

## 1. Key Distributions

- Binomial :  $f(x) = \binom{n}{x} p^x (1-p)^{n-x}$   $\frac{nx}{\theta x} \frac{\theta^2}{2}$
  - Poisson :  $f(x) = \frac{x^k}{x!} \exp(-\lambda)$ ,  $\lambda > 0$
  - Exponential :  $f(x) = \lambda \cdot \exp(-\lambda x)$   $1/\lambda$   $1/\lambda^2$
- B) Exponential family

$$f(y|x) = h(y) \exp(\eta(\theta) T(y) - A(\eta(\theta)))$$

canonical form:  $\eta(\theta) = C_1 \theta + C_2$

## 2. Method of moments

- Set theoretical moment equal to sample moments

$$E[X^k] \text{ set equal to } \frac{1}{n} \sum x_i^k$$

$\rightarrow$  can also at moment around center

### Expected value

cont.  $E[X^k] = \int_{-\infty}^{\infty} x^k \cdot f(x) dx$

disc.  $E[X^k] = \sum_i x_i^k p(x)$

-  $\text{Var}[X] = E[(X - E[X])^2] = E[X^2] - E[X]^2$

- Corr:  $\rho = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X) \text{Var}(Y)}}$  where  $\text{Cov}(X, Y) = E[XY] - E[X]E[Y]$

## 3. a) Likelihood + log-Likelihood

$$L(\theta) = f(x, \theta) = \prod_{i=1}^n f(x_i; \theta), \theta \in \Theta \quad \text{common}$$

$$\ell(\theta) = \log(L(\theta)) = \sum_{i=1}^n \log(f(x_i; \theta)), \theta \in \Theta$$

$\rightarrow$  (log-) likelihood without constant as known as kernel

## b) MLE

$$\hat{\theta}_{ML} = \underset{\theta}{\operatorname{argmax}} L(\theta) = \underset{\theta}{\operatorname{argmax}} \ell(\theta) = \underset{\theta}{\operatorname{argmax}} \tilde{L}(\theta) = \underset{\theta}{\operatorname{argmax}} \tilde{\ell}(\theta)$$

c) Relative Likelihood  $\tilde{L}(\theta) = \frac{L(\theta)}{L(\hat{\theta}_{ML})}$

$$\tilde{\ell}(\theta) = \ell(\theta) - \ell(\hat{\theta}_{ML})$$

## d) One-to-one Transformations

- $\tilde{\ell}_{\varphi}(\theta) = h(\hat{\theta}_{ML})$  If  $\varphi = h(\theta)$

$$\cdot S_{\varphi}(\varphi) = S_{\theta}(\theta) \frac{d h^{-1}(\varphi)}{d \varphi}$$

$$\cdot I_{\varphi}(\varphi) = I_{\theta}(\theta) \left[ \frac{d h^{-1}(\varphi)}{d \varphi} \right]^2 - S_{\theta}(\theta) \frac{d^2 h^{-1}(\varphi)}{d \varphi^2}$$

$$\cdot I_{\varphi}(\hat{\theta}_{ML}) = I_{\theta}(\hat{\theta}_{ML}) \left[ \frac{d h(\hat{\theta}_{ML})}{d \theta} \right]^{-2} \text{ diff. first! then plug in}$$

$$\cdot J_{\varphi}(\theta) = J_{\theta}(\theta) \cdot \left[ \frac{d h(\theta)}{d \theta} \right]^{-2}$$

## e) Quadratic approx. of log-Likelihood

$$\ell(\theta) \approx \ell(\hat{\theta}_{ML}) + \underbrace{S_{\theta}(\hat{\theta}_{ML})(\theta - \hat{\theta}_{ML})}_{\geq 0} - \frac{1}{2} I_{\theta}(\hat{\theta}_{ML})(\theta - \hat{\theta}_{ML})^2$$

$\rightarrow$  simplification using rel. log-Likelihood

$$\tilde{\ell}(\theta) \approx -\frac{1}{2} I_{\theta}(\hat{\theta}_{ML})(\theta - \hat{\theta}_{ML})^2$$

## f) Likelihood Ratio

$$\lambda_{x_{1:n}}(\theta_1, \theta_2) = \frac{L(\theta_1; x_{1:n})}{L(\theta_2; x_{1:n})} = \frac{\tilde{L}(\theta_1, x_{1:n})}{\tilde{L}(\theta_2, x_{1:n})}$$

$\hookrightarrow$  cf compared to  $\tilde{L}(\hat{\theta}_{ML}; x_{1:n}) \rightarrow$  can recover  $\tilde{L}(\theta)$

## 4. Sufficient Statistics

$T$  with realization  $t = h(x_{1:n})$  is sufficient for  $\theta$  iff:

$$f(x_{1:n}|t) = \frac{f(t|x_{1:n}) \cdot f(x_{1:n})}{f(t)} \perp \theta$$

## b) Factorization Theorem

$$f(x_{1:n}; \theta) = g_1(t, \theta) \cdot g_2(x_{1:n}) \quad \begin{matrix} \leftarrow \text{can also} \\ \text{e.g. sample} \end{matrix}$$

## c) Minimal Sufficiency

•  $T$  is min. sufficient if it can be written as a function of any other sufficient statistic

$\rightarrow T_1, T_2$  min. sufficient for  $\theta$

$$\rightarrow T_1 = g(T_2), \quad T_2 = g^{-1}(T_1)$$

• Nec. and suff. criterion for  $T$  to be min. suff. is that  $h(x_{1:n}) \stackrel{d}{=} h(\tilde{x}_{1:n})$  iff:

$$\lambda_{x_{1:n}}(\theta_1, \theta_2) = \lambda_{\tilde{x}_{1:n}}(\theta_1, \theta_2) \rightarrow L(\theta) \text{ as min. suff.}$$

## 5. Unbiasedness and Consistency

• Estimator ( $\hat{T}_n = h(X_{1:n})$ ) vs. Estimate ( $\hat{t}_n = h(x_{1:n})$ )  
 $\hat{T}$  suitable statistic whose realization is a point estimate  $\hat{t}$  suitable value for  $\theta$  realization acc. to data and model

• Unbiasedness:  $E(\hat{T}_n(\hat{\theta})) = \theta, \forall \theta \in \Theta, n \in \mathbb{N}$  (asymptotic)

• Consistency:  $\hat{T}_n(\hat{\theta}) \rightarrow \theta$  for  $n \rightarrow \infty$

• MS Consistency:  $E[(\hat{T}_n(\hat{\theta}) - \theta)^2] \xrightarrow{n \rightarrow \infty} 0$  (unbiased  $\hat{\theta}$  consistent)

• Jensen's Inequality:  $E[g(x)] \geq g(E[x])$  linear  $\geq$  convex  $g$

• CLT: To derive asympt. distributions for mean ( $\bar{x}$ ) or sum ( $\sum x_i$ ) of R.V.s

$$\frac{\sqrt{n}(\bar{x} - E[\bar{x}])}{\sqrt{\text{Var}[\bar{x}]}} \xrightarrow{D} N(0, 1)$$

## 6. Further Likelihood properties / F.I.

a) E and Var of Score (under FRC)  
 $E[S(\theta; X)] = 0$   
 $\text{Var}[S(\theta; X)] = J(\theta)$

b) Exp. F.I. from a random sample:  
 $J_{\text{true}}(\theta) = n \cdot J(\theta)$

## 7. Test and Confidence Interval

### a) Confidence Interval:

$$P(T_L \leq \theta \leq T_U) = \gamma \quad \forall \theta \in \Theta$$

b) Pivof: statistic depending on data ( $x_{1:n}$ ) and  $\theta$  whose distribution is independent of  $\theta$   $\rightarrow$  pivotal distribution  
 $\hookrightarrow$  approx. Pivof: first show asymptotic MS consistency!

### c) Tests: quantity evidence against $H_0$ (using p-value)

- P-value:  $P(\text{obtaining equal or more extreme result } T \text{ than observed } | H_0 \text{ is true})$

- Type I:  $P(T(x_{1:n}) \geq c_{\text{ref}} | H_0) = \alpha$  (sign. level)

- Type II:  $P(T(x_{1:n}) \leq c_{\text{ref}} | H_1) = \beta \rightarrow 1 - \beta$  (power)

$\rightarrow$  Power curve: shows  $P_{1-\beta}$  (power) for diff. levels of  $\theta$

### d) Statistics (Z, Wald, Score, LR)

• Wald:  $\sqrt{I(\hat{\theta}_{ML})} (\hat{\theta}_{ML} - \theta_0) \stackrel{\text{or}}{\sim} N(0, 1)$   
 $\text{or } J(\hat{\theta}_{ML}) \text{ or } \text{or } I/J(\hat{\theta})$

① CI:  $[ \hat{\theta}_{ML} \pm z_{1-\alpha/2} \cdot \text{se}(\hat{\theta}_{ML}) ]$

② CI:  $\{ \theta : \frac{(\hat{\theta}_{ML} - \theta)^2}{I(\theta)} \leq z_{1-\alpha/2}^2 \} \rightarrow \text{solve for } \theta$

$\hookrightarrow$  solve quad. eq:  $x = -b \pm \sqrt{b^2 - 4ac}$

• Score  $S(\theta; x_{1:n}) \stackrel{\text{a}}{\sim} N(0, 1)$  (large. consistent)  
 $\text{FRC}$  hold  $\frac{S(\theta; x_{1:n})}{\sqrt{J(\theta; x_{1:n})}} \stackrel{\text{a}}{\sim} N(0, 1)$   
 $\text{or } I(\theta) \text{ or } I/J(\hat{\theta}_{ML})$

$\rightarrow$  Same CI as ②

### e) Likelihood Ratio

$$W = 2 \cdot \log \left[ \frac{L(\hat{\theta}_{ML})}{L(\theta)} \right] = -2 \tilde{\ell}(\theta) \stackrel{\text{a}}{\sim} \chi^2_{p,p}$$

$$\rightarrow \text{CI}: \{ \theta : \tilde{L}(\theta) \geq \exp \left[ -\frac{1}{2} \chi^2_{p,p} \right] \}$$

$$= \tilde{\ell}(\theta) \geq -\frac{1}{2} \chi^2_{p,p}$$

$\hookrightarrow$  use LR to compute CI / Region for test when we have rel. log-Likelihood

e) Common critical values with Z-score:  
 10%: 1.64; 5%: 1.96; 1%: 2.58

8. Distribution of the MLE  $\rightarrow$  Need FRC

a)  $\hat{\theta}_{ML} \sim N(\theta_0, J_{\theta\theta}(\theta_0)^{-1})$  as  $n \uparrow$ ,  
 or  $n \cdot J(\theta_0)^{-1}, I(\theta_0)^{-1}$

$\rightarrow$  MLE is asympt. unbiased + consistent

b) Cramér-Rao Lower bound  $\rightarrow$  shows efficiency  
 for an arbitrary consistent estimator, under FRC  
 $\text{Var}(T(x_{1:n})) \geq \frac{g'(\theta)^2}{J(\theta)} \left( = \frac{1}{J(\theta)} \text{ if } g(\theta) = \theta \right)$   
 b/c under CRLB for  $\text{Var}$

$\rightarrow$  an asympt. unbiased estimator that attains the CRLB is efficient

b) Efficiency:  $\text{eff}(T) = \frac{\partial}{\partial \theta} E(T(x_{1:n}))^2 \leq 1$

c) SE of the MLE

$$\text{se}(\hat{\theta}_{ML}) = 1/\sqrt{J(\hat{\theta}_{ML})} \text{ or } 1/\sqrt{J(\hat{\theta}_{ML}) \cdot \text{Var}(T(x_{1:n}))}$$

d) Delta Method (for SE in general)

$$\text{se}(h(\hat{\theta})) = \text{se}(\hat{\theta}) \cdot \left| \frac{d h(\hat{\theta})}{d \hat{\theta}} \right|$$

$\hookrightarrow$  can use to calc Wald CI for transformation

e) Variance Stabilizing Transformation

$$\phi = h(\theta) \propto \int^{\theta} J_{\theta}(u)^{\frac{1}{2}} du$$

$\phi$  is a var. stabilizing trans of  $\theta$

9. Multiparameter Model

a)  $L(\vec{\theta}) \rightarrow l(\vec{\theta}) \rightarrow S(\vec{\theta}) = \frac{\partial l(\vec{\theta})}{\partial \vec{\theta}} \dots \hat{\theta}_{ML} = (\hat{\theta}_1 \dots \hat{\theta}_p)^T$

b) Fisher Info Matrix:  
 $I(\vec{\theta}) = - \left( \frac{\partial^2 l(\vec{\theta})}{\partial \theta_i \partial \theta_j} \right), 1 \leq i, j \leq p$  sym.  $p \times p$  matrix  
 + pos. definite  
 $\hookrightarrow$  exp. of the Hessian

c) SE:  $\sqrt{[I(\hat{\theta}_{ML})]_{ii}}$  (ith diagonal component)

d) SE of a difference: split data into two independent parts (e.g. obs. 0 and 1 with two LCs)  
 $\rightarrow$  SE of  $\hat{\delta}_{ML} = \hat{\theta}_{ML} - \hat{\theta}_{ML}$  has form:

$$\text{se}(\hat{\delta}_{ML}) = \sqrt{J_{\theta\theta}(\hat{\theta}_{ML})^{-1} + I_{\theta\theta}(\hat{\theta}_{ML})^{-1}}$$

e) Score Statistic:  $J_{\theta\theta}(\hat{\theta})^{-\frac{1}{2}} S(\hat{\theta}|x_{1:n})^T \sim N_p(0, I_p)$   
 or  $I(\hat{\theta}|x_{1:n}) \text{ or } \hat{\theta}_{ML}$

f) Multivariate Delta Method: Transformation  $g(\hat{\theta}_{ML})$ ,  
 for  $\hat{\theta}_{ML} \sim N_p(\hat{\theta}, I(\hat{\theta}_{ML})^{-1})$   $R^p \rightarrow R^q$

$$\rightarrow g(\hat{\theta}_{ML}) \sim N_q(g(\hat{\theta}), D(\hat{\theta}_{ML}) I(\hat{\theta}_{ML})^{-1} D(\hat{\theta}_{ML})^T)$$

where  $D(\theta)$  denotes  $q \times p$  Jacobian (1st deriv) of  $g(\theta)$

$\rightarrow$  sqrt of Var term are sum of components of  $g(\hat{\theta}_{ML})$

10) Profile Likelihood

a)  $L_p(\vec{\theta}) = \max_{\theta} L(\vec{\theta}; \vec{x}) = L(\vec{\theta}; \hat{\vec{\theta}}_{ML}(\vec{\theta}))$   $\hookrightarrow$  max thp multivariate and

b) Estimated Likelihood of  $\theta$ :  $L_e(\vec{\theta}) = L(\vec{\theta}; \hat{\vec{\theta}}_{ML})$

$\rightarrow$  ignores uncertainty on  $\vec{\theta}$ ! also generally not a really quantifies the uncertainty w.r.t.  $\vec{\theta}$

c) + Profile log-Likelihood & relative profile(log-)Likelihood

11) Generalized Likelihood Ratio Statistic

- LR stat analogous to above  $\chi^2 \sim \chi^2_p$
- For  $\vec{\theta}$  ( $q$ -dim) parameter of interest vector  
 $\vec{\eta}$  ( $r$ -dim) nuisance parameter vector

1. get joint likelihood  $L(\vec{\theta}; \vec{x})$

2. profile likelihood for parameters of interest  $L_p(\vec{\theta}) = \max_{\vec{\eta}} L(\vec{\theta}; \vec{\eta})$

3. LR:  $W_p := -2 \ln \frac{L_p(\vec{\theta}_{ML})}{L_p(\vec{\theta})}$  (via profile log-.)  
 $= 2 \ln \left[ \frac{L_p(\vec{\theta}_{ML})}{L_p(\vec{\theta})} \right] \sim \chi^2_q$

$\hookrightarrow$  can use to test  $H_0: \vec{\theta} = \vec{\theta}_0 \rightarrow$  reject if  $W_p$  is large ( $> \chi^2_{q, 1-\alpha}$ )

12) Regression:  $\vec{y} = X \vec{\beta} + \vec{\epsilon}$  (with  $X$  full rank)

- $\hat{\beta}_{LS} = \arg \min_{\vec{\beta}} (\vec{y} - X\vec{\beta})^T (\vec{y} - X\vec{\beta})^T = (X^T X)^{-1} X^T \vec{y}$
- $\vec{y} \sim N_n(X\vec{\beta}, \sigma^2 I_n)$
- $\hat{\beta} \sim N_{p \times 1}(\vec{\beta}, \sigma^2 (X^T X)^{-1})$
- $\hat{\epsilon} = X\hat{\beta} \sim N_n(X\vec{\beta}, \sigma^2 X(X^T X)^{-1} X^T)$

12) Model Selection

a) AIC =  $-2 l(\hat{\theta}_{ML}) + 2p$  <sup>total # parameters estimated</sup>

b) BIC =  $-2 l(\hat{\theta}_{ML}) + p \log(n)$

13) Prediction

a) Plug-in Prediction (Distribution)

$$f(y) = f(y; \hat{\theta}_{ML}) \leftarrow \text{just plug in} \quad \text{+ ignore uncertainty of } \hat{\theta}$$

b) Extended Likelihood

$$L(\theta, y) = f(x_{1:n}, y; \theta) = f(y; \theta) \cdot \prod_{i=1}^n f(x_i; \theta)$$

$\rightarrow$  focus on the kernel

c) Predictive Likelihood

$$L_{\text{pred}}(y) = \max_{\theta} L(\theta, y) = L(\hat{\theta}(y), y)$$

$\hookrightarrow$  corr. distribution  $f_p(y) = L_{\text{pred}}(y) / \int L_{\text{pred}}(y) dy$

d) Assessment of Predictions

• PIT  $C_{\theta,y} = F_{C_{\theta,y}} \sim U(0,1)$

• CRPS =  $\int \{F(t) - \mathbb{I}_{[C_{\theta,y} < y]}(t)\}^2 dt$

lower is better!  $\mathbb{I}$  indicator with  $=1$  for  $g_0 \leq t$

14) Bayesian Inference (parameter is R.V.)

a)  $f(\theta|x) = \frac{f(x|\theta) \cdot f(\theta)}{f(x)}$

$$= \underline{L(\theta)} \propto f(x|\theta) \cdot f(\theta)$$

b) Post. estimator:  $E[\theta|x] = \int \theta f(\theta|x) d\theta$

• Post. mode:  $\text{Mod}(\theta) = \arg \max_{\theta} f(\theta|x)$

• Post. median: any number that satisfies  $\int_{-\infty}^a f(\theta|x) d\theta = 0.5 \& \int_a^{+\infty} f(\theta|x) d\theta = 0.5$

c) Improper Prior:  $\int_{\Theta} f(\theta) d\theta = \infty$  or  $\sum_{\theta} f(\theta) = \infty$

d) Jeffrey's Prior:  $f(\theta) \propto \sqrt{J(\theta)}$   $\hookrightarrow$  take small proportion

e) Credible Interval:  $\int_{\theta_L}^{\theta_U} f(\theta|x_{\text{min}}) d\theta = \gamma$

$\hookrightarrow$  R.V.  $\theta|x$  contained with prob.  $\gamma$

$\hookrightarrow$  just take  $\frac{1-\gamma}{2}$  quantiles of posterior

$\hookrightarrow$  HPD interval for equi-tailed = Wald Type Interval

f) Bayesian Testing:  $P(C|H_0 \wedge y_{\text{obs}}) = \frac{P(C|y_{\text{obs}} \mid H_0) \cdot P(H_0)}{P(C|y_{\text{obs}} \mid H_1) \cdot P(H_1)}$   $\hookrightarrow$  Bayes Factor

g) Conjugates: Likelihood  $\beta_{\text{con}} \text{, prior } \beta_{\text{prior}} \rightarrow \text{Post. } \beta_{\text{post}}$   
 $\text{Post. } \beta_{\text{exp}} \rightarrow \beta_{\text{post}} \sim \text{Exp}(\beta_{\text{prior}} + \text{Post. } \beta_{\text{exp}})$   
 $\text{Post. } \beta_{\text{norm}} \rightarrow \beta_{\text{post}} \sim N(\beta_{\text{prior}} + \text{Post. } \beta_{\text{norm}}, \text{Post. } \sigma^2)$   
 $\text{Post. } \beta_{\text{unif}} \rightarrow \beta_{\text{post}} \sim \text{Unif}(\beta_{\text{prior}}, \beta_{\text{post}})$

h) Normal, Normal Model ( $y|u \sim N(u; u \sim N)$ )

$$f(u|y) \propto \exp \left( -\frac{1}{2} \left( \frac{u - \bar{y}}{\sigma^2} + \frac{y - \mu}{\sigma^2} \right)^2 \right) \left( u - \left( \frac{n\bar{y} + 1}{n+1} \right) \right)^{\frac{n}{2}} \left( \frac{\bar{y}}{\sigma^2} + \frac{1}{\sigma^2} \right)^{\frac{n+1}{2}}$$

$\hookrightarrow$  precision of data + precision of prior  $\hookrightarrow$  posterior mean