# Last weeks

- ▶ likelihood

- ▶ marginal and conditional likelihood

- ▶ profile likelihood

- ▶ adjusted profile likelihood

- ▶ composite likelihood

# This week

- semiparametric likelihoods

- nonparametric likelihoods

- consistency of maximum likelihood estimators

- comments on problem sets

# Survival Data: single sample

- Model: $f(t), h(t), 1 - F(t), H(t)$
  density, hazard, survivor function, cumulative hazard
- Data: $(t_1, \delta_1), \ldots, (t_n, \delta_n)$
  - $t_i$ an observed time
  - $\delta_i = 1$ if $t_i$ a true failure time, 0 if $t_i$ is a censoring time
- random censorship assumption
- parametric inference:

$$L(\theta; \underline{t}, \underline{\delta}) = \sum_{i=1}^{n} \delta_i \log h(t_i; \theta) - H(t_i; \theta)$$

- examples:
  - $h(t; \lambda) = \lambda$
  - $h(t; \theta = (\lambda, \alpha)) = \lambda t^{\alpha}$
  - $f(t; \theta = \nu, \mu) = \text{Gamma}(\nu, \mu)$
  - $f(t; \theta = (\mu, \sigma^2)) = \log \text{Normal}(\mu, \sigma^2)$ ...

# Parametric regression models

- ▶ Data: $(t_i, \delta_i, \underline{x}_j), \ldots, j = 1, \ldots, n$

- ▶ Likelihood function:

$$L(\theta; \underline{t}, \underline{\delta}) = \sum_{i=1}^{n} \delta_i \log h(t_i; \theta) - H(t_i; \theta)$$

- ▶ Example: Exponential distribution
    - ▶ $h(t; \beta) = \exp(x_i^T \beta)$, for example
    - ▶ $\ell(\beta) = \sum_{i=1}^{n} \delta_i x_i^T \beta - \exp(x_i^T \beta) t_i$
    - ▶ usual maximum likelihood theory applies

- ▶ Example: Weibull distribution
    - ▶ $h(t; \theta) = h(t; \beta, \alpha) = \exp(x_i^T \beta) t^\alpha$
    - ▶ $\theta = (\beta, \alpha)$
    - ▶ usual maximum likelihood theory applies

# Semi-parametric regression models

- proportional hazards model:

$$h(t; x, \beta) = h_0(t) \exp(x^T \beta)$$

- $h_0(t)$ unknown

-
$$\frac{h(t; x)}{h(t; 0)} = \exp(x^T \beta), \quad \text{does not depend on } t$$

-
$$1 - F(t; x) = \{1 - F_0(t)\}^{\exp(x^T \beta)}$$

- survivor functions can never cross
- $x^T \beta = x_1 \beta_1 + \cdots + x_p \beta_p, \quad$ no constant term

## Estimation of $\beta$

- partial likelihood

$$L_{part}(\beta) = \prod_{i=1}^{n} \left( \frac{\exp(x_i^T \beta)}{\sum_{k \in \mathcal{R}_i} \exp(x_k^T \beta)} \right)^{\delta_i}$$

- $\mathcal{R}_i$ risk set at time $t_i^-$; number of units with $t_i \geq t_i$

- derived in SM §10.8 as approximately a profile likelihood ($h_0(\cdot)$ maximized out)

- $\hat{\beta}$ estimated by maximizing partial log-likelihood $\ell_{part}(\beta) = \log L_{part}(\beta)$

- estimated standard error from $-\ell''_{part}(\hat{\beta})$

# ... partial likelihood

- $\hat{\beta}$ estimated by maximizing partial log-likelihood
  $\ell_{part}(\beta) = \log L_{part}(\beta)$

- estimated standard error from $-\ell_{part}''(\hat{\beta})$

- usual asymptotic theory applies: $\hat{\beta} \overset{.}{\sim} N(\beta, -\ell_{part}''(\hat{\beta}))$

- special property of this model: components of the score vector are uncorrelated

- no need to compute analogue of Godambe information

- there could be loss of efficiency in estimating $\beta$; this loss has been shown to be small in a wide range of settings

- general treatment of likelihood inference for semi-parametric models    Murphy and van der Waart, 2000

# Semi-parametric regression models

- for example, $E(y_i) = \mu_i(\theta) = x_i^T \beta + m(t_i)$, $\quad \text{Var}(y_i) = \sigma^2$
- $m(\cdot)$ a 'smooth' function of covariates $t$

- least squares

$$\min_{\beta, m(\cdot)} \sum_{i=1}^{n} \{y_i - x_i^T \beta - m(t_i)\}^2$$

- without constraint on $m(\cdot)$, minimum will be 0, thus

$$\min_{\beta, m(\cdot)} \sum_{i=1}^{n} \{y_i - x_i^T \beta - m(t_i)\}^2 - \frac{1}{2} \lambda \int \{m''(t)\}^2 dt$$

- equivalent to

$$\min_{\beta, m(\cdot)} \sum_{i=1}^{n} \{y_i - x_i^T \beta - m(t_i)\}^2 + \lambda m^T K m$$

for suitable $n \times n$ matrix $K$

## ... semi-parametric regression models

- $$\min_{\beta, m(\cdot)} \sum_{i=1}^{n} \{y_i - x_i^T \beta - m(t_i)\}^2 - \frac{1}{2}\lambda \int \{m''(t)\}^2 dt$$

- extend to generalized linear model

$$h\{\mathsf{E}(y_i)\} = x_i^T \beta + m(t_i) = \eta_i$$

- penalized log-likelihood

$$\min_{\beta, m(\cdot)} \sum_{i=1}^{n} \ell_i(\eta_i) - \frac{1}{2}\lambda \int \{m''(t)\}^2 dt$$

Green, 1987; SM, §10.7

# Nonparametric likelihood

- likelihood functions for infinite-dimensional parameters can be tricky
- for example, given $y_1, \ldots, y_n$ i.i.d. with distribution function $F(\cdot)$ and density function $f(\cdot)$
- the nonparametric maximum likelihood estimator of $F(\cdot)$ is

$$F_n(t) = \frac{1}{n} \sum_{i=1}^{n} \mathbf{1}(Y_i \leq t), \quad t \in \mathbb{R}$$

- this is a cumulative distribution function, although discrete
- the nonparametric maximum likelihood estimator of $f(\cdot)$ is not a density function
- unless we put some constraints on the class of densities over which we maximize
- for example, might require $f(x)$ to be log concave: $f(x) = \exp\{\eta(x)\}, \eta$ concave     Balabdaoui et al, 2009

# Empirical likelihood

- $y_1, \ldots, y_n$ i.i.d. with distribution function $F_0(\cdot)$
- define

$$L(F) = \prod_{i=1}^{n} \{F(y_i) - F(y_i^-)\}$$

- maximized at $F_n$, empirical c.d.f.
- empirical likelihood ratio

$$R(F) = \frac{L(F)}{L(F_n)}$$

- suppose $T(F_0)$ is a function of interest, e.g. $\mu = \int x dF(x)$
- maximizing $R(F)$, subject to $\mu$ fixed, is equivalent to

$$\max_{w_1, \ldots, w_n} \prod_{i=1}^{n} w_i, \text{ subject to } \sum_{i=1}^{n} w_i y_i = \mu, \sum_{i=1}^{n} w_i = 1, w_i \geq 0, \forall i$$

Owen, 1988; 2001

## ... empirical likelihood

- 
$$\max_{w_1,\ldots,w_n} \prod_{i=1}^{n} w_i, \text{ subject to } \sum_{i=1}^{n} w_i y_i = \mu, \sum_{i=1}^{n} w_i = 1, w_i \geq 0, \forall i$$

- likelihood ratio confidence intervals are valid

$$-2 \log R(F_0) \xrightarrow{\mathcal{L}} \chi_1^2, \quad n \to \infty$$

- parameter of interest, $\mu \in \mathbb{R}$
- nuisance parameter $w = (w_1, \ldots, w_n)$
- generalized to many more complex situations

Hjort et al. 2009

# Those pesky regularity conditions

- two proofs of the consistency of the maximum likelihood estimator
- Wald, 1949 – the log-likelihood is maximized in expectation at the true value; apply Jensen's inequality to conclude $\hat{\theta}$ must converge to the true value
- requires the parameter space to be compact

- Cramer, 1946 – there exist solutions to the score equation that are consistent
- Taylor series expansion of $\log f(y; \theta)$
- if the likelihood function is maximized in the interior of the parameter space, the m.l.e. is one of these solutions
- if the score equation has only one root, the m.l.e. is consistent

## Non-standard cases

- true parameter $\theta_0$ on the boundary of the parameter space
- example: $y_{ij} = \mu + b_i + \epsilon_{ij}, \quad b_i \sim N(0, \sigma_b^2), \epsilon_{ij} \sim N(0, \sigma^2)$
- if $\sigma_b^2 = 0$, no difference between groups; this is a boundary point of the parameter space

- non-identifiability; two different $\theta_1$, $\theta_2$ for which $f(y; \theta_1) = f(y; \theta_2)$
- example $f(y; \theta) = pN(\mu_1, 1) + (1 - p)N(\mu_2, 1)$
- if $\mu_1 = \mu_2$, then $p$ is not identifiable
- if $p = 0$ ,then $\mu_1$ is not identifiable
- likelihood ratio test of, e.g. $H_0 : p = 0$ will not be asymptotically $\chi^2$

# ... non-standard cases

- multi-modal log-likelihoods
- in principle, find all the stationary points, and choose that corresponding to the maximum
- in practice, may not be feasible
- example: feed-forward neural networks;

- support of the distribution depends on the parameter
- example $U(0, \theta)$; $n(y_{(n)} - \theta) \xrightarrow{\mathcal{L}}$ Exponential
- example $f(y; \theta) = \lambda \exp\{-\lambda(y - \mu)\}$

SM, §4.6; BNC94, §3.8; Cox, Ch. 7

## ... non-standard cases

- singular information matrix: $\text{var}_{\theta_0}\{U(\theta_0)\} \equiv 0$
- usual Taylor series expansions do not apply; need to go to higher order terms
- might be fixable by re-parameterization

- Example: skew-normal distribution
- $Z \sim SKN(\alpha) : f_Z(z; \alpha) = 2\phi(z)\Phi(\alpha z)$
- three-parameter version: $Y = \xi + \omega Z$
- information matrix is singular, at $\alpha = 0$
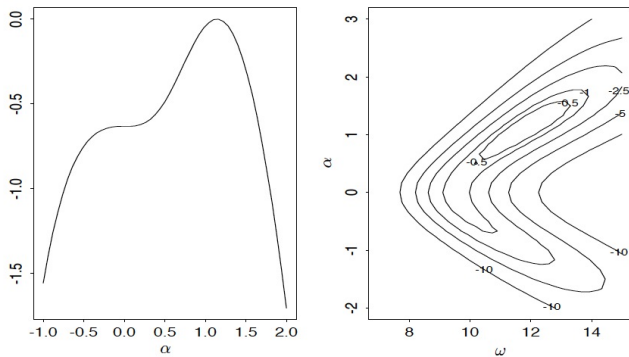- can be fixed by reparametrization to $(\mu, \sigma, \alpha)$

Azzalini, 1999; 2011

Figure 2: *Twice relative profile loglikelihood of $\alpha$ (left) and contour levels of the similar function of $(\omega, \alpha)$ (right) for the Otis data, when the direct parametrization is used*

Azzalini, 1999

## Problems – Week 4

1. Suppose $y = y_1, \ldots, y_n$ are independent and identically distributed from a distribution with density $f(y; \theta) = \prod_{i=1}^{n} f_1(y_i; \theta)$, $\theta \in R$. Further let $g(y; \theta) = \sum_{i=1}^{n} g_1(y_i; \theta)$ be an *unbiased estimating equation* for $\theta$, satisfying $E_\theta\{g(y_i; \theta)\} = 0$ for all $\theta$. The estimate defined by $g(y; \tilde{\theta}_g) = 0$ has asymptotic variance $G^{-1}(\theta) = H^{-1}(\theta) J(\theta) H^{-1}(\theta)$, where $H(\theta) = -E_\theta\{\nabla_\theta g(y_1; \theta)\}$ and $J(\theta) = \text{var}_\theta\{g(y_1; \theta)\}$. The estimating equation is called *optimal* if it has the largest possible value of $G(\theta)$.

   Show that $G(\theta) \leq i_1(\theta)$, where $i_1(\theta)$ is the expected Fisher information in a single observation. This implies that the score equation is the optimum estimating equation.

   Two fun facts that you don't need to prove:

   (a) The multivariate version of this is that $i_1(\underline{\theta}) - G(\underline{\theta})$ is non-negative definite (but you don't need to show this).

   (b) In the autoregressive model

   $$y_i = \theta y_{i-1} + \epsilon_i, \quad i = 1, \ldots, n$$

   where $y_0$ is a constant and $\epsilon_i$ are i.i.d. $N(0, \sigma^2)$, show the equation

   $$\Sigma y_i y_{i-1} - \theta \Sigma y_i^2 = 0$$

   is an unbiased estimating equation obtaining the lower bound.

# Problems – Week 4

2. Suppose $Z \sim \sum_{r=1}^{m} \mu_r X_r^2$, where $X_1, \ldots, X_m$ are independent observations from a $N(0,1)$ distribution. If all the $\mu_r$ were equal, the distribution of $Z$ would be proportional to a $\chi_m^2$. *Satterthwaite's approximation* (Satterthwaite, 1946) to the distribution of $Z$ is $a\chi_b^2$, where $a$ and $b$ are chosen so that $\mathrm{E}(Z)$ and $\mathrm{var}(Z)$ are equal to the mean and variance of a $a\chi_b^2$ random variable. This idea can also be used to approximate a non-central $\chi^2$ distribution, and arises in the distribution of quadratic forms in unbalanced analysis of variance.

   (a) Find expressions for $a$ and $b$, in terms of $\mu_1, \ldots, \mu_m$.

   (b) Illustrate the approximation numerically in a simple example with, say, $m = 5, 10$. You can choose the values of $\mu_r$ in any way you like, but one possibility is to simulate a random vector from $N(0, A)$ for some choice of $A \neq I$; then $X^T X$ will (I think), have the distribution you are looking for. The function `mvrnorm` in the `MASS` library simulates multivariate normal random variables.

# Problems – Week 3

1. Suppose $Y_1, \ldots, Y_n$ are independent and identically distributed from a model $f(y; \theta), y \in R, \theta \in R$, and that $\pi(\theta)$ is a proper prior density (with respect to Lebesgue measure on $R$). Denote by $\hat{\theta}_\pi$ the posterior mode:

$$\hat{\theta}_\pi = \arg\sup_\theta \pi(\theta \mid y)$$

which we assume is obtained as the unique root of the equation

$$\frac{d}{d\theta} \log \pi(\hat{\theta}_\pi \mid y) = 0. \tag{1}$$

Denote by $\tilde{\theta}$ the posterior mean:

$$\tilde{\theta} = \int \theta \pi(\theta \mid y) d\theta.$$

Show that

$$\hat{\theta}_\pi - \hat{\theta} = O_p(\frac{1}{n}), \text{ and } \tilde{\theta} - \hat{\theta} = O_p(\frac{1}{n}),$$

where $\hat{\theta}$ is the maximum likelihood estimator of $\theta$.

# Problems – Week 3

2. Consider a linear regression model

$$y_i = x_i^T \beta + \epsilon_i, \quad i = 1, \dots, n$$

where $x_i$ and $\beta$ are $p \times 1$ vectors, and $\epsilon_i$ are i.i.d. $N(0, \sigma^2)$. Compare the log-likelihood ratio statistics for inference about $\beta$, based on the

(a) profile log-likelihood $w(\beta) = 2\{\ell_p(\hat{\beta}) - \ell_p(\beta)\}$,

(b) adjusted profile log-likelihood $w_A(\beta) = 2\{\ell_A(\hat{\beta}_A) - \ell_A(\beta)\}$, and

(c) modified profile log-likelihood $w_M(\beta) = 2\{\ell_M(\hat{\beta}_M) - \ell_M(\beta)\}$,

where

$$\ell_p(\beta) = \ell(\beta, \hat{\sigma}_\beta^2), \quad \ell_A(\beta) = \ell_p(\beta) - \frac{1}{2}\log|j_{\sigma^2\sigma^2}(\beta, \hat{\sigma}_\beta^2)|, \text{ and } \ell_M(\sigma^2) = \ell_A(\beta) + \log|\frac{d\hat{\sigma}^2}{d\hat{\sigma}_\beta^2}|,$$

and $\hat{\beta}_A$, $\hat{\beta}_M$ are the adjusted and modified maximum likelihood estimators, respectively.

# Problems – Week 2

1. Suppose $Y_1, \ldots, Y_n$ are i.i.d. with density

$$f_{Y_i}(y; \mu) = \frac{1}{\mu} \exp(-\frac{y}{\mu}), y > 0, \mu > 0.$$

Show that the leading term in the saddlepoint approximation to the density of $\bar{Y} = \hat{\mu}$ reproduces the gamma density, with $\Gamma(n)$ replaced by Stirling's approximation to it. Deduce that the renormalized saddlepoint approximation is exact.

# Problems – Week 2

(a) Suppose that $Y_{i1}$ and $Y_{i2}$ are independent observations from exponential distributions with means $\psi\lambda_i$ and $\psi/\lambda_i$, respectively, $i = 1, \ldots, n$. Show that the maximum likelihood estimator of $\psi$ is not consistent, but converges in probability to $(\pi/4)\psi$.

(b) A modification to the profile likelihood to account for estimation of nuisance parameters was proposed in Cox & Reid (1987):

$$\ell_m(\psi) = \ell(\psi, \hat{\lambda}_\psi) - \frac{1}{2} \log |j_{\lambda\lambda}(\psi, \hat{\lambda}_\psi)|,$$

where $\lambda = (\lambda_1, \ldots, \lambda_n)$ and $\hat{\lambda}_\psi$ is the constrained maximum likelihood estimator of $\lambda$. This is to be computed using a parametrization of the nuisance parameter that is *orthogonal* to the parameter of interest $\psi$, with respect to expected Fisher information. (The correction term $\frac{1}{2} \log |j_{\lambda\lambda}(\psi, \hat{\lambda}_\psi)|$ is not invariant to reparameterizations,) Show for the exponential case that $\lambda$ is orthogonal to $\psi$, and that the value of $\psi$ that solves $\ell'_m(\psi) = 0$, $\hat{\psi}_m$, say, converges to $(\pi/3)\psi$.

# Problems – Week 1

1. *Orthogonal nuisance parameters.* In a model $f(y; \theta)$ with $\theta = (\psi, \lambda)$, the component parameter $\psi$ and $\lambda$ are orthogonal (with respect to Fisher information) if $i_{\psi\lambda}(\theta) = 0$.

   (a) Suppose we have a sample $y_1, \ldots, y_n$ from the density $f(y; \theta)$. Show that

   $$\hat{\lambda}_\psi = \hat{\lambda} + O_p(n^{-1/2}),$$

   whereas if $\psi$ and $\lambda$ are orthogonal that

   $$\hat{\lambda}_\psi = \hat{\lambda} + O_p(n^{-1}).$$

   (b) Assume $y_i$ follows an exponential distribution with mean $\lambda e^{-\psi x_i}$, where $x_i$ is known. Find conditions on the sequence $\{x_i, i = 1, \ldots, n\}$ in order that $\lambda$ and $\psi$ are orthogonal with respect to expected Fisher information. Find an expression for the constrained maximum likelihood estimate $\hat{\lambda}_\psi$ and show the effect of parameter orthogonality on the form of the estimate.

# Problems – Week 1

2. *Sufficient statistics (CH Exercise 2.2).* Find the log-likelihood function for a sample of size $n$ from an $AR(1)$ process:

$$y_t = \mu + \rho(y_{t-1} - \mu) + \epsilon_t, \quad \epsilon_t (i.i.d.) \sim N(0, \sigma^2), \quad t = 1, ..., n,$$

where $|\rho| < 1$, as a function of $\theta = (\mu, \sigma^2, \rho)$ and $y_0$. Write down the likelihood for data $y_1, \ldots, y_n$ in the cases where the initial value $y_0$ is

(a) a given constant;

(b) normally distributed with mean $\mu$ and variance $\sigma^2/(1 - \rho^2)$;

(c) assumed equal to $y_n$,

and give the sufficient statistic for each case.

# Extra notes for HW1, 3

## Notes to help

**LTCC/Reid: Derivation of limiting results: scalar parameter**   November 6, 2012

Using the notation on the handout from November 5, ("week1-handout.pdf"), here is a moderately rigorous proof of the results

$$\sqrt{n}(\hat{\theta} - \theta) = \frac{1}{\sqrt{n}} i_1^{-1}(\theta) U(\theta)\{1 + o_p(1)\}, \tag{1}$$

$$2\{\ell(\hat{\theta}) - \ell(\theta)\} = (\hat{\theta} - \theta)^T i(\theta)(\hat{\theta} - \theta)\{1 + o_p(1)\}. \tag{2}$$

The vector case is unchanged, except for tedious notational changes in Taylor's theorem with remainder, although of course we need the dimension of $\theta$ fixed as $n \to \infty$.

For (1), we have

$$\ell'(\hat{\theta}) = \ell'(\theta) + (\hat{\theta} - \theta)\ell''(\theta) + \frac{1}{2}(\hat{\theta} - \theta)^2 \ell'''(\theta_n^*),$$

$$-\frac{\ell'(\theta)}{\ell''(\theta)} = (\hat{\theta} - \theta)\{1 + \frac{1}{2}(\hat{\theta} - \theta)\frac{\ell'''(\theta_n^*)}{\ell''(\theta)}\},$$

$$\frac{\frac{1}{\sqrt{n}}\ell'(\theta)}{-\ell''(\theta)/n} \cdot \frac{i_1(\theta)}{i_1(\theta)} = \sqrt{n}(\hat{\theta} - \theta)\{1 - \frac{1}{2}(\hat{\theta} - \theta)\frac{\ell'''(\theta_n^*)/n}{-\ell''(\theta)/n}\},$$

$$\frac{\frac{1}{\sqrt{n}}\ell'(\theta)}{i_1(\theta)}\left(\frac{i_1(\theta)}{-\ell''(\theta)/n}\right) = \sqrt{n}(\hat{\theta} - \theta)\{1 + Z_n\}.$$

The term in brackets on the LHS of the last line converges in probability to 1, by the WLLN, so can be written $1 + o_p(1)$. The remainder term $Z_n$ converges in probability to 0, because we assume $\hat{\theta} \overset{p}{\to} \theta$, so that $\theta_n^* \overset{p}{\to} \theta$, because $|\theta_n^* - \theta| < |\hat{\theta} - \theta|$. Also $\frac{1}{n}\ell'''(\theta_n^*) \overset{p}{\to} E\{\ell'''(\theta; Y)\}$ which we assume is finite (p.281 of CH, for example); similarly $-\frac{1}{n}\ell''(\theta) \overset{p}{\to} i_1(\theta)$, so $Z_n = o_p(1)O_p(1) = o_p(1)$. Then we can move over the LHS term as

$$\frac{1}{\sqrt{n}}\frac{\ell'(\theta)}{i_1(\theta)}\{1 + o_p(1)\} = \sqrt{n}(\hat{\theta} - \theta)$$