$  version 1

$$\begin{aligned} \bar{\mathbf{I}}(\mathbf{Y}, \boldsymbol{\theta}) &= n^{-1} \sum_{i=1}^n \left\{ -\frac{\partial}{\partial \boldsymbol{\theta}} \mathbf{s}(Y_i, \boldsymbol{\theta}) \right\} \\ &= n^{-1} \sum_{i=1}^n \left\{ -\frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} \log f(Y_i; \boldsymbol{\theta}) \right\} \end{aligned}$$

3. estimated  $I(\cdot)$  version 2 (usually less efficient)

$$\begin{aligned} \bar{\mathbf{I}}^*(\mathbf{Y}, \boldsymbol{\theta}) &= n^{-1} \sum_{i=1}^n \mathbf{s}(Y_i, \boldsymbol{\theta})^{\otimes 2} \\ &= n^{-1} \sum_{i=1}^n \left\{ \frac{\partial}{\partial \boldsymbol{\theta}^T} \log f(Y_i; \boldsymbol{\theta}) \right\}^{\otimes 2} \end{aligned}$$

Note: when  $Y_i$  are independent but not identically distributed. The  $f(\cdot; \boldsymbol{\theta}) \implies f_i(\cdot; \boldsymbol{\theta})$

4. estimated  $I(\cdot)$  version 1: ( $Y_i$  are independent but not identically distributed)

$$\begin{aligned} \bar{\mathbf{I}}(\mathbf{Y}, \boldsymbol{\theta}) &= \frac{1}{n} \sum_{i=1}^n \left\{ -\frac{\partial}{\partial \boldsymbol{\theta}} \mathbf{s}_i(Y_i, \boldsymbol{\theta}) \right\} \\ &= \frac{1}{n} \sum_{i=1}^n \left\{ -\frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} \log f_i(Y_i; \boldsymbol{\theta}) \right\} \end{aligned}$$

the according average expected information matrix (also called the average Fisher information matrix)

$$\bar{\mathbf{I}}(\boldsymbol{\theta}) = E\{\bar{\mathbf{I}}(\mathbf{Y}, \boldsymbol{\theta})\} = \frac{1}{n} \sum_{i=1}^n E \left\{ -\frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} \log f_i(Y_i; \boldsymbol{\theta}) \right\} = \frac{1}{n} \sum_{i=1}^n E \left[ \left\{ \frac{\partial}{\partial \boldsymbol{\theta}^T} \log f_i(Y_i; \boldsymbol{\theta}) \right\} \left\{ \frac{\partial}{\partial \boldsymbol{\theta}} \log f_i(Y_i; \boldsymbol{\theta}) \right\} \right]$$

in the i.i.d case  $\bar{\mathbf{I}}(\boldsymbol{\theta}) = \mathbf{I}(\boldsymbol{\theta})$

5. estimated  $I(\cdot)$  version 2: ( $Y_i$  are independent but not identically distributed)

$$\begin{aligned} \bar{\mathbf{I}}^*(\mathbf{Y}, \boldsymbol{\theta}) &= \frac{1}{n} \sum_{i=1}^n \mathbf{s}_i(Y_i, \boldsymbol{\theta}) \mathbf{s}_i(Y_i, \boldsymbol{\theta})^T \\ &= \frac{1}{n} \sum_{i=1}^n \left\{ \frac{\partial}{\partial \boldsymbol{\theta}^T} \log f_i(Y_i; \boldsymbol{\theta}) \right\} \left\{ \frac{\partial}{\partial \boldsymbol{\theta}} \log f_i(Y_i; \boldsymbol{\theta}) \right\} \end{aligned}$$

Actually, we discussed two types of information here

1. the total information:

$$\mathbf{I}_T(\boldsymbol{\theta}) = -E \left\{ \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} \log L(\boldsymbol{\theta} | \mathbf{Y}) \right\}$$

2. the observed total information

$$\mathbf{I}_T(\mathbf{Y}, \boldsymbol{\theta}) = -\frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} \log L(\boldsymbol{\theta} | \mathbf{Y})$$

Example 1:  $N(\mu, \sigma^2)$

$$\log f(y; \mu, \sigma) = \text{constant} - \log \sigma - \frac{1}{2\sigma^2}(y - \mu)^2$$

$$\mathbf{I}(\mu, \sigma) = \frac{1}{\sigma^2} \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix}$$

more specifically

$$\ell'(\theta) = \frac{\partial}{\partial \theta} \log f(y; \mu, \sigma) = \left( -\frac{1}{\sigma} + \frac{(y - \mu)^2}{\sigma^3}, \frac{y - \mu}{\sigma^2} \right)^T$$

then

$$\text{var}(\ell'(\theta)) = \begin{bmatrix} \text{var}\left(-\frac{1}{\sigma} + \frac{(y - \mu)^2}{\sigma^3}\right) & \text{cov}\left(-\frac{1}{\sigma} + \frac{(y - \mu)^2}{\sigma^3}, \frac{y - \mu}{\sigma^2}\right) \\ \text{cov}\left(-\frac{1}{\sigma} + \frac{(y - \mu)^2}{\sigma^3}, \frac{y - \mu}{\sigma^2}\right) & \text{var}\left(\frac{y - \mu}{\sigma^2}\right) \end{bmatrix}$$

then by the results of normal moments we get the answer

Order	Non-central moment	Central moment
1	$\mu$	0
2	$\mu^2 + \sigma^2$	$\sigma^2$
3	$\mu^3 + 3\mu\sigma^2$	0
4	$\mu^4 + 6\mu^2\sigma^2 + 3\sigma^4$	$3\sigma^4$
5	$\mu^5 + 10\mu^3\sigma^2 + 15\mu\sigma^4$	0
6	$\mu^6 + 15\mu^4\sigma^2 + 45\mu^2\sigma^4 + 15\sigma^6$	$15\sigma^6$
7	$\mu^7 + 21\mu^5\sigma^2 + 105\mu^3\sigma^4 + 105\mu\sigma^6$	0
8	$\mu^8 + 28\mu^6\sigma^2 + 210\mu^4\sigma^4 + 420\mu^2\sigma^6 + 105\sigma^8$	$105\sigma^8$

Example 2: recall for the normal error regression models

$$\begin{aligned} L(\beta, \sigma | \{Y_i, \mathbf{x}_i\}_{i=1}^n) &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp \left\{ -\frac{(Y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2}{2\sigma^2} \right\} \\ &= \left( \frac{1}{\sqrt{2\pi}\sigma} \right)^n \exp \left\{ -\sum_{i=1}^n \frac{(Y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2}{2\sigma^2} \right\} \end{aligned}$$

then

$$\ell = -\log \sqrt{2\pi} - \log \sigma - \frac{(Y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2}{2\sigma^2}$$

it follows that

$$\ell' = \left( \frac{\partial \ell}{\partial \sigma}, \frac{\partial \ell}{\partial \boldsymbol{\beta}} \right) = \left( -\frac{1}{\sigma} + \frac{(Y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2}{\sigma^3}, \frac{\mathbf{x}_i^T (Y_i - \mathbf{x}_i^T \boldsymbol{\beta})}{\sigma^2} \right)$$

and

$$\ell'' = \begin{bmatrix} \frac{1}{\sigma^2} - 3 \frac{(Y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2}{\sigma^4} & -2 \frac{\mathbf{x}_i^T (Y_i - \mathbf{x}_i^T \boldsymbol{\beta})}{\sigma^3} \\ -2 \frac{\mathbf{x}_i^T (Y_i - \mathbf{x}_i^T \boldsymbol{\beta})}{\sigma^3} & -\frac{\mathbf{x}_i^T \mathbf{x}_i}{\sigma^2} \end{bmatrix}$$

taking expectation and organize we get

$$\bar{\mathbf{I}}(\beta, \sigma) = \frac{1}{\sigma^2} \begin{pmatrix} \mathbf{X}^T \mathbf{X} / n & 0 \\ 0 & 2 \end{pmatrix}$$

then invert  $\mathbf{I}_T(\beta, \sigma) = n\bar{\mathbf{I}}(\beta, \sigma)$  we get

$$\begin{aligned} \text{avar}(\hat{\beta}) &= (\mathbf{X}^T \mathbf{X})^{-1} \sigma^2 \\ \text{avar}(\hat{\sigma}) &= \frac{\sigma^2}{2n} \end{aligned}$$

Similarly, for  $Y_i = g(\mathbf{x}_i, \boldsymbol{\beta}) + e_i$

$$\bar{\mathbf{I}}(\boldsymbol{\beta}, \sigma) = \frac{1}{\sigma^2} \begin{pmatrix} \mathbf{G}^T \mathbf{G} / n & 0 \\ 0 & 2 \end{pmatrix}$$

where  $[\mathbf{G}]_{ij} = \partial g(\mathbf{x}_i, \boldsymbol{\beta}) / \partial \beta_j$

*Example 3: GLM with canonical link* ( $g(\mu_i) = \theta_i = \mathbf{x}_i^T \boldsymbol{\beta}$ )

$$\begin{aligned} \mathbf{S}(\boldsymbol{\beta}, \phi) &= \sum_{i=1}^n \frac{\{Y_i - b'(\mathbf{x}_i^T \boldsymbol{\beta})\} \mathbf{x}_i}{a_i(\phi)} = \sum_{i=1}^n \frac{(Y_i - \mu_i) \mathbf{x}_i}{a_i(\phi)} \\ -\frac{\partial}{\partial \boldsymbol{\beta}} \mathbf{S}(\boldsymbol{\beta}, \phi) &= \sum_{i=1}^n \frac{b''(\mathbf{x}_i^T \boldsymbol{\beta})}{a_i(\phi)} \mathbf{x}_i^{\otimes 2} = \sum_{i=1}^n \frac{\text{var}(Y_i)}{a_i(\phi)^2} \mathbf{x}_i^{\otimes 2} \\ &\equiv \mathbf{X}^T \mathbf{V} \mathbf{X} = n \bar{\mathbf{I}}(\mathbf{Y}, \boldsymbol{\beta}) = n \bar{\mathbf{I}}(\boldsymbol{\beta}) \end{aligned}$$

where  $\mathbf{V} = \text{diag}\{\text{Var}(Y_1)/a_1(\phi)^2, \dots, \text{Var}(Y_n)/a_n(\phi)^2\}$ . In this special case, average Fisher information is equal to average observed information (cause  $\mathbf{X}^T \mathbf{V} \mathbf{X}$  is irrespective of  $Y$ ). When  $\Phi$  is unknown, the full average Fisher information is

$$\bar{\mathbf{I}}(\boldsymbol{\beta}, \phi) = \begin{pmatrix} \mathbf{X}^T \mathbf{V} \mathbf{X} / n & 0 \\ 0 & \bar{I}(\phi) \end{pmatrix}$$

thus estimating  $\Phi$  does not increase the variability of  $\boldsymbol{\beta}$  estimation asymptotically in canonical link GLM models. Cause the diagonal of  $\bar{\mathbf{I}}(\boldsymbol{\beta}, \phi)$  is 0 and therefore we have

$$\bar{I}^{-1}(\boldsymbol{\beta}, \phi) = \begin{pmatrix} n [X^T V X]^{-1} & 0 \\ 0 & \bar{I}^{-1}(\phi) \end{pmatrix}$$

However, in general

$$\begin{pmatrix} I_{11} & \mathbf{I}_{12} \\ & \mathbf{I}_{22} \end{pmatrix}^{-1} = \begin{pmatrix} (I_{11} - \mathbf{I}_{12} \mathbf{I}_{22}^{-1} \mathbf{I}_{21})^{-1} & -I_{11}^{-1} \mathbf{I}_{12} (\mathbf{I}_{22} - \mathbf{I}_{21} I_{11}^{-1} \mathbf{I}_{12})^{-1} \\ & (\mathbf{I}_{22} - \mathbf{I}_{21} I_{11}^{-1} \mathbf{I}_{12})^{-1} \end{pmatrix}$$

$$\mathbf{I}_{12} \mathbf{I}_{22}^{-1} \mathbf{I}_{21} \geq 0 \implies (I_{11} - \mathbf{I}_{12} \mathbf{I}_{22}^{-1} \mathbf{I}_{21})^{-1} \geq I_{11}^{-1}$$

which means that when adding parameters to a model, the diagonal elements of the inverse information matrix are always no less than the corresponding elements of the simpler model unless  $I_{12} = 0$ . This is called *variance inflation*

### 2.8.1 Transformed and modeled parameters

If we have  $f(y; \boldsymbol{\theta})$  where  $\boldsymbol{\theta}$ : b-dimensional and therefore  $I_{\boldsymbol{\theta}} : b$ . We also have  $\boldsymbol{\beta}$ : s-dimensional and  $s \leq b$ ,  $\boldsymbol{\theta} = g(\boldsymbol{\beta})$  then

$$\begin{aligned} \ell(\boldsymbol{\beta}) &= \log f\{y; \mathbf{g}(\boldsymbol{\beta})\} \\ \frac{\partial}{\partial \boldsymbol{\beta}} \ell(\boldsymbol{\beta}) &= \frac{\partial}{\partial \boldsymbol{\theta}} \log f(y; \boldsymbol{\theta}) \Big|_{\boldsymbol{\theta}=\mathbf{g}(\boldsymbol{\beta})} \frac{\partial \mathbf{g}(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \\ \mathbf{I}_{\boldsymbol{\beta}}(\boldsymbol{\beta}) &= \left\{ \frac{\partial \mathbf{g}(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \right\}^T \mathbf{I}_{\boldsymbol{\theta}}\{\mathbf{g}(\boldsymbol{\beta})\} \frac{\partial \mathbf{g}(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \end{aligned}$$

*Example 1:* in the normal setting, let  $\boldsymbol{\theta} = (\mu, \sigma)^T$  and  $\boldsymbol{\beta} = (\mu, \sigma^2)^T$  then

$$\begin{aligned} \frac{\partial \mathbf{g}(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} &= \text{diag}\{1, 1/(2\sigma)\} \\ \mathbf{I}_{\boldsymbol{\beta}}(\boldsymbol{\beta}) &= \begin{pmatrix} 1 & 0 \\ 0 & 1/(2\sigma) \end{pmatrix} \frac{1}{\sigma^2} \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & 1/(2\sigma) \end{pmatrix} \\ &= \frac{1}{\sigma^2} \begin{pmatrix} 1 & 0 \\ 0 & 1/(2\sigma^2) \end{pmatrix} \end{aligned}$$

## 2.9 Methods for maximizing the likelihood

1. Profile likelihood: maximize in a sequential way and therefore achieve dimension reduction

$$\max_{\boldsymbol{\theta}_1, \boldsymbol{\theta}_2} L(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) = \max_{\boldsymbol{\theta}_1} L\left\{\boldsymbol{\theta}_1, \hat{\boldsymbol{\theta}}_2(\boldsymbol{\theta}_1)\right\}, \quad \hat{\boldsymbol{\theta}}_2(\boldsymbol{\theta}_1) = \arg \max_{\boldsymbol{\theta}_2 | \boldsymbol{\theta}_1} L(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2)$$

Example 1:  $Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} \text{Gamma}(\alpha, \beta)$

$$f_Y(y; \alpha, \beta) = \frac{1}{\Gamma(\alpha)\beta^\alpha} y^{\alpha-1} e^{-y/\beta}$$

$$\ell(\alpha, \beta) = -n \log \Gamma(\alpha) - n\alpha \log \beta + (\alpha - 1) \sum_i \log Y_i - \frac{\sum_i Y_i}{\beta}$$

$$\frac{\partial}{\partial \beta} \ell(\alpha, \beta) = -\frac{n\alpha}{\beta} + \frac{\sum_i Y_i}{\beta^2} \implies \hat{\beta}(\alpha) = \bar{Y}/\alpha$$

$$\begin{aligned} \ell\{\alpha, \hat{\beta}(\alpha)\} &= -n \log \Gamma(\alpha) - n\alpha(\log \bar{Y} - \log \alpha) \\ &\quad + (\alpha - 1) \sum_i \log Y_i - n\alpha \end{aligned}$$

## 2. Newton methods

$$\begin{aligned} \mathbf{0} = \mathbf{S}(\boldsymbol{\theta}) &\approx \mathbf{S}(\boldsymbol{\theta}^{(\nu)}) + \left\{ \frac{\partial}{\partial \boldsymbol{\theta}} \mathbf{S}(\boldsymbol{\theta}) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^{(\nu)}} \right\} (\boldsymbol{\theta} - \boldsymbol{\theta}^{(\nu)}) \\ &= \mathbf{S}(\boldsymbol{\theta}^{(\nu)}) - \mathbf{I}_T(\mathbf{Y}, \boldsymbol{\theta}^{(\nu)}) (\boldsymbol{\theta} - \boldsymbol{\theta}^{(\nu)}) \\ \boldsymbol{\theta}^{(\nu+1)} &= \boldsymbol{\theta}^{(\nu)} + \mathbf{I}_T(\mathbf{Y}, \boldsymbol{\theta}^{(\nu)})^{-1} \mathbf{S}(\boldsymbol{\theta}^{(\nu)}) \end{aligned}$$

(a) start with initial value  $\boldsymbol{\theta}^{(0)}$

(b) update from current  $\boldsymbol{\theta}^{(\nu)}$  to obtain  $\boldsymbol{\theta}^{(\nu+1)}$

(c) stop if  $\|\mathbf{S}(\boldsymbol{\theta}^{(\nu+1)})\|$  or  $\|\boldsymbol{\theta}^{(\nu+1)} - \boldsymbol{\theta}^{(\nu)}\|$  is sufficiently small. Otherwise keep updating

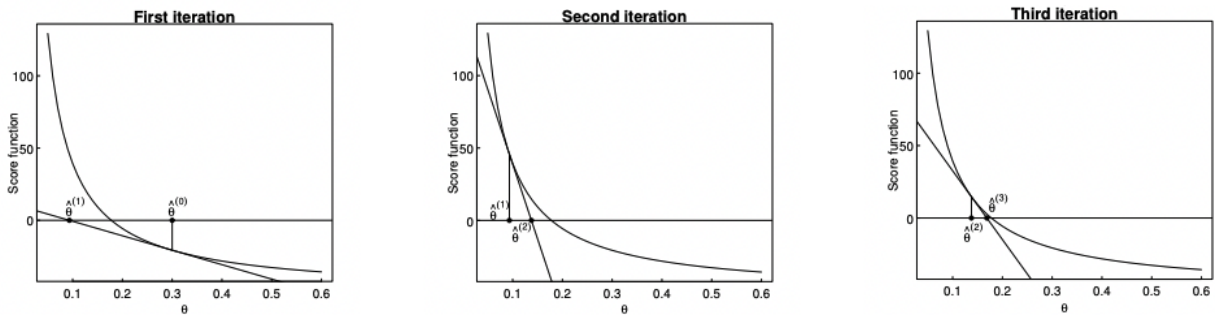


Figure 1: Newton-Raphson method (source)

If the first derivative is not well behaved in the neighborhood of a particular root, the method may overshoot. To prevent overshooting: At each step, ensure improvement of  $\boldsymbol{\theta}^{(\nu+1)}$  over  $\boldsymbol{\theta}^{(\nu)}$  by repeatedly halving the step size.

Under certain conditions, Newton methods have *local quadratic convergence*, which means that for some  $c > 0$

$$\|\boldsymbol{\theta}^{(\nu+1)} - \hat{\boldsymbol{\theta}}_{\text{MLE}}\| \leq c \|\boldsymbol{\theta}^{(\nu)} - \hat{\boldsymbol{\theta}}_{\text{MLE}}\|^2$$

*proof:* (although I am not totally agree with this proof)

$$\begin{aligned} 0 = S(\hat{\boldsymbol{\theta}}_{\text{MLE}}) &= S(\boldsymbol{\theta}^{(\nu)}) - I_T(\mathbf{Y}, \boldsymbol{\theta}^{(\nu)}) (\hat{\boldsymbol{\theta}}_{\text{MLE}} - \boldsymbol{\theta}^{(\nu)}) + \frac{1}{2} S''(\boldsymbol{\theta}^{(\nu)}) (\hat{\boldsymbol{\theta}}_{\text{MLE}} - \boldsymbol{\theta}^{(\nu)})^2 \\ \hat{\boldsymbol{\theta}}_{\text{MLE}} - \boldsymbol{\theta}^{(\nu+1)} &= \frac{1}{2} \frac{S''(\boldsymbol{\theta}^{(\nu)})}{I_T(\mathbf{Y}, \boldsymbol{\theta}^{(\nu)})} (\hat{\boldsymbol{\theta}}_{\text{MLE}} - \boldsymbol{\theta}^{(\nu)})^2 \end{aligned}$$

Thus, the local quadratic convergence holds under the following conditions:

- $I_T(\mathbf{Y}, \boldsymbol{\theta}) \neq 0$  in a neighborhood of  $\hat{\boldsymbol{\theta}}_{\text{MLE}}$
- $S''(\boldsymbol{\theta})$  is bounded
- $\boldsymbol{\theta}^{(\nu)}$  is sufficiently close to  $\hat{\boldsymbol{\theta}}_{\text{MLE}}$  (so that we can eliminate  $S(\boldsymbol{\theta}^{(\nu)})$ )

3. Fisher scoring:  $\mathbf{I}_T(\mathbf{Y}, \boldsymbol{\theta}^{(\nu)})$  in the above method replaced by its expectation  $\mathbf{I}_T(\boldsymbol{\theta}^{(\nu)})$

4. One-step estimator: typically  $\widehat{\boldsymbol{\theta}}_{\text{MLE}} - \boldsymbol{\theta} = O_p(n^{-1/2})$  (from asymptotic normality). If one starts with  $\boldsymbol{\theta}^{(0)}$  such that  $\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta} = O_p(n^{-1/2})$ , then

$$\boldsymbol{\theta}^{(1)} - \widehat{\boldsymbol{\theta}}_{\text{MLE}} = O_p(n^{-1})$$

under regularity conditions (by the local quadratic convergence property).

5. EM algorithm (**nice reading**): view observed data  $Y$  as incomplete, with  $Z$  missing. Write log joint likelihood of “complete” data as  $\ell_C(\boldsymbol{\theta} \mid \mathbf{Y}, \mathbf{Z})$

(a) E step: calculate

$$\begin{aligned} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(\nu)}, \mathbf{Y}) &= E_{\boldsymbol{\theta}^{(\nu)}} \{ \ell_C(\boldsymbol{\theta} \mid \mathbf{Y}, \mathbf{Z}) \mid \mathbf{Y} \} \\ &= \int \ell_C(\boldsymbol{\theta} \mid \mathbf{Y}, \mathbf{z}) f_{\mathbf{Z} \mid \mathbf{Y}}(\mathbf{z} \mid \mathbf{Y}, \boldsymbol{\theta}^{(\nu)}) d\mathbf{z} \end{aligned}$$

(b) maximize  $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(\nu)}, \mathbf{Y})$  with respect to  $\boldsymbol{\theta}$  to obtain  $\boldsymbol{\theta}^{\nu+1}$

*Example 1: 2-component mixtures*

$$Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} f(y; \boldsymbol{\theta}) = pf_1(y; \mu_1, \sigma_1) + (1-p)f_2(y; \mu_2, \sigma_2)$$

where  $f_1, f_2$  are normal densities

$$\begin{aligned} \boldsymbol{\theta} &= (\mu_1, \sigma_1, \mu_2, \sigma_2, p)^T \\ \ell(\boldsymbol{\theta}) &= \sum_{i=1}^n \log \{ pf_1(Y_i; \mu_1, \sigma_1) + (1-p)f_2(Y_i; \mu_2, \sigma_2) \} \end{aligned}$$

it's hard to maximize directly. Therefore, turn to the EM algorithm

$$\begin{aligned} Y_i &= Z_i X_{1i} + (1 - Z_i) X_{2i} \\ X_{11}, \dots, X_{1n} &\stackrel{\text{iid}}{\sim} N(\mu_1, \sigma_1^2) \\ X_{21}, \dots, X_{2n} &\stackrel{\text{iid}}{\sim} N(\mu_2, \sigma_2^2) \\ Z_1, \dots, Z_n &\stackrel{\text{iid}}{\sim} \text{Bernoulli}(p) \end{aligned}$$

join likelihood of complete data  $(Y, Z)$

$$\begin{aligned} L_C(\boldsymbol{\theta} \mid \mathbf{Y}, \mathbf{Z}) &= \prod_{i=1}^n \{ pf_1(Y_i; \mu_1, \sigma_1) \}^{Z_i} \times \{ (1-p)f_2(Y_i; \mu_2, \sigma_2) \}^{(1-Z_i)} \\ \ell_C(\boldsymbol{\theta} \mid \mathbf{Y}, \mathbf{Z}) &= \sum_{i=1}^n \{ Z_i \log f_1(Y_i; \mu_1, \sigma_1) + (1 - Z_i) \log f_2(Y_i; \mu_2, \sigma_2) + Z_i \log p + (1 - Z_i) \log(1 - p) \} \end{aligned}$$

cause

$$\begin{aligned} p(y, z = 1) &= p(y \mid z = 1)p(z = 1) = pf_1 \\ p(y, z = 0) &= p(y \mid z = 0)p(z = 0) = (1 - p)f_2 \end{aligned}$$

• E-step

$$\begin{aligned} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(\nu)}, \mathbf{Y}) &= E_{\boldsymbol{\theta}^{(\nu)}} \{ \ell_C(\boldsymbol{\theta} \mid \mathbf{Y}, \mathbf{Z}) \mid \mathbf{Y} \} \\ &= \sum_{i=1}^n \left\{ w_i^{(\nu)} \log f_1(Y_i; \mu_1, \sigma_1) + (1 - w_i^{(\nu)}) \log f_2(Y_i; \mu_2, \sigma_2) + w_i^{(\nu)} \log p + (1 - w_i^{(\nu)}) \log(1 - p) \right\} \end{aligned}$$

where

$$\begin{aligned} w_i^{(\nu)} &= E_{\boldsymbol{\theta}^{(\nu)}}(Z_i \mid Y_i) \\ &= \frac{p^{(\nu)} f_1(Y_i; \mu_1^{(\nu)}, \sigma_1^{(\nu)})}{p^{(\nu)} f_1(Y_i; \mu_1^{(\nu)}, \sigma_1^{(\nu)}) + (1 - p^{(\nu)}) f_2(Y_i; \mu_2^{(\nu)}, \sigma_2^{(\nu)})} \end{aligned}$$

substituting normal densities

$$\begin{aligned}
Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(\nu)}, \mathbf{Y}) &= \sum_{i=1}^n \left[ w_i^{(\nu)} \left\{ -\log \sigma_1 - \frac{(Y_i - \mu_1)^2}{2\sigma_1^2} \right\} \right. \\
&\quad \left. + (1 - w_i^{(\nu)}) \left\{ -\log \sigma_2 - \frac{(Y_i - \mu_2)^2}{2\sigma_2^2} \right\} \right. \\
&\quad \left. + w_i^{(\nu)} \log p + (1 - w_i^{(\nu)}) \log(1 - p) \right] \\
&\quad + \text{const.}
\end{aligned}$$

- M-step: maximizing  $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(\nu)}, \mathbf{Y})$  is simple. When the “complete” data likelihood has the form of an exponential family, M-step is straightforward.

*Example 2: right-censored data*

Let  $(Y_1, \delta_1), \dots, (Y_n, \delta_n)$  be i.i.d where  $Y_i = \min(X_i, R_i)$  and  $\delta_i = I(X_i \leq R_i)$ .

$$\begin{aligned}
X_i &\sim f(x; \boldsymbol{\theta}) = \sigma^{-1} \exp(-x/\sigma) \\
L(\boldsymbol{\theta} | \mathbf{Y}, \boldsymbol{\delta}) &\propto \prod_{i=1}^n f(Y_i; \boldsymbol{\theta})^{\delta_i} \{1 - F(Y_i; \boldsymbol{\theta})\}^{1-\delta_i} \\
&= \sigma^{-n_u} \exp\left(-\sum_{i=1}^n Y_i/\sigma\right)
\end{aligned}$$

cause  $F(Y_i; \boldsymbol{\theta}) = 1 - \exp(-x/\sigma)$ , where  $n_u$  is the number of uncensored observations. The likelihood is easy to maximize. Nonetheless, let  $Z_i = X_i$

$$\begin{aligned}
\ell(\boldsymbol{\theta} | \mathbf{Y}, \boldsymbol{\delta}, \mathbf{X}) &= \sum_{i=1}^n \log f(X_i; \boldsymbol{\theta}) \\
&= \sum_{i=1}^{n_u} \log f(X_i; \boldsymbol{\theta}) + \sum_{i=n_u+1}^n \log f(X_i; \boldsymbol{\theta})
\end{aligned}$$

- E-step

$$\begin{aligned}
Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(\nu)}, \mathbf{Y}, \boldsymbol{\delta}) &= \sum_{i=1}^{n_u} \log f(Y_i; \boldsymbol{\theta}) + \sum_{i=n_u+1}^n E_{\boldsymbol{\theta}^{(\nu)}} \{\log f(X_i; \boldsymbol{\theta}) | Y_i, \delta_i\} \\
&= -n \log \sigma - \sigma^{-1} \sum_{i=1}^{n_u} Y_i - \sigma^{-1} \sum_{i=n_u+1}^n E_{\sigma^{(\nu)}}(X_i | X_i > Y_i) \\
&= -n \log \sigma - \sigma^{-1} \sum_{i=1}^n Y_i - \sigma^{-1} (n - n_u) \sigma^{(\nu)}
\end{aligned}$$

things are simplified here by assuming  $R_i$  is not random

- M-step

$$\sigma^{(\nu+1)} = n^{-1} \left\{ \sum_{i=1}^n Y_i + (n - n_u) \sigma^{(\nu)} \right\}$$

as  $\nu \rightarrow \infty$ , we expect  $\sigma^{(\nu+1)} \rightarrow \hat{\sigma}_{\text{MLE}}$  and  $\sigma^{(\nu)} \rightarrow \hat{\sigma}_{\text{MLE}}$  thus

$$\hat{\sigma}_{\text{MLE}} = n^{-1} \left\{ \sum_{i=1}^n Y_i + (n - n_u) \hat{\sigma}_{\text{MLE}} \right\} \implies \hat{\sigma}_{\text{MLE}} = \sum_{i=1}^n Y_i / n_u$$

### 2.9.1 Why does EM work?

- Jensen’s inequality: for a random variable  $X$  and convex function  $\psi(\cdot)$ ,  $\psi\{E(X)\} \leq E\{\psi(X)\}$
- for  $\mathbf{Y} \sim f(\mathbf{y}; \boldsymbol{\theta}_0)$ ,  $E_{\boldsymbol{\theta}_0}\{\log f(\mathbf{Y}; \boldsymbol{\theta})\}$  is maximized at  $\boldsymbol{\theta} = \boldsymbol{\theta}_0$  *proof*:  $\psi(x) = -\log(x)$  is a convex function for  $x \in (0, \infty)$ . Therefore, by Jensen’s inequality

$$-\log \left[ E_{\boldsymbol{\theta}_0} \left\{ \frac{f(\mathbf{Y}; \boldsymbol{\theta})}{f(\mathbf{Y}; \boldsymbol{\theta}_0)} \right\} \right] \leq -E_{\boldsymbol{\theta}_0} \left[ \log \left\{ \frac{f(\mathbf{Y}; \boldsymbol{\theta})}{f(\mathbf{Y}; \boldsymbol{\theta}_0)} \right\} \right]$$

note that

$$E_{\theta_0} \left\{ \frac{f(Y; \theta)}{f(Y; \theta_0)} \right\} = \int \frac{f(Y; \theta)}{f(Y; \theta_0)} f(Y; \theta_0) dy = 1$$

therefore the left side is 0

$$\begin{aligned} -E_{\theta_0} \left[ \log \left\{ \frac{f(Y; \theta)}{f(Y; \theta_0)} \right\} \right] &\geq 0 \Rightarrow E_{\theta_0} \left[ \log \left\{ \frac{f(Y; \theta_0)}{f(Y; \theta)} \right\} \right] \leq 0 \\ \Rightarrow E_{\theta_0} [\log f(Y; \theta)] &\leq E_{\theta_0} [\log f(Y; \theta_0)] \Rightarrow E_{\theta_0} \{\log f(\mathbf{Y}; \boldsymbol{\theta})\} \text{ is maximized at } \boldsymbol{\theta} = \boldsymbol{\theta}_0 \end{aligned}$$

- now go back to the EM algorithm

$$\ell(\boldsymbol{\theta} | \mathbf{Y}) \equiv \log f(\mathbf{Y} | \boldsymbol{\theta}) = \log f(\mathbf{Y}, \mathbf{Z} | \boldsymbol{\theta}) - \log f(\mathbf{Z} | \mathbf{Y}, \boldsymbol{\theta})$$

cause  $f(y, z; \theta) = f(z | y; \theta) f(y; \theta)$ . Taking expectation on both sides treating  $Z$  as a random variable with density  $f(\mathbf{Z} | \mathbf{Y}, \boldsymbol{\theta}^{(\nu)})$

$$\ell(\boldsymbol{\theta} | \mathbf{Y}) = Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(\nu)}, \mathbf{Y}) - E_{\boldsymbol{\theta}^{(\nu)}} \{\log f(\mathbf{Z} | \mathbf{Y}, \boldsymbol{\theta}) | \mathbf{Y}\}$$

where based on the previous step  $E_{\boldsymbol{\theta}^{(\nu)}} \{\log f(\mathbf{Z} | \mathbf{Y}, \boldsymbol{\theta}) | \mathbf{Y}\}$  is maximized at  $\boldsymbol{\theta}^{(\nu)}$ . In performing the EM algorithm, we ensure that  $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(\nu)}, \mathbf{Y}) \geq Q(\boldsymbol{\theta}^{(\nu)}, \boldsymbol{\theta}^{(\nu)}, \mathbf{Y})$  by maximizing  $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(\nu)}, \mathbf{Y})$ . Therefore  $\ell(\boldsymbol{\theta}^{(\nu+1)} | \mathbf{Y}) \geq \ell(\boldsymbol{\theta}^{(\nu)} | \mathbf{Y})$ .

### 2.9.2 Calculating observed info matrix after EM

The general way is

$$\begin{aligned} \mathbf{S}_{\mathbf{Y}}(\mathbf{Y}, \boldsymbol{\theta}) &= \frac{\partial}{\partial \boldsymbol{\theta}^T} \log f_{\mathbf{Y}}(\mathbf{Y}; \boldsymbol{\theta}) = \frac{f'_{\mathbf{Y}}(\mathbf{Y}; \boldsymbol{\theta})^T}{f_{\mathbf{Y}}(\mathbf{Y}; \boldsymbol{\theta})} \\ \mathbf{I}_{\mathbf{Y}}(\mathbf{Y}, \boldsymbol{\theta}) &= -\frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} \log f_{\mathbf{Y}}(\mathbf{Y}; \boldsymbol{\theta}) \\ &= -\frac{f''_{\mathbf{Y}}(\mathbf{Y}; \boldsymbol{\theta})}{f_{\mathbf{Y}}(\mathbf{Y}; \boldsymbol{\theta})} + \mathbf{S}_{\mathbf{Y}}(\mathbf{Y}, \boldsymbol{\theta})^{\otimes 2} \end{aligned}$$

Actually, there is a better way than direct computation. Note for the complete data  $f_{\mathbf{Y}, \mathbf{Z}}(\mathbf{y}, \mathbf{z}; \boldsymbol{\theta})$

$$\begin{aligned} \mathbf{S}_{\mathbf{Y}, \mathbf{Z}}(\mathbf{Y}, \mathbf{Z}, \boldsymbol{\theta}) &= \frac{f'_{\mathbf{Y}, \mathbf{Z}}(\mathbf{Y}, \mathbf{Z}; \boldsymbol{\theta})^T}{f_{\mathbf{Y}, \mathbf{Z}}(\mathbf{Y}, \mathbf{Z}; \boldsymbol{\theta})} \\ \mathbf{I}_{\mathbf{Y}, \mathbf{Z}}(\mathbf{Y}, \mathbf{Z}, \boldsymbol{\theta}) &= -\frac{f''_{\mathbf{Y}, \mathbf{Z}}(\mathbf{Y}, \mathbf{Z}; \boldsymbol{\theta})}{f_{\mathbf{Y}, \mathbf{Z}}(\mathbf{Y}, \mathbf{Z}; \boldsymbol{\theta})} + \mathbf{S}_{\mathbf{Y}, \mathbf{Z}}(\mathbf{Y}, \mathbf{Z}, \boldsymbol{\theta})^{\otimes 2} \end{aligned} \tag{4}$$

given that the complete density is related to the density we are interested in, we have and based on 4 we have

$$\begin{aligned} f_{\mathbf{Y}}(\mathbf{y}; \boldsymbol{\theta}) &= \int f_{\mathbf{Y}, \mathbf{Z}}(\mathbf{y}, \mathbf{z}; \boldsymbol{\theta}) dz \\ \mathbf{S}_{\mathbf{Y}}(\mathbf{Y}, \boldsymbol{\theta}) &= \frac{\partial}{\partial \boldsymbol{\theta}^T} \log \int f_{\mathbf{Y}, \mathbf{Z}}(\mathbf{Y}, \mathbf{z}; \boldsymbol{\theta}) dz \\ &= \frac{\int f'_{\mathbf{Y}, \mathbf{Z}}(\mathbf{Y}, \mathbf{z}; \boldsymbol{\theta})^T dz}{\int f_{\mathbf{Y}, \mathbf{Z}}(\mathbf{Y}, \mathbf{z}; \boldsymbol{\theta}) dz} \\ &= \frac{\int \mathbf{S}_{\mathbf{Y}, \mathbf{Z}}(\mathbf{Y}, \mathbf{z}, \boldsymbol{\theta}) f_{\mathbf{Y}, \mathbf{Z}}(\mathbf{Y}, \mathbf{z}; \boldsymbol{\theta}) dz}{\int f_{\mathbf{Y}, \mathbf{Z}}(\mathbf{Y}, \mathbf{z}; \boldsymbol{\theta}) dz} \\ &= E_{\boldsymbol{\theta}} \{\mathbf{S}_{\mathbf{Y}, \mathbf{Z}}(\mathbf{Y}, \mathbf{Z}, \boldsymbol{\theta}) | \mathbf{Y}\} \\ \mathbf{I}_{\mathbf{Y}}(\mathbf{Y}, \boldsymbol{\theta}) &= -\frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} \log \int f_{\mathbf{Y}, \mathbf{Z}}(\mathbf{Y}, \mathbf{z}; \boldsymbol{\theta}) dz \\ &= -\frac{\int f''_{\mathbf{Y}, \mathbf{Z}}(\mathbf{Y}, \mathbf{z}; \boldsymbol{\theta}) dz}{\int f_{\mathbf{Y}, \mathbf{Z}}(\mathbf{Y}, \mathbf{z}; \boldsymbol{\theta}) dz} + \mathbf{S}_{\mathbf{Y}}(\mathbf{Y}, \boldsymbol{\theta})^{\otimes 2} \\ &= E_{\boldsymbol{\theta}} \{\mathbf{I}_{\mathbf{Y}, \mathbf{Z}}(\mathbf{Y}, \mathbf{Z}, \boldsymbol{\theta}) | \mathbf{Y}\} \\ &\quad - E_{\boldsymbol{\theta}} \{\mathbf{S}_{\mathbf{Y}, \mathbf{Z}}(\mathbf{Y}, \mathbf{Z}, \boldsymbol{\theta})^{\otimes 2} | \mathbf{Y}\} \\ &\quad + \mathbf{S}_{\mathbf{Y}}(\mathbf{Y}, \boldsymbol{\theta})^{\otimes 2} \end{aligned}$$

where  $\mathbf{I}_{\mathbf{Y}, \mathbf{Z}}$  and  $\mathbf{Z}_{\mathbf{Y}, \mathbf{Z}}$  only need to be computed at the last iteration of the EM procedure where  $\mathbf{S}_{\mathbf{Y}}(\mathbf{Y}, \hat{\boldsymbol{\theta}}) = \mathbf{0}$ .

## 2.10 Uniqueness of MLE

*Converge to boundary:* sequence  $\theta^{(1)}, \theta^{(2)}, \dots$  in  $\Theta$  is said to converge to boundary  $\partial\Theta$  if for every compact set (If a closed set  $A$  is bounded, then  $A$  is a compact set)  $K \subset \Theta$ , there exists  $k_0 \geq 1$  such that  $\theta^{(k)} \notin K \quad \forall k \geq k_0$  (if  $\Theta = R^b$ , then this definition is equivalent to  $\lim_{k \rightarrow \infty} \|\theta^{(k)}\| = \infty$ )

*Constant on boundary:* a real-valued function  $f$  defined on  $\Theta$  is said to be constant on boundary  $\partial\Theta$  if  $\lim_{k \rightarrow \infty} f(\theta^{(k)}) = c$  for every sequence  $\theta^{(k)}$  in  $\Theta$  converging to  $\partial\Theta$ ;  $c$  could be  $\pm\infty$ . Written as  $\lim_{\theta \rightarrow \partial\Theta} f(\theta) = c$ .

**Theorem 2.1.** for  $\Theta \subset R^b, b \geq 1$  which is a connected open set. If

1.  $\ell(\theta)$  is twice continuously differentiable with  $\lim_{\theta \rightarrow \partial\Theta} \ell(\theta) = c$ , where constant  $c$  is either a real number or  $-\infty$
2.  $\mathbf{I}_T(\mathbf{Y}, \theta) = -\partial^2 \ell(\theta) / \partial \theta \partial \theta^T$  (observed total information matrix) is positive definite at every point  $\theta \in \Theta$  for which  $\partial \ell(\theta) / \partial \theta^T = \mathbf{0}$

then

1. the critical point is unique and is the MLE
2.  $\ell(\theta) > c, \forall \theta \in \Theta$

Recall: Note a matrix  $A$  is positive definite if  $\mathbf{x}^T A \mathbf{x} > 0, \forall \mathbf{x}$ .  $A$  is positive definite if and only if

$$a_{11} > 0; \begin{vmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{vmatrix} > 0; \begin{vmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{vmatrix} > 0; \dots; |A| > 0$$

*Note:* the constancy on boundary is not really necessary for the scalar parameter but is need for other cases

- with a scalar parameter, the second conditional that minus Hessian is positive definite at all stationary points ( $x$  at which  $f'(x) = 0$ )  $\Rightarrow$  absence of local minima  $\Rightarrow$  absence of multiple local maxima (multiple local maxima cannot occur without local minima)
- counter-example showing the necessity of “constancy on boundary”

$$\begin{aligned} \ell &= g(x, y) = -e^{-2y} - e^{-y} \sin x \\ \frac{\partial}{\partial x} g(x, y) &= -e^{-y} \cos x \\ \frac{\partial}{\partial y} g(x, y) &= 2e^{-2y} + e^{-y} \sin x \end{aligned}$$

solution:  $x = (2k + 1.5)\pi, k = 1, 2, \dots$ , and  $y = \log 2$

$$\begin{aligned} -\nabla^2 g(x, y) &= \begin{pmatrix} -e^{-y} \sin x & -e^{-y} \cos x \\ -e^{-y} \cos x & 4e^{-2y} + e^{-y} \sin x \end{pmatrix} \\ &\stackrel{\text{at stationary points}}{=} \begin{pmatrix} 0.5 & 0 \\ 0 & 0.5 \end{pmatrix} \end{aligned}$$

which is positive definite  $\Rightarrow$  all solutions corresponding to maxima, which is the solutions are all maxima and there are countably infinite of them cause the “constancy on boundary” is violated

**Theorem 2.2.** for  $\Theta \subset R^b, b \geq 1$  which is a connected open set and  $\ell(\theta)$  is twice continuously differentiable. if

1.  $\partial \ell(\theta) / \partial \theta^T = \mathbf{0}$  has at least one solution
2.  $\mathbf{I}_T(\mathbf{Y}, \theta)$  is positive definite  $\forall \theta \in \Theta$

then

1.  $\ell(\theta)$  is concave
2. the solution is unique and is the MLE

*Example 1:* normal location-scale model

$$Y_1, \dots, Y_n \sim \text{iid } N(\mu, \sigma^2), \Theta = (-\infty, \infty) \times (0, \infty)$$



$$\begin{aligned}\ell(\mu, \sigma) &= \text{const} - n \log \sigma - \frac{1}{2\sigma^2} \sum_{i=1}^n (Y_i - \mu)^2 \\ \frac{\partial}{\partial \boldsymbol{\theta}^T} \ell(\mu, \sigma) &= \begin{pmatrix} \frac{1}{\sigma^2} \sum_{i=1}^n (Y_i - \mu) \\ -\frac{n}{\sigma} + \frac{1}{\sigma^3} \sum_{i=1}^n (Y_i - \mu)^2 \end{pmatrix} \\ \mathbf{I}_T(\mu, \sigma) &= \begin{pmatrix} \frac{n}{\sigma^2} & \frac{2}{\sigma^3} \sum_{i=1}^n (Y_i - \mu) \\ \frac{2}{\sigma^3} \sum_{i=1}^n (Y_i - \mu) & -\frac{n}{\sigma^2} + \frac{3}{\sigma^4} \sum_{i=1}^n (Y_i - \mu)^2 \end{pmatrix}\end{aligned}$$

only solution to  $\partial \ell(\mu, \sigma) / \partial \boldsymbol{\theta}^T$  is  $\hat{\mu} = \bar{Y}$ ,  $\hat{\sigma} = s_n$ . Use 2.1 we can show that 1)  $\mathbf{I}_T(\bar{Y}, s_n) = \text{diag}(n/s_n^2, 2n/s_n^2)$  is positive definite, and 2)  $\lim_{\boldsymbol{\theta} \rightarrow 2\Theta} \ell(\boldsymbol{\theta}) = -\infty$  (HW3 2.58)

2.2 is not applicable cause  $\mathbf{I}_T(\mu, \sigma)$  is not positive definite everywhere.  $|\mathbf{I}_T(\mu, \sigma)| = n^2 \sigma^{-6} \{3s_n^2 - \sigma^2 - (\bar{Y} - \mu)^2\}$  could be negative.

*Example 2:* Exponential threshold model (combine with the profile likelihood approach)

$$f(y; \mu, \sigma) = \frac{1}{\sigma} \exp\left(-\frac{y - \mu}{\sigma}\right), \quad \mu \leq y < \infty$$

with  $Y_i$ 's are i.i.d and  $\Theta = (-\infty, \infty) \times (0, \infty)$

$$L(\mu, \sigma | \mathbf{Y}) = \sigma^{-n} \exp\left(-\sum_{i=1}^n \frac{Y_i - \mu}{\sigma}\right) \prod_{i=1}^n I(\mu \leq Y_i)$$

which is not differentiable in  $\mu$  everywhere. By using the profile likelihood method. For each  $\sigma$ ,  $L(\mu, \sigma | \mathbf{Y})$  is maximized at  $\hat{\mu} = Y_{(1)} \Rightarrow$  log-profile likelihood  $\ell^*(\sigma) = -n \log \sigma - \sigma^{-1} \sum_i (Y_i - Y_{(1)})$

$$\begin{aligned}\frac{\partial}{\partial \sigma} \ell^*(\sigma) &= -n\sigma^{-1} + n\sigma^{-2} (\bar{Y} - Y_{(1)}) \implies \hat{\sigma} = \bar{Y} - Y_{(1)} \text{ if } \bar{Y} > Y_{(1)} \\ -\frac{\partial^2}{\partial \sigma^2} \ell^*(\sigma) &= -n\sigma^{-2} + 2n\sigma^{-3} (\bar{Y} - Y_{(1)})\end{aligned}$$

is not always positive thus 2.2 does not apply. Nonetheless, 2.1 does

1.  $-\partial^2 \ell^*(\sigma) / \partial \sigma^2 \big|_{\sigma=\hat{\sigma}} = n\hat{\sigma}^{-2} > 0$
2.  $\lim_{\sigma \rightarrow 0} \ell^*(\sigma) = \lim_{\sigma \rightarrow \infty} \ell^*(\sigma) = -\infty$

*Example 3:* example to show that existence and uniqueness of the MLE is not always a given

$Y_1, \dots, Y_n \stackrel{i.i.d}{\sim}$  mixture of normals

$$f(y; \mu, \sigma, p) = p\phi(y - \mu) + (1 - p) \frac{1}{\sigma} \phi\left(\frac{y - \mu}{\sigma}\right)$$

the log-likelihood assuming  $\mu = Y_1$

$$\ell = \log \left\{ p\phi(0) + (1 - p) \frac{1}{\sigma} \phi(0) \right\} + \sum_{i=2}^n \log \left\{ p\phi(Y_i - Y_1) + (1 - p) \frac{1}{\sigma} \phi\left(\frac{Y_i - Y_1}{\sigma}\right) \right\}$$

as  $\sigma \rightarrow 0$  the first term  $\rightarrow \infty$  and the rest are bounded. Thus, MLE does not exist in the strict sense. Furthermore, there exist multiple local maxima.

Nonetheless, local maxima satisfying the likelihood score tend to behave well. EM algorithm can find at least one of them.

### Uniqueness of the MLE for exponential family

$X \sim$  minimal exponential family

$$f(x; \boldsymbol{\theta}) = h(x) \exp \left\{ \sum_{i=1}^s g_i(\boldsymbol{\theta}) T_i(x) - B(\boldsymbol{\theta}) \right\} = h(x) \exp \left\{ \sum_{i=1}^k \eta_i T_i(\boldsymbol{\eta}) - A(\boldsymbol{\eta}) \right\}$$

where  $g(\boldsymbol{\theta})$  is 1-to-1, twice differentiable in  $\Theta$  and  $\Theta$  is an open subset of  $R^s$ . Then if there is at least one solution to the transformed likelihood equation  $E_{\boldsymbol{\theta}}\{\mathbf{T}(X)\} = \mathbf{T}(x)$  then with 2.2 the solution is unique and is the MLE.

With the canonical representation, the likelihood score is

$$\frac{\partial}{\partial \boldsymbol{\eta}^T} A(\boldsymbol{\eta}) - \mathbf{T}(x) = 0 \implies \frac{\partial}{\partial \boldsymbol{\eta}^T} A(\boldsymbol{\eta}) \equiv E_{\boldsymbol{\eta}} \mathbf{T}(X) = \mathbf{T}(x)$$

the minus Hessian is

$$\frac{\partial^2}{\partial \boldsymbol{\eta} \partial \boldsymbol{\eta}^T} A(\boldsymbol{\eta}) \equiv \mathbf{I}(\boldsymbol{\eta}) = \text{Var}_{\boldsymbol{\eta}}\{\mathbf{T}(X)\}$$

which is positive definite as  $T$  is affinely independent. Then based on invariant property so is  $\hat{\boldsymbol{\theta}}_{\text{MLE}} = g^{-1}(\hat{\boldsymbol{\eta}}_{\text{MLE}})$

**Families with truncation or threshold parameters**

$$x_1, \dots, x_n \stackrel{i.i.d.}{\sim} f(x; \boldsymbol{\theta}, \mu_1, \mu_2)$$

$$f(x; \boldsymbol{\theta}, \mu_1, \mu_2) = c(\mu_1, \mu_2, \boldsymbol{\theta}) d(x, \boldsymbol{\theta}) \quad \mu_1 < x < \mu_2$$

for fixed  $\boldsymbol{\theta}$ ,  $(X_{(1)}, X_{(n)})$  is minimal sufficient for  $(\mu_1, \mu_2)$ . Conditional on  $(X_{(1)}, X_{(n)}) = (x_{(1)}, x_{(n)})$ , sample values between  $X_{(1)}$  and  $X_{(n)}$  say

$$z_1, \dots, z_{n-2} \stackrel{i.i.d.}{\sim} q(z, \boldsymbol{\theta}) = \frac{d(z, \boldsymbol{\theta}) I(x_{(1)} < z < x_{(n)})}{\int_{x_{(1)}}^{x_{(n)}} d(z, \boldsymbol{\theta}) dz}$$

if  $d(x, \boldsymbol{\theta})$  has an exponential family form then so does  $q(z, \boldsymbol{\theta})$

Similar results for

$$f(x; \boldsymbol{\theta}, \mu) = c_1(\mu, \boldsymbol{\theta}) d_1(x, \boldsymbol{\theta}) \quad \mu < x$$

$$f(x; \boldsymbol{\theta}, \mu) = c_2(\mu, \boldsymbol{\theta}) d_2(x, \boldsymbol{\theta}) \quad x < \mu$$

*Example 1:*

$$f(x; \theta, \mu) = \theta e^{-\theta(x-\mu)} I(\mu < x)$$

conditional on  $X_{(1)} = x_{(1)}$ , observations larger than  $x_{(1)}$  have density

$$q(z, \theta) = \theta e^{-\theta(z-x_{(1)})} I(x_{(1)} < z) = e^{-\theta z + \theta x_{(1)} + \log(\theta)} I(x_{(1)} < z)$$

### 3 Likelihood-based tests and confidence regions

Let's first look at the simplest scalar parameter case

$$H_0 : \theta = \theta_0 \text{ vs. } H_a : \theta \neq \theta_0$$

- Wald test

$$T_W = \frac{(\hat{\theta}_{\text{MLE}} - \theta_0)^2}{\left\{ I_T(\hat{\theta}_{\text{MLE}}) \right\}^{-1}}$$

- Likelihood ratio test

$$T_{\text{LR}} = -2 \left\{ \ell(\theta_0) - \ell(\hat{\theta}_{\text{MLE}}) \right\}$$

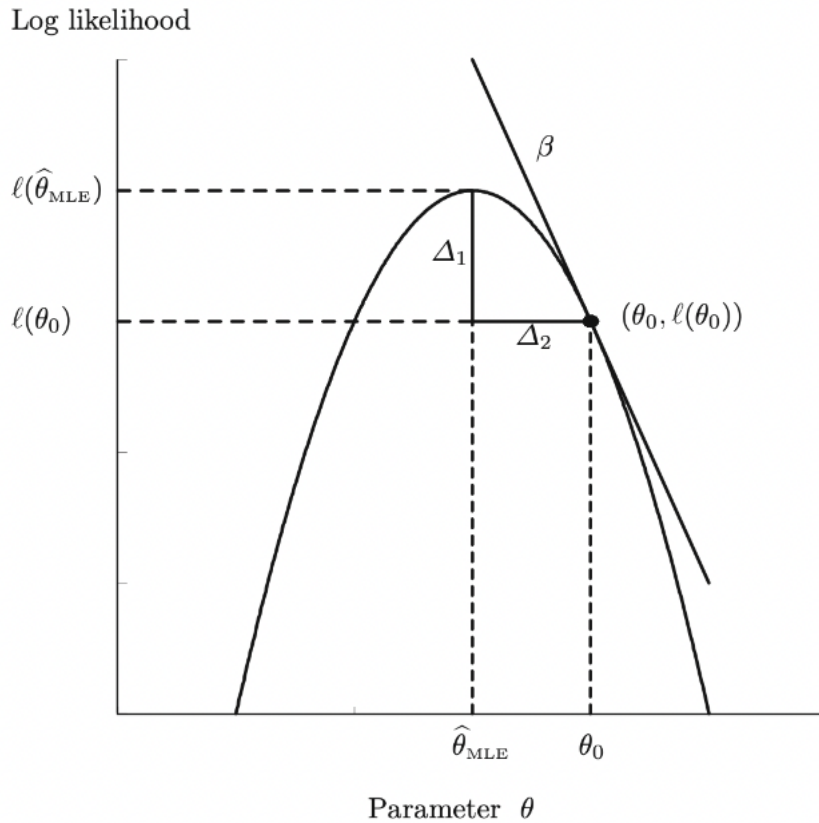
- Score test

$$T_S = \frac{S(\theta_0)^2}{I_T(\theta_0)}$$

Under  $H_0$ , asymptotically they are all  $\chi_1^2$ . Under local alternatives, they have identical asymptotic non-central  $\chi^2$  distribution.

The idea behind those tests are graphically illustrated below

Figure 2: Graphical representation of the relationships between Wald, Score and Likelihood Ratio test statistics



The Likelihood Ratio test statistic is a multiple of the difference,  $\Delta_1$ ; the Wald test statistic is a multiple of the squared difference,  $\Delta_2^2$ ; and the Score test statistic is a multiple of the squared slope  $\beta^2$  (this figure is from [the textbook](#))

### 3.1 Simple null hypothesis

By "simple" we mean don't have nuisance parameter which is we are interested in the whole parameter vector.

$$H_0 : \theta = \theta_0 \text{ vs. } H_a : \theta \neq \theta_0 \quad b\text{-dimensional}$$

- Wald test

$$T_W = \left( \hat{\theta}_{MLE} - \theta_0 \right)^T \mathbf{I}_T \left( \hat{\theta}_{MLE} \right) \left( \hat{\theta}_{MLE} - \theta_0 \right)$$

where  $\mathbf{I}_T \left( \hat{\theta}_{MLE} \right)$  may be replace by  $\mathbf{I}_T \left( \theta_0 \right)$  or  $\mathbf{I}_T \left( \mathbf{Y}, \hat{\theta}_{MLE} \right)$

- Likelihood ratio test

$$T_{LR} = -2 \left\{ \ell \left( \theta_0 \right) - \ell \left( \hat{\theta}_{MLE} \right) \right\}$$

By taylor expansion

$$\ell \left( \theta_0 \right) = \ell \left( \hat{\theta}_{MLE} \right) + \mathbf{S} \left( \hat{\theta}_{MLE} \right)^T \left( \theta_0 - \hat{\theta}_{MLE} \right) - \frac{1}{2} \left( \theta_0 - \hat{\theta}_{MLE} \right)^T \mathbf{I}_T \left( \mathbf{Y}, \hat{\theta}_{MLE} \right) \left( \theta_0 - \hat{\theta}_{MLE} \right) + \text{residual}$$

here  $\mathbf{S} \left( \hat{\theta}_{MLE} \right)^T = 0$  cause this is how MLE is derived. Therefore

$$T_W + \delta_n \text{ with } \delta_n \xrightarrow{p} 0$$

where  $\delta_n$  is the residual

- Score test

$$T_S = \mathbf{S} \left( \theta_0 \right)^T \mathbf{I}_T \left( \theta_0 \right)^{-1} \mathbf{S} \left( \theta_0 \right)$$

Under  $H_0$ ,  $E \{ \mathbf{S} \left( \theta_0 \right) \} = 0$  and  $\text{var} \{ \mathbf{S} \left( \theta_0 \right) \} = \mathbf{I}_T \left( \theta_0 \right)$ . By central limit theorem and continuous mapping theorem  $T_S \xrightarrow{d} \chi_b^2$ .

Example 1:

$$\begin{aligned}
Y_1, \dots, Y_n &\stackrel{\text{i.i.d.}}{\sim} N(\mu, 1), H_0 : \mu = \mu_0 \\
\ell(\mu) &= -\frac{n}{2} \log(2\pi) - \frac{1}{2} \sum_i (Y_i - \mu)^2 \\
S(\mu) &= \frac{\partial}{\partial \mu} \ell(\mu) = \sum_i (Y_i - \mu) \implies \hat{\mu} = \bar{Y} \\
I_T(\mathbf{Y}, \mu) &= -\frac{\partial}{\partial \mu} S(\mu) = n
\end{aligned}$$

we got the test statistics as follows

$$\begin{aligned}
T_W &= (\bar{Y} - \mu_0) (n) (\bar{Y} - \mu_0) = n (\bar{Y} - \mu_0)^2 \\
T_S &= \left\{ \sum_i (Y_i - \mu_0) \right\} n^{-1} \left\{ \sum_i (Y_i - \mu_0) \right\} = T_W \\
T_{LR} &= -2 \left\{ -\frac{1}{2} \sum_i (Y_i - \mu_0)^2 + \frac{1}{2} \sum_i (Y_i - \bar{Y})^2 \right\} = T_W
\end{aligned}$$

Example 2:  $Y_1, \dots, Y_n \stackrel{\text{i.i.d.}}{\sim} \text{Bernoulli}(p)$ .  $H_0 : p = p_0$ . Let  $X = \sum_i Y_i$

$$\begin{aligned}
L &= \prod_{i=1}^n p^{y_i} (1-p)^{(1-y_i)} \\
\ell &= \sum_{i=1}^n y_i \log(p) + (1-y_i) \log(1-p) = X \log(p) + (n-X) \log(1-p) \\
S(p) &= \frac{X}{p} - \frac{n-X}{1-p} = \frac{X-np}{p(1-p)} \implies \hat{p} = \frac{X}{n} \implies X = n\hat{p} \\
I_T(\mathbf{Y}, p) &= \frac{X}{p^2} + \frac{n-X}{(1-p)^2} \quad I_T(p) = \frac{n}{p(1-p)} \\
T_W &= (\hat{p} - p_0) \frac{n}{\hat{p}(1-\hat{p})} (\hat{p} - p_0) = \frac{n(\hat{p} - p_0)^2}{\hat{p}(1-\hat{p})} \\
T_S &= S(p_0) I_T(p_0)^{-1} S(p_0) = \frac{n(\hat{p} - p_0)^2}{p_0(1-p_0)} \quad X \text{ is replaced by } n\hat{p} \\
T_{LR} &= -2 [X \log(p_0/\hat{p}) + (n-X) \log\{(1-p_0)/(1-\hat{p})\}]
\end{aligned}$$

### 3.2 Composite null hypothesis

In this case,  $\theta$  is partitioned and we are only interested in part of it.

$$\boldsymbol{\theta}_{b \times 1} = \begin{pmatrix} \boldsymbol{\theta}_1 \\ r \times 1 \\ \boldsymbol{\theta}_2 \\ (b-r) \times 1 \end{pmatrix}$$

$H_0 : \boldsymbol{\theta}_1 = \boldsymbol{\theta}_{10}$  vs.  $H_a : \boldsymbol{\theta}_1 \neq \boldsymbol{\theta}_{10}$ , with  $\boldsymbol{\theta}_2$  as nuisance.

- Wald test

$$\mathbf{I}_T(\hat{\boldsymbol{\theta}}_{\text{MLE}}) \equiv \hat{\mathbf{I}}_T = \begin{pmatrix} \hat{\mathbf{I}}_{T,11} & \hat{\mathbf{I}}_{T,12} \\ \hat{\mathbf{I}}_{T,21} & \hat{\mathbf{I}}_{T,22} \end{pmatrix}$$

Under  $H_0$ ,  $\text{Avar}(\hat{\boldsymbol{\theta}}_1 - \boldsymbol{\theta}_{10})$  is the upper (1,1) element of  $\hat{\mathbf{I}}_T^{-1}$ , given by

$$\left( \hat{\mathbf{I}}_{T,11} - \hat{\mathbf{I}}_{T,12} \hat{\mathbf{I}}_{T,22}^{-1} \hat{\mathbf{I}}_{T,21} \right)^{-1}$$

therefore

$$T_W = (\hat{\boldsymbol{\theta}}_1 - \boldsymbol{\theta}_{10})^T \left( \hat{\mathbf{I}}_{T,11} - \hat{\mathbf{I}}_{T,12} \hat{\mathbf{I}}_{T,22}^{-1} \hat{\mathbf{I}}_{T,21} \right) (\hat{\boldsymbol{\theta}}_1 - \boldsymbol{\theta}_{10}) \stackrel{d}{\rightarrow} \chi_r^2$$

be careful here, don't write the middle part as  $(\hat{\mathbf{I}}_{T,11})^{-1}$

- Likelihood ratio test

$$T_{LR} = -2 \log \left\{ \frac{\sup_{\theta \in H_0} L(\theta | \mathbf{Y})}{\sup_{\theta \in \Theta} L(\theta | \mathbf{Y})} \right\} = -2 \{ \ell(\tilde{\theta}) - \ell(\hat{\theta}) \}$$

where  $\sup_{\theta \in H_0} L(\theta | \mathbf{Y})$  is the restricted MLE and  $\sup_{\theta \in \Theta} L(\theta | \mathbf{Y})$  is the unrestricted MLE

- Score test

First derive the  $H_0$ -restricted MLE

$$\tilde{\theta} = \begin{pmatrix} \theta_{10} \\ \tilde{\theta}_2 \end{pmatrix} \quad \tilde{\theta}_2 = \arg \max_{\theta_2} \ell(\theta_{10}, \theta_2)$$

$$\mathbf{I}_T(\tilde{\theta}_{MLE}) \equiv \tilde{\mathbf{I}}_T$$

$$\mathbf{S}(\theta) = \begin{pmatrix} \mathbf{S}_1(\theta) \\ \mathbf{S}_2(\theta) \end{pmatrix} = \begin{pmatrix} \frac{\partial}{\partial \theta_1^T} \ell(\theta) \\ \frac{\partial}{\partial \theta_2^T} \ell(\theta) \end{pmatrix}$$

$$T_S = \mathbf{S}(\tilde{\theta})^T \tilde{\mathbf{I}}_T^{-1} \mathbf{S}(\tilde{\theta})$$

$$= \begin{pmatrix} \mathbf{S}_1(\tilde{\theta})^T, \mathbf{0}^T \end{pmatrix} \begin{pmatrix} \tilde{\mathbf{I}}_{T,11} & \tilde{\mathbf{I}}_{T,12} \\ \tilde{\mathbf{I}}_{T,21} & \tilde{\mathbf{I}}_{T,22} \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{S}_1(\tilde{\theta}) \\ \mathbf{0} \end{pmatrix}$$

$$= \mathbf{S}_1(\tilde{\theta})^T \left( \tilde{\mathbf{I}}_{T,11} - \tilde{\mathbf{I}}_{T,12} \tilde{\mathbf{I}}_{T,22}^{-1} \tilde{\mathbf{I}}_{T,21} \right)^{-1} \mathbf{S}_1(\tilde{\theta})$$

where  $\mathbf{S}_2(\tilde{\theta}) = \mathbf{0}$  by definition of MLE

**Composite null hypotheses of general form**  $H_0 : \mathbf{h}(\theta) = \mathbf{0}$   
 $r \times 1$

$$\mathbf{H}(\theta) = \frac{\partial}{\partial \theta} \mathbf{h}(\theta), \quad r \leq b$$

which needs to be of *full rank*

Below are some examples to help understand the definition of  $h$  and  $H$  *Example 1: bi-variate normal data*

$$\theta = (\mu_1, \mu_2, \sigma_1, \sigma_2, \rho)^T$$

$$H_0 : \mu_1 = \mu_2$$

$$h(\theta) = \mu_1 - \mu_2$$

$$\mathbf{H}(\theta) = (1, -1, 0, 0, 0)$$

*Example 2: linear hypotheses*  $\mathbf{K}^T \beta = \mathbf{m}$

$$\mathbf{h}(\beta) = \mathbf{K}^T \beta - \mathbf{m}$$

*Example 3: partitioned-vector hypothesis*  $H_0 : \theta_1 = \theta_{10}$

$$\mathbf{h}(\theta) = \theta_1 - \theta_{10}$$

Now let's jump in to the test statistics

- Wald test

$$\hat{\theta} \sim AN(\theta, \mathbf{I}_T(\theta)^{-1})$$

$$\mathbf{h}(\hat{\theta}) \sim AN\{\mathbf{h}(\theta), \mathbf{H}(\theta) \mathbf{I}_T(\theta)^{-1} \mathbf{H}(\theta)^T\}$$

$$T_W = \mathbf{h}(\hat{\theta})^T \left\{ \mathbf{H}(\hat{\theta}) \mathbf{I}_T(\hat{\theta})^{-1} \mathbf{H}(\hat{\theta})^T \right\}^{-1} \mathbf{h}(\hat{\theta})$$

from the first line to second line we used Delta method one problem is that the test statistic varies with reparameterization and choice of  $h$ , e.g.,  $h(\mu_1, \mu_2) = \mu_1 - \mu_2$  and  $h(\mu_1, \mu_2) = \mu_1 / \mu_2 - 1$ .

- Score test

$$T_S = \mathbf{S}(\tilde{\theta})^T \tilde{\mathbf{I}}_T^{-1} \mathbf{S}(\tilde{\theta})$$

where  $\tilde{\theta}$  maximizes the likelihood subject to  $\mathbf{h}(\theta) = \mathbf{0}$ .

Actually a better way to do this is using the Lagrange multiplier for

$$\max \ell(\theta) \text{ subject to } \mathbf{h}(\theta) = \mathbf{0}$$

we will focus on  $\ell(\boldsymbol{\theta}) - \mathbf{h}(\boldsymbol{\theta})_{r \times 1}^T \boldsymbol{\lambda}$

$$\mathbf{S}(\boldsymbol{\theta}) - \mathbf{H}(\boldsymbol{\theta})^T \boldsymbol{\lambda} = \mathbf{0}$$

$$\mathbf{h}(\boldsymbol{\theta}) = \mathbf{0}$$

Denote the solution by  $\tilde{\boldsymbol{\theta}}$  and  $\tilde{\boldsymbol{\lambda}}$  then we get  $\mathbf{S}(\tilde{\boldsymbol{\theta}}) = \mathbf{H}(\tilde{\boldsymbol{\theta}})^T \tilde{\boldsymbol{\lambda}}$  and the score test statistic turns into

$$T_S = \tilde{\boldsymbol{\lambda}}^T \mathbf{H}(\tilde{\boldsymbol{\theta}}) \tilde{\mathbf{I}}_{\mathbf{T}}^{-1} \mathbf{H}(\tilde{\boldsymbol{\theta}})^T \tilde{\boldsymbol{\lambda}}$$

Below is a comparison of the three method

	invariant to re-parameterization and the choice of $\mathbf{h}(\cdot)$		computing MLE's	
	unrestricted		restricted	
Wald	N		Y	N
Score	Y		N	Y
LR	Y		Y	Y

Nonetheless, LR can be convenient for nested models. Typically, Wald does not have as good type I error as score and LR. And Wald and score are easier to adjust in the case of model misspecification.

*Example 1:* Normal location-scale model  $Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2)$  and  $H_0 : \mu = \mu_0$  vs.  $H_a : \mu \neq \mu_0$   $\sigma$  unrestricted

$$\ell(\mu, \sigma) = \text{constant} - n \log \sigma - \frac{1}{2\sigma^2} \sum_i (Y_i - \mu)^2$$

$$\mathbf{S}(\mu, \sigma) = \sum_i \left[ \begin{array}{c} \sigma^{-2} (Y_i - \mu) \\ \sigma^{-3} \left\{ (Y_i - \mu)^2 - \sigma^2 \right\} \end{array} \right] \Rightarrow \hat{\mu} = \bar{Y} \quad \hat{\sigma}^2 = s_n^2 = n^{-1} \sum_i (Y_i - \bar{Y})^2, \quad \mathbf{I}_{\mathbf{T}}(\mu, \sigma) = \text{diag}(n/\sigma^2, 2n/\sigma^2)$$

$$\tilde{\mu} = \mu_0 \quad \tilde{\sigma}^2 = n^{-1} \sum_i (Y_i - \mu_0)^2 = s_n^2 + (\bar{Y} - \mu_0)^2$$

$$T_W = (\bar{Y} - \mu_0) \frac{n}{s_n^2} (\bar{Y} - \mu_0) = \frac{n(\bar{Y} - \mu_0)^2}{s_n^2} = \frac{n}{n-1} \frac{n(\bar{Y} - \mu_0)^2}{ns_n^2/(n-1)} = \frac{n}{n-1} t^2 \text{ where } t \text{ means } t\text{-distribution with } df = 1$$

$$T_S = \left\{ \frac{1}{\tilde{\sigma}^2} \sum_i (Y_i - \mu_0) \right\} \frac{\tilde{\sigma}^2}{n} \left\{ \frac{1}{\tilde{\sigma}^2} \sum_i (Y_i - \mu_0) \right\} = \frac{n(\bar{Y} - \mu_0)^2}{\tilde{\sigma}^2} = \frac{nT_W}{n + T_W}$$

$$T_{LR} = n \log \frac{\tilde{\sigma}^2}{\hat{\sigma}^2} = n \log \left\{ 1 + \frac{(\bar{Y} - \mu_0)^2}{s_n^2} \right\} = n \log \left\{ 1 + \frac{T_W}{n} \right\}$$

Since  $\frac{x}{1+x} \leq \log(1+x) \leq x$  for  $x > -1$

$$T_S \leq T_{LR} \leq T_W$$

Using exact distributions, they are all equivalent to  $t$  test (can be transformed to  $t$  test). Using the asymptotic  $\chi_1^2$  critical values,  $T_W$  is more liberal.

Example 2: Score test for multinomial data = Pearson  $\chi^2$  test

$$\begin{aligned}
(N_1, \dots, N_k) &\sim \text{Multinomial}(n, p_1, \dots, p_k) \\
\mathbf{p} &= (p_1, \dots, p_{k-1})^T \quad p_k = 1 - p_1 - \dots - p_{k-1} \\
f(n_1, \dots, n_k) &= \frac{n!}{n_1! \dots n_k!} p_1^{n_1} \dots p_k^{n_k} \\
\ell(\mathbf{p}) &= \sum_{i=1}^{k-1} N_i \log p_i + N_k \log \left( 1 - \sum_{i=1}^{k-1} p_i \right) + \text{const.} \\
\mathbf{S}(\mathbf{p}) &= \left( \frac{N_1}{p_1} - \frac{N_k}{p_k}, \dots, \frac{N_{k-1}}{p_{k-1}} - \frac{N_k}{p_k} \right)^T \\
\mathbf{I}_T(\mathbf{p}) &= n \text{diag} \left( \frac{1}{p_1}, \dots, \frac{1}{p_{k-1}} \right) + \frac{n}{p_k} \mathbf{1}^{\otimes 2} \\
\mathbf{I}_T(\mathbf{p})^{-1} &= \{ \text{diag}(\mathbf{p}) - \mathbf{p}^{\otimes 2} \} / n \\
T_S &= \mathbf{S}(\tilde{\mathbf{p}})^T \mathbf{I}_T(\tilde{\mathbf{p}})^{-1} \mathbf{S}(\tilde{\mathbf{p}}) \\
&= n^{-1} \sum_{i=1}^k \left( \frac{N_i}{\tilde{p}_i} - \frac{N_k}{\tilde{p}_k} \right)^2 \tilde{p}_i \\
&\quad - n^{-1} \left\{ \sum_{i=1}^k \left( \frac{N_i}{\tilde{p}_i} - \frac{N_k}{\tilde{p}_k} \right) \tilde{p}_i \right\}^2 \\
&= n^{-1} \sum_{i=1}^k \left\{ \frac{N_i}{\tilde{p}_i} - \frac{N_k}{\tilde{p}_k} - \left( n - \frac{N_k}{\tilde{p}_k} \right) \right\}^2 \tilde{p}_i \\
&= \sum_{i=1}^k \frac{(N_i - n\tilde{p}_i)^2}{n\tilde{p}_i}
\end{aligned}$$

in the last step, the variance equality

$$\sum a_i^2 p_i - \left( \sum a_i p_i \right)^2 = \sum \left( a_i - \sum a_i p_i \right)^2 p_i \text{ is used.}$$

Example 3: Testing for Hardy-Weinberg equilibrium

Multinomial  $k = 3 : p_{AA}, p_{Aa}, p_{aa}$

Under the equilibrium,  $p_{AA} = p_A^2, p_{Aa} = 2p_A(1 - p_A),$

$p_{aa} = (1 - p_A)^2$

restricted MLE:  $\tilde{p}_A = (2N_{AA} + N_{Aa}) / (2n)$

$$\begin{aligned}
T_S &= \frac{(N_{AA} - n\tilde{p}_A^2)^2}{n\tilde{p}_A^2} + \frac{\{N_{Aa} - 2n\tilde{p}_A(1 - \tilde{p}_A)\}^2}{2n\tilde{p}_A(1 - \tilde{p}_A)} \\
&\quad + \frac{\{N_{aa} - n(1 - \tilde{p}_A)^2\}^2}{n(1 - \tilde{p}_A)^2}
\end{aligned}$$

the degree of freedom is 1 cause

$$h(P_{AA}, P_{Aa}) = P_{Aa} - 2P_{AA}^{1/2} (1 - P_{AA}^{1/2}) \Rightarrow r = 1$$

### 3.3 Confidence interval

$$\boldsymbol{\theta}_{b \times 1} = \begin{pmatrix} \boldsymbol{\theta}_1 \\ r \times 1 \\ \boldsymbol{\theta}_2 \\ (b - r) \times 1 \end{pmatrix}$$

$T(\boldsymbol{\theta}_1) =$  test stat with  $\boldsymbol{\theta}_1$  as the null value (actually the true value)

$$C_{1-\alpha} = \{ \boldsymbol{\theta}_1 : T(\boldsymbol{\theta}_1) \leq \chi_r^2(1 - \alpha) \}$$

Asymptotically, (random) confidence region contains the true parameter value  $\boldsymbol{\theta}_{10}$  with probability  $1 - \alpha$  if  $\boldsymbol{\theta}_{10}$  is the true value. If  $\boldsymbol{\theta}_{10} \notin C_{1-\alpha}$  then we reject the null hypothesis

$$P(\boldsymbol{\theta}_{10} \in C_{1-\alpha}) = P\{T(\boldsymbol{\theta}_{10}) \leq \chi_r^2(1 - \alpha)\}$$

*Example 1: Binomial model*

We previously derived the test statistics for the binomial model as  $T_W = \frac{n(\hat{p}-p)^2}{\hat{p}(1-\hat{p})}$  the solving

$$T_W = \frac{n(\hat{p}-p)^2}{\hat{p}(1-\hat{p})} \leq \chi_1^2(1-\alpha) = z_{1-\alpha/2}^2 \implies \left( \hat{p} - z_{1-\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}, \hat{p} + z_{1-\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \right)$$

### 3.4 Nonstandard hypothesis testing problems

If  $\frac{\partial}{\partial \theta} \ell(\theta | \mathbf{Y})|_{\theta=\hat{\theta}_{MLE}} \neq 0$  in probability, typically the asymptotic normality of  $\hat{\theta}_{MLE}$  is no longer true and the 3 likelihood-based tests do not have a limiting  $\chi_r^2$  null distribution. *Example 1: Exponential threshold model*

$$Y_1, \dots, Y_n \stackrel{iid}{\sim} f(y; \mu) = \begin{cases} e^{-(y-\mu)} & y > \mu \\ 0 & \text{otherwise} \end{cases}$$

and the MLE is  $\hat{\mu}_{MLE} = Y_{(1)}$ . In this case, the score equation is not 0.

**Null hypotheses on the boundary of the parameter space**

$$\begin{aligned} Y_1, \dots, Y_n &\stackrel{iid}{\sim} N(\mu, 1) \quad \mu \geq \mu_0 \\ \hat{\mu}_{MLE} &= \max(\bar{Y}, \mu_0) \\ H_0 : \mu &= \mu_0 \quad \text{vs.} \quad H_a : \mu > \mu_0 \\ T_W &= n(\hat{\mu}_{MLE} - \mu_0)^2 \\ &= \begin{cases} n(\bar{Y} - \mu_0)^2 & \bar{Y} \geq \mu_0 \\ 0 & \text{otherwise} \end{cases} \\ T_S &= n(\bar{Y} - \mu_0)^2 \sim \chi_1^2 \text{ under } H_0 \\ T_{LR} &= T_W \end{aligned}$$

When a null hypothesis value, say  $\theta_0$  lies on the boundary of the parameter space, then maximum likelihood estimators are often truncated at that boundary because by definition  $\hat{\theta}_{MLE}$  must lie in the parameter space of  $\theta$ . Thus  $\hat{\theta}_{MLE}$  equal to the boundary value  $\theta_0$  with positive probability and correspondingly  $T_{LR}$  zero for those cases. The result is that the limiting distribution of  $T_{LR}$  is a mixture of a point mass at zero and a chi-squared distribution.

In this case,  $\bar{Y}$  used to be greater than  $\mu_0$  with probability  $\frac{1}{2}$ . Therefore  $T_w$  or  $T_{LR}$  has a null distribution that is an equal mixture of a point mass at 0 and a  $\chi_1^2$ :

$$Z^2 I(Z > 0), \quad Z \sim N(0, 1)$$

with this modified limiting distribution. The confidence interval shall be modified accordingly: level-  $\alpha$  test: reject  $H_0$  with critical value being  $1 - 2\alpha$  quantile of  $\chi_1^2$

$$\begin{aligned} P \{ Z^2 I(Z > 0) \geq \chi_1^2(1 - 2\alpha) \} &= P \{ Z^2 \geq \chi_1^2(1 - 2\alpha), Z > 0 \} \\ &= \frac{1}{2} P \{ Z^2 \geq \chi_1^2(1 - 2\alpha) \} = \alpha \end{aligned}$$

Although the score test statistics still follows the chi-square distribution. By following standard approach we can reject if

$$\bar{Y} < \mu_0 - \sqrt{\chi_1^2(\alpha)/n} \text{ or } \bar{Y} > \mu_0 + \sqrt{\chi_1^2(\alpha)/n}$$

however, the  $\bar{Y} < \mu_0 - \sqrt{\chi_1^2(\alpha)/n}$  part doesn't make sense cause the parameter space is  $[\theta_0, \infty)$  thus a natural solution is to reject when  $\sqrt{n}(\bar{Y} - \mu_0) > z_\alpha$

For cases that are not in the same form as this example, we can try to transfer the test statistics to be like that in this example.

*Example 2: Isotonic regression*

There are  $k$  independent normal samples of size  $n_1, \dots, n_k$ , each i.i.d with means  $\mu_1, \dots, \mu_k$ , respectively, and common variance  $\sigma^2$ , satisfying

$$\mu_1 \leq \mu_2 \leq \dots \leq \mu_k$$

The MLE minimizes  $\sum_{i=1}^k n_i (\bar{Y}_i - \mu_i)^2$  subject to  $\mu_1 \leq \mu_2 \leq \dots \leq \mu_k$ . We'd like to test

$$\begin{aligned} H_0 : & \mu_1 = \dots = \mu_k \\ H_a : & \mu_1 \leq \dots \leq \mu_k \text{ with at least 1 strict inequality} \end{aligned}$$



we can reparameterize as

$$\begin{aligned} & \mu_1 \\ \Delta\mu_2 &= \mu_2 - \mu_1 \\ & \vdots \\ \Delta\mu_k &= \mu_k - \mu_{k-1} \end{aligned}$$

Then, this results in null hypothesis on the boundary:

$$\begin{aligned} H_0 : & \Delta\mu_2 = \dots = \Delta\mu_k = 0 \\ H_a : & \Delta\mu_2 \geq 0, \dots, \Delta\mu_k \geq 0 \text{ with at least 1 strict inequality} \end{aligned}$$

## 4 Bayesian methods

★ this part will not be tested

★ this part of the note is combined with notes from *BIS 567 Bayesian Statistics* at Yale to enhance understanding

### 4.1 Introduction

- Frequentist approach
  - the unknown parameter  $\theta$  is assume to be constant
  - data  $\mathbf{Y} \sim f(\mathbf{y}; \theta)$  is considered random
  - estimation can be via MLE, M-estimation, ...
  - hypothesis testing can be likelihood-based,...
  - confidence interval is derived by inverting a test statistic
- Bayesian approach
  - the unknown parameter  $\theta$  is from  $\theta \sim \pi(\theta)$ , the prior distribution
  - data is from  $\mathbf{Y} \sim f(\mathbf{y} | \theta)$
  - estimation is via the posterior density

$$\pi(\theta | \mathbf{Y} = \mathbf{y}) = \frac{f(\mathbf{y} | \theta)\pi(\theta)}{\int f(\mathbf{y} | \theta)\pi(\theta)d\theta}$$

point estimation can use ,e.g., posterior mean

- $1 - \alpha$  credible region: parameter space with posterior probability  $1 - \alpha$

*Marginal density (prior predictive density) of Y*

$$\mathbf{Y}, m(\mathbf{y}) = \int f(\mathbf{y} | \theta)\pi(\theta)d\theta$$

*Posterior predictive density*

$$\begin{aligned} m(\mathbf{y}_{new} | \mathbf{Y}) &= \int f(\mathbf{y}_{new} | \theta, \mathbf{Y}) \pi(\theta | \mathbf{Y})d\theta \\ &= \int f(\mathbf{y}_{new} | \theta) \pi(\theta | \mathbf{Y})d\theta \text{ when } \mathbf{Y}_{new} \perp \mathbf{Y} \end{aligned}$$

Where does the prior  $\pi(\theta)$  come from?

- Subjective Bayesian: personal uncertainty about  $\theta$
- From previous knowledge: previous information about  $\theta$  (Bayesian analysis is used to combine previous info with current data)
- For convenience: convenient technical density to employ the Bayesian machinery

Example 1:  $Y \sim \text{binomial}(n; p)$

$$\begin{aligned}
 f(y | p) &= \binom{n}{p} p^y (1-p)^{n-y} \quad y = 0, 1, \dots, n \\
 \pi(p) &= \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} p^{\alpha-1} (1-p)^{\beta-1} \text{ (Beta)} \\
 m(y) &= \int f(y | p) \pi(p) dp = \frac{B(y + \alpha, n - y + \beta)}{B(\alpha, \beta)} \binom{n}{y} \\
 \pi(p | \mathbf{Y}) &= \frac{f(y | p) \pi(p)}{m(y)} = \frac{\binom{n}{p} p^y (1-p)^{n-y} \cdot \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} p^{\alpha-1} (1-p)^{\beta-1}}{m(y)} \\
 &= \frac{\binom{n}{p} \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} p^{y+\alpha-1} (1-p)^{n-y+\beta-1}}{m(y)} \sim \text{beta}(Y + \alpha, n - Y + \beta)
 \end{aligned}$$

the posterior mean is

$$E(p | \mathbf{Y}) = \frac{Y + \alpha}{\alpha + \beta + n} = \frac{\alpha + \beta}{\alpha + \beta + n} \frac{\alpha}{\alpha + \beta} + \frac{n}{\alpha + \beta + n} \frac{Y}{n}$$

which shrinks the MLE toward prior mean. Bayesian estimator may have smaller MSE than the MLE. When  $\alpha = \beta = 0$ , the prior is improper

$$\implies \text{posterior beta}(Y, n - Y) \text{ is proper if } 1 \leq Y \leq n - 1 \implies \text{posterior mean } Y/n$$

Example 2: Normal  $(\theta, \sigma_Y^2)$

$$\begin{aligned}
 Y_1, \dots, Y_n &\stackrel{\text{iid}}{\sim} N(\theta, \sigma_Y^2) \quad \sigma_Y^2 \text{ known} \\
 \text{prior } \theta &\sim N(\mu_0, \sigma_0^2) \quad \mu_0, \sigma_0^2 \text{ known} \\
 \text{posterior } \theta | \mathbf{Y} &\sim N\left(\frac{\tau_0 \mu_0 + \tau_n \bar{Y}}{\tau_0 + \tau_n}, \frac{1}{\tau_0 + \tau_n}\right)
 \end{aligned}$$

where precisions  $\tau_0 = 1/\sigma_0^2$  and  $\tau_n = n/\sigma_Y^2$   
95% credible interval

$$\left( \frac{\tau_0 \mu_0 + \tau_n \bar{Y}}{\tau_0 + \tau_n} - \frac{1.96}{\sqrt{\tau_0 + \tau_n}}, \frac{\tau_0 \mu_0 + \tau_n \bar{Y}}{\tau_0 + \tau_n} + \frac{1.96}{\sqrt{\tau_0 + \tau_n}} \right)$$

A nice feature of independent data is that one can sequentially update the prior for each additional datum, or all at once. For example, suppose  $Y_1, \dots, Y_n$  are independent with respective densities  $f_i(y_i | \theta), i = 1, \dots, n$ . Then, the posterior is

$$\begin{aligned}
 \pi(\theta | \mathbf{Y}) &= \frac{\pi(\theta) \prod_{i=1}^n f_i(Y_i | \theta)}{\int \pi(\theta) \prod_{i=1}^n f_i(Y_i | \theta) d\theta} \\
 &= \frac{\pi(\theta | Y_1) \prod_{i=2}^n f_i(Y_i | \theta)}{\int \pi(\theta | Y_1) \prod_{i=2}^n f_i(Y_i | \theta) d\theta} \\
 &= \frac{\pi(\theta | Y_1, \dots, Y_k) \prod_{i=k+1}^n f_i(Y_i | \theta)}{\int \pi(\theta | Y_1, \dots, Y_k) \prod_{i=k+1}^n f_i(Y_i | \theta) d\theta}
 \end{aligned}$$

Here  $\pi(\theta | Y_1)$  is the posterior from using the prior and  $Y_1$ . It then is used as the prior for the remaining data. Sufficient statistics often make calculations easier. Because of the factorization theorem, if a sufficient statistic for  $\theta$  exists, then the posterior depends on the data only through the sufficient statistic.

Estimation and inference: Frequentist uses sampling distribution, whereas Bayesian uses posterior density.

## 4.2 Bayesian estimator from decision theory perspective

Loss function  $L\{\theta, \delta(\mathbf{Y})\}$ : non-negative function of the true parameter and an estimator  $\delta(\mathbf{Y})$ . For example:

1. squared error loss:  $L\{\theta, \delta(\mathbf{Y})\} = \|\theta - \delta(\mathbf{Y})\|^2$
2. absolute error loss:  $L\{\theta, \delta(\mathbf{Y})\} = |\theta - \delta(\mathbf{Y})|$
3. 0-1 loss:  $L\{\theta, \delta(\mathbf{Y})\} = \begin{cases} 0 & \text{if } \theta = \delta(\mathbf{Y}) \\ 1 & \text{otherwise} \end{cases}$

*Risk*: is the average loss over  $Y$

$$R(\boldsymbol{\theta}, \boldsymbol{\delta}) = \int L\{\boldsymbol{\theta}, \boldsymbol{\delta}(\mathbf{y})\}f(\mathbf{y} | \boldsymbol{\theta})d\mathbf{y}$$

which is not a single number, but rather a function of  $\boldsymbol{\theta}$  Several standard frequentist approaches regarding the risk:

1. minimize  $R(\boldsymbol{\theta}, \boldsymbol{\delta})$ , with squared error loss, over the class of unbiased estimators  $\implies$  if exists, minimum variance unbiased estimator (MVUE)
2. minimax estimator

$$\inf_{\boldsymbol{\delta}} \sup_{\boldsymbol{\theta}} R(\boldsymbol{\theta}, \boldsymbol{\delta})$$

maximize the risk over  $\boldsymbol{\theta}$  first then seek an estimator that minimizes the maximum risk

Bayesian approach (take another average over  $\boldsymbol{\theta}$ )

$$R_{\text{Bayes}}(\boldsymbol{\pi}, \boldsymbol{\delta}) = \int R(\boldsymbol{\theta}, \boldsymbol{\delta})\boldsymbol{\pi}(\boldsymbol{\theta})d\boldsymbol{\theta}$$

Bayesian estimator  $\boldsymbol{\delta}_{\text{Bayes}}$  minimizes  $R_{\text{Bayes}}(\boldsymbol{\pi}, \boldsymbol{\delta})$

We can also minimize the *posterior risk* instead

$$\rho\{\boldsymbol{\pi}, \boldsymbol{\delta}(\mathbf{Y})\} = \int L\{\boldsymbol{\theta}, \boldsymbol{\delta}(\mathbf{Y})\}\boldsymbol{\pi}(\boldsymbol{\theta} | \mathbf{Y})d\boldsymbol{\theta}$$

But typically  $\arg \min \rho\{\boldsymbol{\pi}, \boldsymbol{\delta}(\mathbf{Y})\} \equiv \arg \min R_{\text{Bayes}}(\boldsymbol{\pi}, \boldsymbol{\delta})$

$$\begin{aligned} R_{\text{Bayes}}(\boldsymbol{\pi}, \boldsymbol{\delta}) &= \iint L\{\boldsymbol{\theta}, \boldsymbol{\delta}(\mathbf{y})\}f(\mathbf{y} | \boldsymbol{\theta})\boldsymbol{\pi}(\boldsymbol{\theta})d\mathbf{y}d\boldsymbol{\theta} \\ &= \iint L\{\boldsymbol{\theta}, \boldsymbol{\delta}(\mathbf{y})\}\boldsymbol{\pi}(\boldsymbol{\theta} | \mathbf{y})m(\mathbf{y})d\mathbf{y}d\boldsymbol{\theta} \\ &= \int \rho\{\boldsymbol{\pi}, \boldsymbol{\delta}(\mathbf{y})\}m(\mathbf{y})d\mathbf{y} \end{aligned}$$

thus minimize  $\rho\{\boldsymbol{\pi}, \boldsymbol{\delta}(\mathbf{Y})\}$  usually minimize  $R_{\text{Bayes}}(\boldsymbol{\pi}, \boldsymbol{\delta})$

Different loss function leads to different parameter estimate

- squared error loss  $\implies$  posterior mean
- absolute error loss  $\implies$  posterior median
- 0-1 loss  $\implies$  posterior mode

Bayes estimators with proper priors are generally not unbiased in the frequentist sense. However, they typically have good risk behavior in the frequentist sense. Interestingly, a key technique for finding minimax estimators starts with a Bayes estimator (see Lehmann and Casella, 1998, Ch. 5). Moreover, an admissible estimator (an estimator not uniformly larger in risk compared to any other estimator) must be a Bayes estimator or the limit of Bayes estimators. Thus, Bayes estimators are not only good in terms of Bayes risk but are often of interest to frequentists willing to sacrifice unbiasedness.

### 4.3 Credible intervals

$1 - \alpha$  *credible interval/ region*: a region of  $\boldsymbol{\theta}$  of posterior prob  $1 - \alpha$

*highest posterior density (HPD) region*: the region with minimized volume

region with equal tail probability is often used for a scalar parameter

### 4.4 Conjugate prior

*Conjugate prior*: when the data  $Y$  has density  $f(y | \boldsymbol{\theta})$  and the prior and posterior are from the same family of densities, we say that the prior is conjugate.

$\boldsymbol{\pi}(\boldsymbol{\theta})$  governed by fixed hyperparameters

$\boldsymbol{\gamma}_{\text{prior}}$   $\boldsymbol{\pi}(\boldsymbol{\theta} | \mathbf{Y})$  has updated hyperparameters  $\boldsymbol{\gamma}_{\text{post}}$ , via a known function of  $\boldsymbol{\gamma}_{\text{prior}}$  and  $\mathbf{Y}$

Some examples:

- beta( $\alpha, \beta$ ) prior for binomial ( $n, p$ ) data
- $N(\mu_0, \sigma_0^2)$  prior for  $N(\boldsymbol{\theta}, \sigma_Y^2)$  data (known  $\sigma_Y^2$ )

- beta( $\alpha, \beta$ ) prior for negative binomial data
- gamma prior for Poisson data
- gamma prior for gamma data
- Pareto prior for Uniform ( $0, \theta$ ) data

Example 1:

$$\begin{aligned} \mathbf{Y} &\sim \text{multinomial}(n; p_1, \dots, p_k) \\ \pi(p_1, \dots, p_k) &\sim \text{Dirichlet}(\alpha_1, \dots, \alpha_k) \\ \pi(\mathbf{p}) &= \frac{\prod_{i=1}^k p_i^{\alpha_i-1}}{B(\boldsymbol{\alpha})}, \quad B(\boldsymbol{\alpha}) = \frac{\prod_{i=1}^k \Gamma(\alpha_i)}{\Gamma(\sum_{i=1}^k \alpha_i)} \\ \pi(\mathbf{p} | \mathbf{Y}) &\propto \prod_{i=1}^k p_i^{N_i} \prod_{i=1}^k p_i^{\alpha_i-1} \\ \implies \pi(\mathbf{p} | \mathbf{Y}) &\sim \text{Dirichlet}(\alpha_1 + N_1, \dots, \alpha_k + N_k) \end{aligned}$$

Any data density having a sufficient statistic of fixed dimension  $\forall n$  has a conjugate prior.

Example 2: exponential family canonical form

$$\begin{aligned} f(\mathbf{y}; \boldsymbol{\eta}) &= h(\mathbf{y}) \exp \left\{ \sum_{i=1}^s \eta_i T_i(\mathbf{y}) - A(\boldsymbol{\eta}) \right\} \\ \text{conjugate } \pi(\boldsymbol{\eta} | \boldsymbol{\gamma}, \lambda) &= K(\boldsymbol{\gamma}, \lambda) \exp \left\{ \sum_{i=1}^s \gamma_i \eta_i - \lambda A(\boldsymbol{\eta}) \right\} \\ \text{posterior} &\sim \pi(\boldsymbol{\eta} | \boldsymbol{\gamma} + \mathbf{T}(\mathbf{y}), \lambda + 1) \end{aligned}$$

## 4.5 Noninformative prior

Truly noninformative may not be possible. For example, assigning prior probability to events  $A, B$ , and  $C = C_1 \cup C_2$ . The choice between the following two cases reflects some kind of information

$$\begin{aligned} P(A) = P(B) = P(C) &= 1/3 \text{ or} \\ P(A) = P(B) = P(C_1) = P(C_2) &= 1/4? \end{aligned}$$

For the case of a location parameter  $\mu$  taking values on  $(-\infty, \infty)$ : improper prior  $\pi(\mu) = 1$  giving equal weight to all values can be justified on a variety of grounds. An improper prior does not have a finite integral. However, there seems to be no philosophical problem with improper priors as long as they lead to proper posteriors.

For the case of a scale parameter  $\sigma$  taking values in  $(0, \infty)$ : Jeffreys' suggestion  $\pi(\sigma) = 1/\sigma$  which has an invariance argument. His invariance argument is that any power transformation of  $\sigma$ , say  $\gamma = \sigma^a$  has via a change-of-variables, the improper density

$$\gamma = \sigma^a \implies \pi(\gamma) = \frac{1}{\gamma^{1/a}} \left| \frac{\gamma^{1/a-1}}{a} \right| = \frac{1}{a\gamma}$$

which is similar in form to  $1/\sigma$

Combining these last two improper priors, a location-scale family with  $(\mu, \sigma) \in (-\infty, \infty) \times (0, \infty)$ , suggests using the improper prior

$$\pi(\mu, \sigma) \propto \frac{1}{\sigma}$$

Now, moving to the case of general parameters on continuous parameter spaces. For Jeffreys prior

- Scalar parameter  $\pi(\theta) \propto I(\theta)^{1/2}$ . For transformed parameter  $\gamma = g(\theta)$ ,

$$\pi(\gamma) \propto \frac{1}{|g' \{g^{-1}(\gamma)\}|} I \{g^{-1}(\gamma)\}^{1/2} = I(\gamma)^{1/2}$$

Example 1: for binomial( $n, p$ ).  $I(p) = n/\{p(1-p)\}$

$$\pi(p) \propto \{p(1-p)\}^{-1/2} \sim \text{beta}(1/2, 1/2)$$

- Vector parameter  $\pi(\boldsymbol{\theta}) \propto |\mathbf{I}(\boldsymbol{\theta})|^{1/2}$ ;  $|\cdot|$  is determinant

*Example 1:* for multinomial( $n; p_1, \dots, p_k$ )

$$\mathbf{I}_T(\mathbf{p}) = n \{ \text{diag}(1/p_1, \dots, 1/p_{k-1}) + \mathbf{1}\mathbf{1}^T/p_k \}$$

$$\pi(\mathbf{p}) \propto \frac{1}{\sqrt{p_1 \cdots p_k}} \sim \text{Dirichlet}(1/2, \dots, 1/2)$$

For  $N(\mu, \sigma^2)$  case with  $\boldsymbol{\theta} = (\mu, \sigma)$ ,  $\mathbf{I}(\boldsymbol{\theta}) = \text{diag}(\sigma^{-2}, 2\sigma^{-2}) \implies \pi(\mu, \sigma) \propto \sigma^{-2}$  which is different from  $\sigma^{-1}$  given earlier. Thus, Jeffreys modified his original proposal in the presence of location parameters say  $\mu_1, \dots, \mu_k$  to

$$\pi(\mu_1, \dots, \mu_k, \boldsymbol{\theta}) \propto |\mathbf{I}(\boldsymbol{\theta})|^{1/2}$$

where  $\mu'_i$ 's are held as fixed.

## 4.6 Normal data examples

### 4.6.1 One sample with unknown mean and variance

Data:  $Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2)$ ,  $\boldsymbol{\theta} = (\mu, \tau = 1/\sigma^2)$  unknown

Prior:  $\pi(\mu | \tau) \sim N(\mu_0, (\tau n_0)^{-1})$

$\pi(\tau) \sim \text{gamma}(\alpha_0, 1/\beta_0)$ ,  $\text{mean} = \alpha_0/\beta_0$

a rough way to estimate the prior variance is treating  $\tau = \alpha_0/\beta_0 \implies$  prior variance of  $\mu \approx (\beta_0/\alpha_0) n_0^{-1}$

$\mu \sim t$  with center  $\mu_0$ , scale  $e^2 = (\beta_0/\alpha_0)/n_0$ , and df  $2\alpha_0$

then the likelihood is

$$(2\pi\sigma^2)^{-n/2} \exp \left[ -\frac{1}{2\sigma^2} \left\{ \sum (Y_i - \bar{Y})^2 + n(\bar{Y} - \mu)^2 \right\} \right]$$

then the posterior is

$$\pi(\mu, \tau | \mathbf{Y}) \propto \tau^{1/2} \exp \left\{ -\frac{\tau n'}{2} (\mu - \mu')^2 \right\} \tau^{\alpha' - 1} e^{-\tau \beta'}$$

$$\mu' = \frac{n_0 \mu_0 + n \bar{Y}}{n_0 + n}$$

$$n' = n_0 + n$$

$$\alpha' = n/2 + \alpha_0$$

$$\beta' = \beta_0 + \frac{\sum (Y_i - \bar{Y})^2}{2} + \frac{1}{2} \left( \frac{n_0 n}{n_0 + n} \right) (\mu_0 - \bar{Y})^2$$

$$\implies \mu, \tau | \mathbf{Y} \sim N(\mu', (\tau n')^{-1}) \times \text{gamma}(\alpha', 1/\beta') \text{ (conjugate prior)}$$

$$\mu | \mathbf{Y} \sim t \text{ with center } \mu', \text{ scale } e^2 = (\beta'/\alpha')/n', \text{ and df } 2\alpha'.$$

posterior predictive density

$$m(y_{n+1} | Y_1, \dots, Y_n) = \int f(y_{n+1} | \boldsymbol{\theta}) \pi(\boldsymbol{\theta} | Y_1, \dots, Y_n) d\boldsymbol{\theta}$$

which is a  $t$  distribution with mean  $\mu'$ , scale  $e^2 = (\beta'/\alpha')(1 + 1/n')$ , and  $df = 2\alpha'$ .

Letting  $\alpha_0 \rightarrow -1/2$ ,  $\beta_0 \rightarrow 0$ , and  $n_0 \rightarrow 0$  s.t.  $(n_0/\beta_0)^{1/2} \rightarrow 1$

$\implies \pi(\mu, \tau) \rightarrow 1/\tau$ , an improper prior

$\implies \mu | \mathbf{Y} \sim t$  with mean  $\bar{Y}$ , scale  $e^2 = s_{n-1}^2/n$ , and df  $n - 1$

$\implies$  Bayes estimator (with squared error loss) is  $\bar{Y}$  credible interval  $\equiv$  usual frequentist  $t$  interval

### 4.6.2 Two samples

$$X_1, \dots, X_m \stackrel{\text{iid}}{\sim} N(\mu_1, \sigma^2), \quad Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} N(\mu_2, \sigma^2), \quad \boldsymbol{\theta} = (\mu_1, \mu_2, \tau = 1/\sigma^2)^T$$

$\Delta = \mu_1 - \mu_2$  is of interest

prior

$$\begin{aligned} \pi(\mu_1 | \tau) &\sim N\left\{\mu_{10}, (\tau m_0)^{-1}\right\} \\ \pi(\mu_2 | \tau) &\sim N\left\{\mu_{20}, (\tau n_0)^{-1}\right\} \\ \pi(\tau) &\sim \text{gamma}(\alpha_0, 1/\beta_0) \end{aligned}$$

with similar derivations as in the last example, we can get the posterior

$$\boldsymbol{\theta} | \mathbf{X}, \mathbf{Y} \sim N\left\{\mu'_1, (\tau m')^{-1}\right\} \times N\left\{\mu'_2, (\tau n')^{-1}\right\} \times \text{gamma}(\alpha', 1/\beta')$$

$$m' = m_0 + m\mu'_1 = (m_0\mu_{10} + m\bar{X}) / (m_0 + m)$$

$$n' = n_0 + n\mu'_2 = (n_0\mu_{20} + n\bar{Y}) / (n_0 + n)$$

$$\alpha' = \alpha_0 + m/2 + n/2$$

$$\beta' = \beta_0 + \frac{\sum (X_i - \bar{X})^2}{2} + \frac{1}{2} \left( \frac{m_0 m}{m_0 + m} \right) (\mu_{10} - \bar{X})^2 + \frac{\sum (Y_i - \bar{Y})^2}{2} + \frac{1}{2} \left( \frac{n_0 n}{n_0 + n} \right) (\mu_{20} - \bar{Y})^2$$

$$\implies \Delta, \tau | \mathbf{X}, \mathbf{Y} \sim N\{\Delta' = \mu'_1 - \mu'_2, 1/(\tau m') + 1/(\tau n')\} \times \text{gamma}(\alpha', 1/\beta')$$

$$\implies \Delta | \mathbf{X}, \mathbf{Y} \sim t \text{ with center } \Delta', \text{ scale }^2 = (\beta'/\alpha')(m' + n') / (m'n'), \text{ and df } 2\alpha'$$

Letting  $\alpha_0 \rightarrow -1, \beta_0 \rightarrow 0, m_0 \rightarrow 0, n_0 \rightarrow 0$  such that  $(m_0 n_0)^{1/2} / \beta_0 \rightarrow 1$

$\rightarrow$  improper prior  $\pi(\boldsymbol{\theta}) = 1/\tau$  (Jeffrey's prior)

$\implies \Delta | \mathbf{X}, \mathbf{Y} \sim t$  with mean  $\bar{X} - \bar{Y}$ , scale  $^2 = s_p^2(1/m + 1/n)$ , and  $df = m + n - 2$ , where  $s_p^2$  is the usual pooled estimated variance.

### 4.6.3 Normal linear model

$$\mathbf{Y}_{(n \times 1)} = \mathbf{X}_{(n \times p)} \boldsymbol{\Delta}_{(p)} + \mathbf{e}_{(n \times 1)}, \quad e_1, \dots, e_n \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$$

Normal-gamma prior:

$$\boldsymbol{\Delta} | \tau = 1/\sigma^2 \sim N(\boldsymbol{\Delta}_0, \boldsymbol{\Sigma}_0^{-1}/\tau), \quad \tau \sim \text{gamma}(\alpha_0, 1/\beta_0)$$

Posterior is also Normal-gamma:

$$\boldsymbol{\Delta} | \tau, \mathbf{Y} \sim N\left(\boldsymbol{\Delta}', (\mathbf{X}^T \mathbf{X} + \boldsymbol{\Sigma}_0)^{-1} / \tau\right), \quad \tau | \mathbf{Y} \sim \text{gamma}(\alpha', 1/\beta')$$

where

$$\boldsymbol{\Delta}' = (\mathbf{X}^T \mathbf{X} + \boldsymbol{\Sigma}_0)^{-1} (\mathbf{X}^T \mathbf{Y} + \boldsymbol{\Sigma}_0 \boldsymbol{\Delta}_0)$$

$$\alpha' = \alpha_0 + n/2$$

$$\beta' = \beta_0 + \frac{1}{2} \left\{ (\mathbf{Y} - \mathbf{X} \boldsymbol{\Delta}')^T \mathbf{Y} + (\boldsymbol{\Delta}_0 - \boldsymbol{\Delta}')^T \boldsymbol{\Sigma}_0 \boldsymbol{\Delta}_0 \right\}$$

$\boldsymbol{\Delta}'$  is a weighted average of the OLS estimator  $(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$  and prior mean  $\boldsymbol{\Delta}_0$

$\boldsymbol{\Delta} | \mathbf{Y}$  is multivariate  $t$  with center  $\boldsymbol{\Delta}'$ , scale matrix  $(\beta'/\alpha') (\mathbf{X}^T \mathbf{X} + \boldsymbol{\Sigma}_0)^{-1}$ , and  $df \ 2\alpha'$ . For a subset of  $\boldsymbol{\Delta}$ , the marginal posterior is also a multivariate  $t$ .

let  $\boldsymbol{\Delta}_0 = \mathbf{0}, \boldsymbol{\Sigma}_0 = d\mathbf{I}_p$  with constant  $d \implies \boldsymbol{\Delta}' = (\mathbf{X}^T \mathbf{X} + d\mathbf{I}_p)^{-1} \mathbf{X}^T \mathbf{Y}$ , the ridge regression estimator.

Letting  $\boldsymbol{\Sigma}_0 \rightarrow \mathbf{0}, \alpha_0 \rightarrow -p/2$ , and  $\beta_0 \rightarrow 0$

$\implies \pi(\boldsymbol{\Delta}, \tau) = 1/\tau$

$\implies \boldsymbol{\Delta}' = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$ , and  $\boldsymbol{\Delta} | \mathbf{Y}$  is a  $p$ -dimensional  $t$  with  $df \ n - p$  and scale matrix  $s^2 (\mathbf{X}^T \mathbf{X})^{-1}$ , where  $s^2 = (n - p)^{-1} \sum (Y_i - \mathbf{x}_i^T \boldsymbol{\Delta}')^2$

## 4.7 Hierarchical Bayes and empirical Bayes

Previously we discussed that case where prior  $\pi(\boldsymbol{\theta})$  has been specified fully, say,  $\pi(\boldsymbol{\theta} | \boldsymbol{\alpha})$  with hyperparameter  $\boldsymbol{\alpha} = \boldsymbol{\alpha}_0$  given. What if unsure about  $\boldsymbol{\alpha}_0$ ? We turn to Hierarchical Bayes.

*Hierarchical Bayes:* Apart from  $\pi(\boldsymbol{\theta} | \boldsymbol{\alpha})$ , we also specify a hyperprior  $h(\boldsymbol{\alpha}) = h(\boldsymbol{\alpha} | \boldsymbol{\gamma}_0)$  with  $\boldsymbol{\gamma}_0$  given. The resulting prior

$$\begin{aligned}\pi(\boldsymbol{\theta}) &= \pi(\boldsymbol{\theta} | \boldsymbol{\gamma}_0) = \int \pi(\boldsymbol{\theta} | \boldsymbol{\alpha}) h(\boldsymbol{\alpha} | \boldsymbol{\gamma}_0) d\boldsymbol{\alpha} \\ m(\mathbf{y} | \boldsymbol{\alpha}) &= \int f(\mathbf{y} | \boldsymbol{\theta}) \pi(\boldsymbol{\theta} | \boldsymbol{\alpha}) d\boldsymbol{\theta} \\ \pi(\boldsymbol{\theta} | \boldsymbol{\alpha}, \mathbf{Y}) &= \frac{f(\mathbf{Y} | \boldsymbol{\theta}) \pi(\boldsymbol{\theta} | \boldsymbol{\alpha})}{m(\mathbf{Y} | \boldsymbol{\alpha})} \\ \pi(\boldsymbol{\alpha} | \mathbf{Y}) &= \frac{m(\mathbf{Y} | \boldsymbol{\alpha}) h(\boldsymbol{\alpha})}{\iint f(\mathbf{Y} | \boldsymbol{\theta}) \pi(\boldsymbol{\theta} | \boldsymbol{\alpha}) h(\boldsymbol{\alpha}) d\boldsymbol{\alpha} d\boldsymbol{\theta}} \\ \pi(\boldsymbol{\theta} | \mathbf{Y}) &= \frac{\int f(\mathbf{Y} | \boldsymbol{\theta}) \pi(\boldsymbol{\theta} | \boldsymbol{\alpha}) h(\boldsymbol{\alpha}) d\boldsymbol{\alpha}}{\iint f(\mathbf{Y} | \boldsymbol{\theta}) \pi(\boldsymbol{\theta} | \boldsymbol{\alpha}) h(\boldsymbol{\alpha}) d\boldsymbol{\alpha} d\boldsymbol{\theta}} \\ &= \int \pi(\boldsymbol{\theta} | \boldsymbol{\alpha}, \mathbf{Y}) \pi(\boldsymbol{\alpha} | \mathbf{Y}) d\boldsymbol{\alpha}\end{aligned}$$

*Empirical Bayes:* use MLE  $\hat{\boldsymbol{\alpha}}$  from marginal likelihood  $m(\mathbf{Y} | \boldsymbol{\alpha})$  (other estimators may also be used).

$$\pi(\boldsymbol{\alpha} | \mathbf{Y}) \text{ highly peaked at } \hat{\boldsymbol{\alpha}} \implies \pi(\boldsymbol{\theta} | \mathbf{Y}) \approx \pi(\boldsymbol{\theta} | \hat{\boldsymbol{\alpha}}, \mathbf{Y}) \text{ (taking } \pi(\hat{\boldsymbol{\alpha}} | \mathbf{Y}) \approx 1)$$

Empirical Bayes posterior has a similar mean, but smaller variability than the full Bayes.

*Example 1:* one-way normal random effects model

$$Y_{ij}, i = 1, \dots, k; j = 1, \dots, n_i$$

Frequentist methods: fixed and random effects ANOVA.

Bayesian with fixed effects: earlier normal linear model. A Bayesian analogue of the random effects model:

$$\begin{aligned}Y_{ij} | \theta_i, \sigma_e^2 &\sim N(\theta_i, \sigma_e^2) \text{ given } \theta_1, \dots, \theta_k, \sigma_e^2, Y_{ij} \text{ 's are mutually independent} \\ \theta_1, \dots, \theta_k | \boldsymbol{\alpha} &= (\mu, \sigma_a) \text{ i.i.d } N(\mu, \sigma_a^2) \\ \pi(\sigma_e^2) &\propto 1/\sigma_e^2 \text{ (Jeffrey's)} \\ h(\mu) &\propto 1 \\ h(\sigma_a) &\propto 1\end{aligned}$$

Of interest:

1. random effects population:  $\mu$  and  $\sigma_a$
2. individual  $\theta_i$

For the empirical Bayes approach, consider a simplified model with  $\sigma_e^2$  known. Obtain  $\hat{\boldsymbol{\alpha}}$  from the marginal likelihood  $m(\mathbf{y} | \boldsymbol{\alpha})$ . The posterior is  $\pi(\boldsymbol{\theta} | \hat{\boldsymbol{\alpha}}, \mathbf{Y})$ .

### 4.7.1 James-Stein estimation

$Y_1, \dots, Y_b$  are independent:  $Y_i \sim N(\theta_i, \sigma_0^2)$  with  $\sigma_0^2$  known.

While the parameters are unconnected in the model, the inferences regarding them are connected through the squared error loss  $\|\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}(\mathbf{Y})\|^2$ . Consider  $\hat{\boldsymbol{\theta}}_{\text{MLE}} = \mathbf{Y}$ ,  $R(\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}) = b\sigma_0^2$ .

Stein (1955) proved the remarkable result that for  $b \geq 3$ ,  $\hat{\boldsymbol{\theta}}_{\text{MLE}}$  is inadmissible. James and Stein (1961) provided a dominating estimator,

$$\hat{\boldsymbol{\theta}}_{\text{JS}} = \left\{ 1 - \frac{(b-2)\sigma_0^2}{\sum_{i=1}^b Y_i^2} \right\} \mathbf{Y}$$

shrinking  $\mathbf{Y}$  toward  $\mathbf{0}$ .

Expected squared error loss

$$\sum_{i=1}^b E \left( \hat{\theta}_{\text{JS},i} - \theta_i \right)^2 = b\sigma_0^2 - (b-2)^2\sigma_0^4 E \left( \frac{1}{\sum_{i=1}^b Y_i^2} \right)$$

which is less than the risk for  $\hat{\theta}_{MLE}$  provided  $b \geq 3$

Empirical Bayes interpretation:

$$\text{prior: } \boldsymbol{\theta} \sim MN(\mathbf{0}, \sigma_a^2 \mathbf{I}_b)$$

$\implies$  posterior:  $\boldsymbol{\theta} | \mathbf{Y} \sim$  normal with mean  $\left(1 - \frac{\sigma_0^2}{\sigma_0^2 + \sigma_a^2}\right) \mathbf{Y}$  marginal  $\mathbf{Y} \sim MN(\mathbf{0}, (\sigma_0^2 + \sigma_a^2) \mathbf{I}_b)$   $\sigma_a^2$  is unknown, but can be estimated from the fact

$$E \left\{ \frac{(b-2)\sigma_0^2}{\sum_{i=1}^b Y_i^2} \right\} = \frac{\sigma_0^2}{\sigma_0^2 + \sigma_a^2}$$

Thus,  $\hat{\theta}_{JS}$  may be viewed as an empirical Bayes estimator (Efron and Morris, 1972).

Alternative prior

$$\boldsymbol{\theta} \sim MN(\mu, \sigma_a^2 \mathbf{I}_b)$$

Bayes posterior mean  $B\mu + (1-B)\mathbf{Y}$ ,  $B = \frac{\sigma_0^2}{\sigma_0^2 + \sigma_a^2}$

$$\text{Marginal } \mathbf{Y} \sim N(\mu, (\sigma_0^2 + \sigma_a^2) \mathbf{I}_b)$$

unbiased estimators  $\hat{\mu} = \bar{Y}$   $\hat{B} = (b-3)\sigma_0^2 / \sum_{i=1}^b (Y_i - \bar{Y})^2$

$$\text{James-Stein estimator } \hat{B}\hat{\mu} + (1-\hat{B})\mathbf{Y},$$

shrinking towards sample mean  $\bar{Y}$  and having expected squared error loss less than  $\mathbf{Y}$  if  $b \geq 4$ .

#### 4.7.2 Meta-analysis applications of hierarchical and empirical Bayes

*Meta analysis:* analysis of a group of studies related to the same question of interest. Each study has a point effect estimate and its associated standard error.

*Hierarchical Bayes:*  $k$  studies.  $i$  th study - effect parameter  $\theta_i$  and data  $Y_i$ . Ignoring first-stage nuisance parameters, 3 levels of the model are

$$f(\mathbf{y} | \boldsymbol{\theta}) = f(y_1, \dots, y_k | \theta_1, \dots, \theta_k) = \prod_{i=1}^k f(y_i | \theta_i)$$

$$\pi(\boldsymbol{\theta} | \boldsymbol{\alpha}) = \prod_{i=1}^k f(\theta_i | \boldsymbol{\alpha})$$

$$h(\boldsymbol{\alpha})$$

The 2 nd and 3 rd stages  $\implies$  prior for  $\boldsymbol{\theta}$  :

$$\pi(\boldsymbol{\theta}) = \int \prod_{i=1}^k f(\theta_i | \boldsymbol{\alpha}) h(\boldsymbol{\alpha}) d\boldsymbol{\alpha}$$

Empirical Bayes:

marginal density of  $\mathbf{Y} \implies$  estimator  $\hat{\boldsymbol{\alpha}}$  for  $\boldsymbol{\alpha}$

$\implies$  estimated prior  $\pi(\boldsymbol{\theta} | \hat{\boldsymbol{\alpha}})$

$\implies$  estimated posterior  $\pi(\boldsymbol{\theta} | \hat{\boldsymbol{\alpha}}, \mathbf{Y})$

$\implies$  empirical Bayes estimate close to posterior mean from the hierarchical Bayes, but the variance is smaller.

*Example 1:* Normal models with known variance

Data:  $(Y_1, V_1), \dots, (Y_k, V_k)$ , where  $Y_i$  is approximately normal with variance estimated by  $V_i$  (but treated as known in analysis)

Frequentist fixed effect approach: all  $Y_i$ 's are approximately unbiased estimators of true effect  $\mu \implies$  weighted average  $\sum_{i=1}^k V_i^{-1} Y_i / \sum_{i=1}^k V_i^{-1}$  is approximately optimal. Assumption is not quite realistic!

Frequentist random effect approach:  $\theta_i \sim N(\mu, \sigma_a^2)$ . The likelihood of  $\mathbf{Y}$  or other methods can be used for the estimation of  $\mu$  and  $\sigma_a$ .

Hierarchical Bayesian:

$$Y_i | \theta_i \sim N(\theta_i, V_i)$$

$$\theta_i | \mu, \sigma_a \sim N(\mu, \sigma_a^2)$$

$$h(\mu, \sigma_a) \propto 1 \text{ (noninformative prior)}$$

Empirical Bayes: use estimator of  $(\mu, \sigma_a)$  (same as those from the frequentist random effect approach) without the hyperprior.

Full Bayes: with the hyperprior, MCMC is needed for the estimation.



## 4.8 Monte Carlo estimation of a posterior

(most materials from Carlin and Louis, 1996)

Main technical problem in Bayesian analysis: obtaining the posterior density and computing summary quantities such as posterior mean, standard deviation, and quantiles.

Which means we need to calculate integrals! As soon as we move away from conjugate priors and/or to hierarchical models.

One approach is to resort to asymptotic results. More focus is on Monte Carlo methods.

### 4.8.1 Noniterative Monte Carlo methods

#### Direct sampling:

Suppose that  $\theta \sim f(\theta)$  and  $\gamma = E\{h(\theta)\} = \int h(\theta)f(\theta)d\theta$  is of interest. Generate  $\theta_1, \dots, \theta_N \stackrel{iid}{\sim} f(\theta)$  and obtain

$$\hat{\gamma} = N^{-1} \sum_{j=1}^N h(\theta_j)$$

#### Indirect sampling:

##### 1. Importance sampling

Consider a posterior expectation

$$E\{h(\theta) \mid \mathbf{Y}\} = \frac{\int h(\theta)f(\mathbf{Y} \mid \theta)\pi(\theta)d\theta}{\int f(\mathbf{Y} \mid \theta)\pi(\theta)d\theta}$$

Suppose that we can roughly approximate the normalized likelihood times prior,  $cf(\mathbf{Y} \mid \theta)\pi(\theta)$ , by some density  $g(\theta)$  from which we can easily sample. Define weight function  $w(\theta) = f(\mathbf{Y} \mid \theta)\pi(\theta)/g(\theta)$

$$E\{h(\theta) \mid \mathbf{Y}\} = \frac{\int h(\theta)w(\theta)g(\theta)d\theta}{\int w(\theta)g(\theta)d\theta} \approx \frac{N^{-1} \sum_{j=1}^N h(\theta_j)w(\theta_j)}{N^{-1} \sum_{j=1}^N w(\theta_j)}$$

where  $\theta_j \stackrel{iid}{\sim} g(\theta)$ , the *importance function*.

The performance depends on how close  $g(\theta)$  resembles  $cf(\mathbf{Y} \mid \theta)\pi(\theta)$ .

##### 2. Rejection sampling

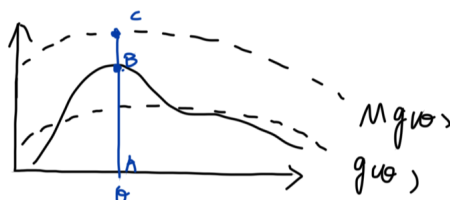
Consider posterior sampling

$$\pi(\theta \mid \mathbf{Y}) = \frac{f(\mathbf{Y} \mid \theta)\pi(\theta)}{\int f(\mathbf{Y} \mid \theta)\pi(\theta)d\theta}$$

Suppose  $\exists M > 0$  and a smooth density  $g(\theta)$ — envelope function - s.t.  $f(\mathbf{Y} \mid \theta)\pi(\theta) < Mg(\theta)$  then

- Generate  $\theta_j \sim g(\theta)$
- Generate  $U \sim \text{Uniform}(0, 1)$
- If  $MUg(\theta_j) < f(\mathbf{Y} \mid \theta_j)\pi(\theta_j)$ , accept  $\theta_j$ ; otherwise reject  $\theta_j$ .
- Return to step (i) and repeat until a desired size is obtained. The members of this sample is an iid sample from  $\pi(\theta \mid \mathbf{Y})$

Figure 3: Graphical representation of the rejection sampling method



The third condition ensures that the density of the sampled  $\theta$  is within  $[A, B]$ . Therefore, more likely to be from the posterior distribution

##### 3. Weighted bootstrap

Again consider the posterior density. Suppose that we have a sample  $\theta_1, \dots, \theta_N$  from approximating

density  $g(\theta)$ .

Define

$$w_i = \frac{f(\mathbf{Y} | \theta_i) \pi(\theta_i)}{g(\theta_i)}, \quad q_i = \frac{w_i}{\sum_{j=1}^N w_j}$$

Draw  $\theta^*$  from the discrete distribution over  $\{\theta_1, \dots, \theta_N\}$  which places mass  $q_i$  at  $\theta_i$ .

In these indirect sampling methods, the prior  $\pi(\theta)$ , if proper, can play a role.

*Example 1:* Suppose  $Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} N(\theta, \sigma^2)$  and  $\pi(\theta) = \text{Cauchy}(\mu, \tau)$  with known  $\sigma^2, \mu$ , and  $\tau$ . Likelihood  $f(\mathbf{Y} | \theta)$  is maximized at  $\hat{\theta} = \bar{Y} \implies M = f(\mathbf{Y} | \hat{\theta})$  in the rejection method, and  $g(\theta) = \pi(\theta)$ . However,  $\pi(\theta)$  is often very flat relative to  $f(\mathbf{Y} | \theta) \implies$  quite inefficient.

## 4.9 MCMC methods

### 4.9.1 Substitution sampling

Consider the 3-stage hierarchical model

$$\begin{aligned} &\text{likelihood } p(Y | \theta) \\ &\text{prior } p(\theta | \eta) \\ &\text{hyperprior } p(\eta) \end{aligned}$$

Assuming the prior is conjugate with the likelihood  $\implies$  marginal distribution  $p(Y | \eta) = \int p(Y | \theta) p(\theta | \eta) d\theta$  easily computed  $\implies$  closed-form posterior  $p(\theta | Y, \eta) = \frac{p(Y|\theta)p(\theta|\eta)}{p(Y|\eta)}$

Also assuming the hyperprior is conjugate with the prior  $\implies$  closed-form  $p(\eta | \theta)$  For inference, we seek the marginal posterior,

$$p(\theta | Y) = \int p(\theta | Y, \eta) p(\eta | Y) d\eta$$

and we also have

$$p(\eta | Y) = \int p(\eta | \theta) p(\theta | Y) d\theta$$

Switching to generic notation, we have a system of two linear integral equations

$$\begin{aligned} p(x) &= \int p(x | y) p(y) dy \\ p(y) &= \int p(y | x) p(x) dx \end{aligned}$$

where  $p(x | y)$  and  $p(y | x)$  are known, and we seek  $p(x)$ . This is a fixed point system:

$$\begin{aligned} p(x) &= \int p(x | y) \int p(y | x') p(x') dx' dy \\ &= \int h(x, x') p(x') dx' \end{aligned}$$

where  $h(x, x') = \int p(x | y) p(y | x') dy$ . *Sampling-based algorithm:*

Draw  $X^{(0)} \sim p_0(x)$

Draw  $Y^{(1)} \sim p(y | X^{(0)}) \sim p_1(y) = \int p(y | x) p_0(x) dx$

Draw  $X^{(1)} \sim p(x | Y^{(1)}) \sim p_1(x) = \int h(x, x') p_0(x') dx'$

Repeat this process:

$$X^{(i)} \stackrel{d}{\sim} X \sim p(x) \quad Y^{(i)} \stackrel{d}{\sim} Y \sim p(y)$$

*A variant:* multiple  $Y_1^{(i)}, \dots, Y_m^{(i)} \stackrel{\text{iid}}{\sim} p(y | X^{(i-1)})$  and single

$$X^{(i)} \sim \hat{p}_i(x) = \frac{1}{m} \sum_{j=1}^m p(x | Y_j^{(i)})$$

$\implies$  automatically produce a smooth estimate of  $p(x)$

*Parallel sampling with multiple chains:* Marginally independent replicates  $\implies$  presumably better  $\hat{p}_i(x)$ . But wasteful.

*Ergodic sampling with a single chain:* Continue for an additional  $m - 1$  iterations after convergence at iteration  $i$ .

$$\hat{p}_i(x) = \frac{1}{m} \sum_{j=1}^m p(x | Y_j^{(i+j-1)})$$

To reduce high autocorrelation, retain only every  $k$  th iteration:

$$\hat{p}_i(x) = \frac{1}{m} \sum_{j=1}^m p\left(x \mid Y_j^{(i+(j-1)k)}\right)$$

Although  $X$  and  $Y$  can conceivably be vectors, sampling from complex multivariate distributions is difficult  $\implies$  Require a  $K$ -variate extension. Consider  $K = 3$  :

$$p(x) = \int p(x, z \mid y)p(y)dydz$$

$$p(y) = \int p(y, x \mid z)p(z)dx dz$$

$$p(z) = \int p(z, y \mid x)p(x)dx dy$$

Closed-form bivariate conditional distributions are unlikely available. With univariate distributions only:

$$p(x) = \int p(x \mid z, y)p(z \mid y)p(y)dy dz$$

$$p(y) = \int p(y \mid x, z)p(x \mid z)p(z)dx dz$$

$$p(z) = \int p(z \mid y, x)p(y \mid x)p(x)dx dy$$

6 of them!

General  $K$ -dimensional problem  $\implies K(K - 1)$  distributions  $\implies$  impractical for large  $K$ .

#### 4.9.2 Gibbs sampling

Gibbs sampler can be viewed as a special case of Metropolis-Hastings. And Metropolis-Hastings algorithm is a generalization of the Metropolis algorithm.

Under mild conditions, between full or complete conditional distributions and joint distribution

$$\{p_i(U_i \mid U_{j \neq i}), i = 1, \dots, K\} \iff p(U_1, \dots, U_K)$$

Suppose these full conditional distributions are available for sampling. Given an arbitrary set of starting values  $(U_1^{(0)}, \dots, U_K^{(0)})$ , the algorithm proceeds as follows:

$$\text{Draw } U_1^{(1)} \sim p_1(U_1 \mid U_2^{(0)}, \dots, U_K^{(0)}),$$

$$\text{Draw } U_2^{(1)} \sim p_2(U_2 \mid U_1^{(1)}, U_3^{(0)}, \dots, U_K^{(0)}),$$

$$\text{Draw } U_K^{(1)} \sim p_K(U_K \mid U_1^{(1)}, \dots, U_{K-1}^{(1)}),$$

completing one iteration of the scheme. After  $t$  such iterations, we obtain  $(U_1^{(t)}, \dots, U_K^{(t)})$ .

$$(U_1^{(t)}, \dots, U_K^{(t)}) \xrightarrow{d} (U_1, \dots, U_K) \sim p(U_1, \dots, U_K)$$

Hierarchical models with conjugate priors and hyperpriors  $\implies$  closed form  $\{p_i(U_i \mid U_{j \neq i}, \mathbf{Y}), i = 1, \dots, K\}$   $\implies$  joint posterior  $p(U_1, \dots, U_K \mid \mathbf{Y})$  by Gibbs sampler.

Marginal posterior  $p(U_i \mid \mathbf{Y})$  can be obtained the same manner as before. For example, with parallel sampling,

$$\hat{p}_t(U_i \mid \mathbf{Y}) = \frac{1}{m} \sum_{j=1}^m p\left(U_i \mid U_{1,j}^{(t)}, \dots, U_{i-1,j}^{(t)}, U_{i+1,j}^{(t)}, \dots, U_{K,j}^{(t)}, \mathbf{Y}\right)$$

which is less variable than a kernel-smoothed estimate.

What about non-conjugate priors? One may consider, for example, an indirect sampling method. But such a solution is not ideal.

### 4.9.3 Metropolis-Hastings algorithm

Target: joint posterior  $p(u)$  for (possibly vector-valued)  $U$

Candidate or proposal density:  $q(v, u)$  such that  $q(\cdot, u)$  is a pdf  $\forall u$ , and  $q(u, v) = q(v, u) \forall u, v$ .

The Metropolis algorithm:

1. Draw  $v \sim q(\cdot, U^{(t-1)})$
2. Compute the density ratio  $r = p(v)/p(u)$
3. If  $r \geq 1$ , set  $U^{(t)} = v$

If  $r < 1$ , set  $U^{(t)} = \begin{cases} v & \text{with prob } r \\ u & \text{with prob } 1 - r \end{cases}$

Mild conditions  $\implies U^{(t)} \xrightarrow{d} U \sim p(\cdot)$ .

$p(\cdot)$  is needed only up to proportionality constant:

$$\pi(\boldsymbol{\theta} \mid \mathbf{Y}) \propto f(\mathbf{Y} \mid \boldsymbol{\theta})\pi(\boldsymbol{\theta})$$

in Bayesian applications.

Why symmetric  $q$  ?

Consider finite state spaces, where the transition kernel can be represented by matrix  $\mathbf{P}$ ;  $P_{ij}$  is the prob of moving from state  $i$  to  $j$ .

The Markov chain has equilibrium distribution

$$\mathbf{d} = (d_1, \dots, d_k)^T$$

if and only if

$$\mathbf{d}^T \mathbf{P} = \mathbf{d}^T.$$

Symmetric  $q \implies$

$$\begin{aligned} d_i P_{ij} &= d_i \left[ \min \left( 1, \frac{d_j}{d_i} \right) Q_{ij} \right] = \min(d_i, d_j) Q_{ij} \\ &= \min(d_i, d_j) Q_{ji} = d_j P_{ji} \end{aligned}$$

$\implies$  chain is reversible  $\implies$

$$(\mathbf{d}^T \mathbf{P})_j = \sum_{i=1}^k d_i P_{ij} = \sum_{i=1}^k d_j P_{ji} = d_j \sum_{i=1}^k P_{ji} = d_j.$$

For the continuous parameter settings, a convenient choice for  $q$  is a  $N(\boldsymbol{\theta}^{(i-1)}, \tilde{\Sigma})$ . In theory, any positive-definite  $\tilde{\Sigma}$  suffices. Care is required in practice:

- too large  $\implies$  large jumps  $\implies$  many candidates far from posterior support and rejected  $\implies$  tendency to "get stuck"
- too small  $\implies$  "baby-stepping"

A simple but important generalization of the Metropolis algorithm was due to Hastings, by dropping the symmetry requirement of  $q$  and redefining

$$r = \frac{p(v)q(u, v)}{p(u)q(v, u)}$$

#### From Yale BIS 567

Metropolis Algorithm

Define  $h(\boldsymbol{\theta}) \equiv f(y \mid \boldsymbol{\theta})f(\boldsymbol{\theta})$  and  $q(\boldsymbol{\theta}^* \mid \boldsymbol{\theta}^{(t-1)})$  is a candidate or proposal density. It must satisfy

1. valid density function for every possible value of  $\boldsymbol{\theta}^{(t-1)}$
2. be symmetric  $q(\boldsymbol{\theta}^* \mid \boldsymbol{\theta}^{(t-1)}) = q(\boldsymbol{\theta}^{(t-1)} \mid \boldsymbol{\theta}^*)$ . For example
  - $q(\boldsymbol{\theta}^* \mid \boldsymbol{\theta}^{(t-1)}) = \text{Uniform}(\boldsymbol{\theta}^{(t-1)} - \delta, \boldsymbol{\theta}^{(t-1)} + \delta)$
  - $q(\boldsymbol{\theta}^* \mid \boldsymbol{\theta}^{(t-1)}) = \text{Normal}(\boldsymbol{\theta}^{(t-1)}, \delta^2)$

$\delta$  is chosen to make the algorithm run more efficiently. A poor choice of  $\delta$  leads to high correlation in the Markov chain.

- (a)  $\delta$  too big: Low acceptance of a new  $\theta$  and high correlation between samples
- (b)  $\delta$  too small: High acceptance of a new  $\theta$  and high correlation between samples

Generally, acceptance around 20% to 50% is desired. Pilot runs are often required to properly tune the algorithm and select a reasonable  $\delta$ .

The algorithm: for  $t = 1, \dots, T$  :

1. draw  $\theta^*$  from  $q(\cdot | \theta^{(t-1)})$
2. compute  $r = h(\theta^*) / h(\theta^{(t-1)}) = \exp[\log\{h(\theta^*)\} - \log\{h(\theta^{(t-1)})\}]$
3. if  $r \geq 1$ , set  $\theta^{(t)} = \theta^*$ ; if  $r < 1$ , set

$$\theta^{(t)} = \begin{cases} \theta^* & \text{with probability } r \\ \theta^{(t-1)} & \text{with probability } 1 - r \end{cases}$$

Under mild conditions (similar to those for the Gibbs sampler), a draw  $\theta^{(t)}$  converges in distribution to a draw from the true posterior density  $f(\theta | \mathbf{Y} = \mathbf{y})$ .

*Intuition:*

Suppose we have a working set of posterior samples  $\{\theta^{(1)}, \dots, \theta^{(s)}\}$ . We want to add a new sample to the set,  $\theta^*$ , which is nearby  $\theta^{(t)}$ . If  $f(\theta^* | \mathbf{Y} = \mathbf{y}) > f(\theta^{(s)} | \mathbf{Y} = \mathbf{y})$ , then we want to add  $\theta^*$  to the set. If  $f(\theta^* | \mathbf{Y} = \mathbf{y}) < f(\theta^{(s)} | \mathbf{Y} = \mathbf{y})$ , then we don't necessarily want to add it to the set. So we accept it with probability  $r$

$$r = \frac{f(\theta^* | \mathbf{Y} = \mathbf{y})}{f(\theta^{(s)} | \mathbf{Y} = \mathbf{y})} = \frac{f(\mathbf{y} | \theta^*) f(\theta^*)}{f(\mathbf{y})} \frac{f(\mathbf{y})}{f(\mathbf{y} | \theta^{(s)}) f(\theta^{(s)})} = \frac{f(\mathbf{y} | \theta^*) f(\theta^*)}{f(\mathbf{y} | \theta^{(s)}) f(\theta^{(s)})}$$

#### 4.9.4 Hybrid forms

Several MCMC algorithms may be combined in a single problem, to take advantage of the strengths of each. Markov kernels  $P_1, \dots, P_m$  all have stationary distribution  $p$ ; they may correspond to, say, one Gibbs sampler and  $m - 1$  Metropolis algorithms.

*Mixer:* At each iteration, kernel  $P_i$  is chosen with prob  $\alpha_i$ , where  $\sum_i \alpha_i = 1$ .

*Cycle:* Each kernel  $P_i$  is used in a prespecified order.

*Metropolis within Gibbs:* In the case that, say, all full conditional distributions are available in closed form except for one, a Metropolis subalgorithm may be embedded in the Gibbs sampler.

## 5 Large sample theory

## 6 M-Estimation (Estimating Equations)

### 6.1 Introduction

$$\sum_{i=1}^n \Psi(Y_i; \theta) = 0$$

$Y_1, \dots, Y_n$  : independent but not necessarily identically distributed

$\theta$  :  $b \times 1$  parameter

$\Psi$  :  $b \times 1$  function

**Example 1.** the mean estimator  $\hat{\theta} = n^{-1} \sum_{i=1}^n Y_i$  is from

$$\sum_{i=1}^n (Y_i - \theta) = 0$$

the deviance from the mean  $\hat{\theta}_1 = n^{-1} \sum_{i=1}^n |Y_i - \bar{Y}|$  is from

$$\sum_{i=1}^n \begin{pmatrix} |Y_i - \theta_2| - \theta_1 \\ Y_i - \theta_2 \end{pmatrix} = \mathbf{0}$$

## 6.2 Basic approach

Suppose  $Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} F$  and the true parameter  $\theta_0$  is defined as a solution to

$$E_F \Psi(Y_1, \theta) = \int \Psi(y, \theta) dF(y) = \mathbf{0}$$

if  $\Psi$  is suitably smooth then by Taylor expansion of

$$\mathbf{G}_n(\theta) = n^{-1} \sum_{i=1}^n \Psi(Y_i, \theta)$$

we have

$$\begin{aligned} 0 &= \mathbf{G}_n(\hat{\theta}) = \mathbf{G}_n(\theta_0) + \mathbf{G}'_n(\theta_0) (\hat{\theta} - \theta_0) + o_p(n^{-1/2}) \\ \sqrt{n} (\hat{\theta} - \theta_0) &= \{-\mathbf{G}'_n(\theta_0)\}^{-1} \sqrt{n} \mathbf{G}_n(\theta_0) + o_p(1) \end{aligned}$$

by WLLN

$$-\mathbf{G}'_n(\theta_0) = -n^{-1} \sum_{i=1}^n \Psi'(Y_i, \theta_0) \xrightarrow{P} E_F \{-\Psi'(Y_1, \theta_0)\} \equiv \mathbf{A}(\theta_0)$$

by CLT

$$\sqrt{n} \mathbf{G}_n(\theta_0) \xrightarrow{d} N(\mathbf{0}, \mathbf{B}(\theta_0)) \text{ where } \mathbf{B}(\theta_0) = E_F \{\Psi(Y_1, \theta_0)^{\otimes 2}\}$$

then by Slutsky's theorem

$$\hat{\theta} \sim AN(\theta_0, \mathbf{V}(\theta_0)/n) \text{ where } \mathbf{V}(\theta_0) = \mathbf{A}(\theta_0)^{-1} \mathbf{B}(\theta_0) \{\mathbf{A}(\theta_0)^{-1}\}^T$$

More generally, the estimating equation does not have to be exactly 0, for example

$$\sum_{i=1}^n \Psi(Y_i, \theta) = o_p(n^{1/2}) \Rightarrow \frac{1}{n} \sum_{i=1}^n \Psi(Y_i, \theta) = o_p(n^{-1/2})$$

The asymptotic distribution of  $\hat{\theta}$  remains the same.

### 6.2.1 Estimation for $\mathbf{A}$ , $\mathbf{B}$ , and $\mathbf{V}$

Empirical estimator of  $\mathbf{A}(\theta_0)$  :

$$\mathbf{A}_n(\mathbf{Y}, \hat{\theta}) = -n^{-1} \sum_{i=1}^n \Psi'(Y_i, \hat{\theta})$$

Empirical estimator of  $\mathbf{B}(\theta_0)$  :

$$\mathbf{B}_n(\mathbf{Y}, \hat{\theta}) = n^{-1} \sum_{i=1}^n \Psi(Y_i, \hat{\theta})^{\otimes 2}$$

then

$$\mathbf{V}_n(\mathbf{Y}, \hat{\theta}) = \mathbf{A}_n(\mathbf{Y}, \hat{\theta})^{-1} \mathbf{B}_n(\mathbf{Y}, \hat{\theta}) \{\mathbf{A}_n(\mathbf{Y}, \hat{\theta})^{-1}\}^T$$

A special case is the MLE, where

$$\mathbf{A}(\theta_0) = \mathbf{B}(\theta_0) = \mathbf{I}(\theta_0) \quad \mathbf{V}(\theta_0) = \mathbf{I}(\theta_0)^{-1}$$

**Example 2.** Sample mean and variance

$$\begin{aligned}\Psi(Y_i, \boldsymbol{\theta}) &= \begin{pmatrix} Y_i - \theta_1 \\ (Y_i - \theta_1)^2 - \theta_2 \end{pmatrix} \\ \mathbf{A}(\boldsymbol{\theta}_0) &= E\{-\Psi'(Y_1, \boldsymbol{\theta}_0)\} = E\begin{pmatrix} 1 & 0 \\ 2(Y_1 - \theta_{10}) & 1 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \\ \mathbf{B}(\boldsymbol{\theta}_0) &= E\{\Psi(Y_1, \boldsymbol{\theta}_0)^{\otimes 2}\} = \begin{pmatrix} \sigma^2 & \mu_3 \\ \mu_3 & \mu_4 - \sigma^4 \end{pmatrix} \\ \mathbf{V}(\boldsymbol{\theta}_0) &= \mathbf{B}(\boldsymbol{\theta}_0) \\ \mathbf{B}_n(\mathbf{Y}, \hat{\boldsymbol{\theta}}) &= n^{-1} \sum_{i=1}^n \begin{pmatrix} (Y_i - \bar{Y})^2 & (Y_i - \bar{Y})^3 \\ (Y_i - \bar{Y})^3 & \{(Y_i - \bar{Y})^2 - s_n^2\}^2 \end{pmatrix}\end{aligned}$$

the estimated  $\hat{\boldsymbol{\theta}} = (\bar{Y}, s_n^2)^T$  is the MLE for the normal density  $f(y; \boldsymbol{\theta}) = \frac{1}{\sqrt{2\pi\theta_2}} \exp\left\{-\frac{(y-\theta_1)^2}{2\theta_2}\right\}$

**Example 3.** Ratio estimator

For i.i.d  $(Y_1, X_1), \dots, (Y_n, X_n)$  and  $\mu_X \neq 0$

$$\begin{aligned}\Psi(Y_i, X_i, \theta) &= Y_i - \theta X_i \\ A(\theta_0) &= E(X_1) = \mu_X \\ B(\theta_0) &= E\{(Y_1 - \theta_0 X_1)^2\} \\ V(\theta_0) &= E\{(Y_1 - \theta_0 X_1)^2\} / \mu_X^2 \\ A_n(\mathbf{Y}, \hat{\boldsymbol{\theta}}) &= \bar{X} \\ B_n(\mathbf{Y}, \hat{\boldsymbol{\theta}}) &= \frac{1}{n} \sum_{i=1}^n \left(Y_i - \frac{\bar{Y}}{\bar{X}} X_i\right)^2 \\ V_n(\mathbf{Y}, \hat{\boldsymbol{\theta}}) &= \frac{1}{\bar{X}^2} \frac{1}{n} \sum_{i=1}^n \left(Y_i - \frac{\bar{Y}}{\bar{X}} X_i\right)^2\end{aligned}$$

The  $\Psi$  functions are not unique, it can even be of different dimensions

$$\begin{aligned}\Psi(Y_i, X_i, \boldsymbol{\theta}) &= \begin{pmatrix} Y_i - \theta_1 \\ X_i - \theta_2 \\ \theta_1 - \theta_3 \theta_2 \end{pmatrix} \\ \mathbf{A}(\boldsymbol{\theta}_0) &= \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ -1 & \theta_{30} & \theta_{20} \end{pmatrix} \\ \mathbf{B}(\boldsymbol{\theta}_0) &= \begin{pmatrix} \sigma_Y^2 & \sigma_{YX} & 0 \\ \sigma_{YX} & \sigma_X^2 & 0 \\ 0 & 0 & 0 \end{pmatrix} \quad \text{singular} \\ \mathbf{V}(\boldsymbol{\theta}_0) &= \mathbf{A}(\boldsymbol{\theta}_0)^{-1} \mathbf{B}(\boldsymbol{\theta}_0) \left\{ \mathbf{A}(\boldsymbol{\theta}_0)^{-1} \right\}^T \\ v_{33} &= \frac{1}{\theta_{20}^2} (\sigma_Y^2 - 2\theta_{30}\sigma_{YX} + \theta_{30}^2\sigma_X^2)\end{aligned}$$

### 6.3 Delta method via M-estimation