

#### Various 'types' of likelihood

- 1. likelihood, marginal and conditional likelihood, profile likelihood, adjusted profile
- 2. semi-parametric likelihood, partial likelihood
- 3. quasi-likelihood, composite likelihood
- 4. empirical likelihood, penalized likelihood

5. (likelihood) inference in high dimensions
6. simulated likelihood, indirect inference

7. bootstrap likelihood, *h*-likelihood, weighted likelihood, pseudo-likelihood, local likelihood, sieve likelihood

misspecified models

#### Presentations

Feb 16 Angela: Cox (2013) Robert: Barndorff-Nielsen and Cox (1979) Shiki: Solomon and Cox (1992)

Feb 23 Hengchao: Rotnitzky et al. (2000) Siyue: De Stavola and Cox (2008) Manuel: Battey and Cox (2018) Ziang: Cox (1975) Libai: Cox €1999)

Feb 16 in SS 1087; Feb 23 online

exercises Jan 26 has details about report structure

#### **High-dimensional inference**

•  $f(y; \theta), y \in \mathbb{R}^{n}; \theta \in \mathbb{R}^{p}, p \text{ large relative to } n, \text{ or } p > n$ • Partial likelihood has p = n - 1 + d, yet usual asymptotic theory applies • Empirical likelihood has p = n - 1, yet usual asymptotic theory applies

"Neyman-Scott problems" have  $\underline{y_{ij}} \sim f(\cdot; \psi, \lambda_i), j = 1, \dots, \underline{m}; i = 1$ 

Padapends an ? note?  
usually 
$$f_n^n \rightarrow 0$$
  $f_n^n \rightarrow 0$ 

## Empirical likelihood

• 
$$Y_1, \dots, Y_n$$
 i.i.d.  $F$   $\mu = E(Y_i) = \int y dF(y)$   
• profile likelihood: maximize  $\prod_{i=1}^n p_i$ , subject to  $p_i \ge 0, \Sigma p_i = 1, \Sigma p_i Y_i = \mu$   
• solution  
• solution  
 $\hat{p}_i(\mu) = \frac{1}{n(1 + \lambda(Y_i - \mu))}$  where  $\lambda$  solves  
 $0 = \frac{1}{n} \sum_{i=1}^n \frac{Y_i - \mu}{1 + \lambda(Y_i - \mu)}$   
• Theorem:  $R(Y) = \{\mu : g(Y; \mu) \le k_\alpha\}$  is an approximate  $1 - \alpha$  confidence interval for  $\mu$ , where  $\Pr\{\chi_1^2 \ge k_\alpha\} = \alpha$   
• if  $E[Y_i]^3 < \infty$ , under  $H_0: \mu = \mu_0$ :  
 $-2 \sum_{i=1}^n \log\{n\hat{p}_i(\mu_0)\} \stackrel{d}{\to} \chi_1^2, \quad n \to \infty$ 

4

 $+2\log\frac{L(\hat{p})}{L(\tilde{p}(\mu))} = 2l_{p}\Pi\frac{(\mu\hat{p}_{i})}{i=n} = \frac{2}{\hat{p}_{i}}\Pi\frac{i}{\hat{p}_{i}} = \frac{1}{\hat{p}_{i}}\Pi\frac{i}{\hat{p}_{i}}$   $= 2l_{p}\Pi\frac{i}{\hat{p}_{i}}(\mu)$   $= 2l_{p}\Pi\frac{i}{\hat{p}_{i}}(\mu)$  $= -2\Sigma b g(x; \mu) \leq (g(x; \mu))$  $\lambda_{r} = \lambda_{r}$ Ś

 $-ZZ \log n \widetilde{p}_{i}(\mu) \xrightarrow{d} \chi_{1}^{2}$ ₱: [m) is defind as the I to ac )]

### ... empirical likelihood

#### Owen 1988; Knight Ch5.6

- proof first shows
- then

• then Taylor series expansion:

$$-2\sum_{i=1}^{n}\log\{n\hat{p}_{i}(\mu_{0})\} \doteq \frac{n(\bar{Y}-\mu_{0})^{2}}{S(\mu_{0})}$$

 $\lambda \neq O_p(n^{-1/2})$ 

 $S(\mu) = \frac{1}{n} \sum_{i=1}^{n} (Y_i - \mu)^2$ 

• see Owen (2001) *Empirical Likelihood* for many generalizations, including proportional hazards model

#### Nevman-Scott problems

Neyman-Scott problems  
where 
$$\sigma^{2} = \frac{1}{qm} \sum_{i=1}^{2} (Y_{i} - Y_{i})^{2} \sum_{i=1}^{m-1} \sum_{i=1}^{2} partori 2003$$
  
 $(Y_{ij} \sim N(\mu_{i}, \sigma_{i}^{2}), \quad j = 1, ..., m; \quad i = 1, ..., q; \quad \theta = (\mu_{1}, ..., \mu_{q}, \sigma^{2})$  Sartori's notation  
 $(Y_{ij} \sim N(\mu, \sigma_{i}^{2}), \quad j = 1, ..., m; \quad i = 1, ..., q; \quad \theta = (\mu, \sigma_{1}^{2}, ..., \sigma_{q}^{2})$   
 $(Y_{ij} \sim Bern(p_{ij}), \quad j = 1, ..., m; \quad i = 1, ..., q; \quad \theta = (\mu, \sigma_{1}^{2}, ..., \sigma_{q}^{2})$   
 $(Y_{ij} \sim Bern(p_{ij}), \quad j = 1, ..., m; \quad i = 1, ..., q; \quad \psi = \log\{\frac{p_{in}(1-p_{in})}{p_{in}(1-p_{in})}\}; \quad \lambda_{1}, ..., \lambda_{q} = \{hp(1-p_{in}), hq(1-p_{in}), hq(1-p$ 

12



Sartori 2003

#### N. Sartori

Table 1: Example 2. Inference about common shape parameter in q gamma samples of size m. Probabilities  $\Phi\{r_P(\psi)\}$  and  $\Phi\{r_M(\psi)\}$  with  $\psi$ such that  $\Phi\{r_C(\psi)\} = 0.05$  in samples with  $\hat{\psi}_C = 1$ . For each q, values in bold face correspond to the smallest m, for  $r_P$  and  $r_M$ , such that the probability is within 0.01 of 0.05.

	т		q = 4	q = 8	q = 16	q = 64	q = 128	0.05
	(3	$r_{\rm P}$	0.190	0.299	0.487	0.952	0.999	1 +
	$\sim$	$r_{\rm M}$	0.053	0.055	0.028	0.070	0.080	tagel
m 200	4	r <sub>P</sub>	0.159	0.239	0.383	0.866	0.988	_
		$r_{\rm M}$	0.052	0.053	0.055	0.062	0.068	
	5	r <sub>P</sub>	0.141	0.206	0.322	0.777	0.962	
		$r_{\rm M}$	0.052	0.052	0.054	0.058	0.062	
	6	r <sub>P</sub>	0.129	0.184	0.281	0.698	0.924	-
		$r_{\rm M}$	0.051	0.052	0.053	0.056	0.029	
	200	r <sub>P</sub>	0.060	0.064	0.071	0.098	0.126	
		$r_{\rm M}$	0.050	0.020	0.020	0.020	0.020	
	400	r <sub>P</sub>	0.056	0.059	0.064	0.081	0.098	հ
		$r_{\rm M}$	0.050	0.020	0.020	0.020	0.020	<u>-</u>
	800	r <sub>P</sub>	0.054	0.057	0.059	0.071	0.080	
		$r_{\rm M}$	0.050	0.020	0.020	0.020	0.020	
	3000	r <sub>P</sub>	0.052	0.053	0.055	0.059	0.064	
		r <sub>M</sub>	0.050	0.020	0.020	0.020	0.020	
	6000	$r_{\rm P}$	0.051	0.053	0.053	0.057	0.060	
				0.0.50			0.0.50	

#### **Increasing dimension asymptotics**



### High-dimensional linear regression

Bühlmann et al 2014

• 
$$y = X\beta + \epsilon$$
,  $E(\epsilon) = 0$ ,  $Cov(\epsilon) = \sigma^2 I$   $p >> n$  running example,  $n = 71, p = 4088$   
 $\hat{\beta}_{ridge} = \arg\min_{\beta} \{\frac{1}{n}(y - X\beta)^{T}(y - X\beta) + \lambda \sum_{j=1}^{p} \beta_{j}^{2}, \}$ 
 $\hat{\beta}_{lasso} = \arg\min_{\beta} \{\frac{1}{n}(y - X\beta)^{T}(y - X\beta) + \lambda \sum_{j=1}^{p} |\beta_{j}|\}$ 

• usual to assume 
$$\sum_{i=1}^{n} x_{ij} = 0, \sum_{i=1}^{n} x_{ij}^2 = 1$$

so units are comparable  $\hat{\beta}_{0} = \bar{y}$  is not 'shrunk'

• Lasso penalty leads to several  $\hat{\beta}_k = 0$ 

sparse solutions

- there are many variations on the penalty term
- $\lambda$  is a tuning parameter

often selected by cross-validation

### ... high-dimensional inference

- Inferential goals (§2.2)
  - (a) prediction of surface  $X\beta$  or  $y_{new} = x_{new}^{\mathrm{T}}\beta$
  - (b) estimation of  $\beta$
  - (c) estimation of  $S = \{j : \beta_j \neq 0\}$

e.g.  $\{j; \hat{\beta}_{j,Lasso}\}$ 

also needs conditions on X

'support set'

- (a): "identifiability of  $\beta$  is not necessary; from this perspective, prediction is often a much easier problem"
- (b): "an identifiability assumption on the design X is required, for example, a restricted eigenvalue condition" not checkable (?)

 $P_n(\hat{s}=s) > 1-\epsilon \longrightarrow P_n(\hat{s} \geq s) > 1-\epsilon$ 

- (c): "ideally, ...  $\hat{S} = S$  with high probability"  $\hat{S}$  estimate of S
  - (c): requires  $\beta_{min}$  condition:  $\min |\beta_j| > c$ ,  $c \sim \sqrt{\log p/n} \le 5$
  - often replaced by (c'): 'screening'  $S \supset S$  with high probability

STA 4508 February 16 2022

11

### ... high-dimensional inference



- re (a) prediction accuracy can be assessed by cross-validation
- note that (a) can be estimable even if p > n, as long as  $X\beta$  small enough dimension

# **Inference about** $\beta$ , p > n

- re (b): inference for  $\beta$ : e.g. *p*-values for testing  $H_{o,j}$ :  $\beta_j = O_{a,j}$
- three methods suggested: multi-sample splitting, debiasing, stability selection

multi-sample splitting: fit the model on random half, say, of observations, leads to  $\hat{S}$ 

- use  $X^{(\hat{S})}$  in fitting to the other half
- $P_j = p$ -value for t-test of  $H_j$  if  $j \in \hat{S}$ , o.w. 1
- $P_{corr,j} = P_j \times |\hat{S}|, j \in \hat{S}, 0.W. 1$  Bonferron corr
- Repeat *B* times and aggregate *P<sup>b</sup><sub>corr</sub>*

STA 4508 February 16 2022

 $P^{b}_{corr,j}$  not independent

0=71



#### Non-linear models

- example  $y_i$  independent,  $E(y_i) = \mu_i(\beta)$ ;  $g(\mu_i) = \beta_o + x_i^{T}\beta$
- Lasso-type 'mle':  $\arg\min\{-\frac{1}{n}\ell(\beta,\beta_0;\mathbf{y}) + \lambda \Sigma_{j=1}|\beta_j|\}$   $\beta = (\beta_1,\ldots,\beta_p)$
- can use multi-sample splitting or stability selection
- a version of de-biasing applies to GLMs, based on weighted least squares

- a marginal approach would fit  $y = \alpha_0 + \alpha_j x^{(j)}$ , one feature at a time
- leading to 4088 *p*-values, and then need techniques for controlling FWER or FDR

#### • Model: $y_i = x_i^{\mathrm{T}}\beta + Z_i$ , $i = 1, \dots, n$

• M-estimation:

 $n, p \rightarrow \infty$ 

$$\sum_{i=1}^{n} x_i \psi(y_i - x_i^{\mathrm{T}} \hat{\beta}) = 0$$
(1)

- result: if  $\psi$  is monotone, and  $p \log(p)/n \rightarrow 0$ , and conditions on X, then

there is a solution of (1) satisfying  $||\hat{eta}-eta||^2=O(p/n)$ 

• "rows of X behave like a sample from a distribution in  $\mathbb{R}^{p^n}$ "

independent

independent

 $n, \overline{p} \rightarrow \overline{\infty}$ 

- Model:  $y_i = x_i^{\mathrm{T}}\beta + Z_i$ ,  $i = 1, \dots, n$
- M-estimation:

$$\sum_{i=1}^{n} x_i \psi(y_i - x_i^{\mathrm{T}} \hat{\beta}) = 0$$
(1)

- result: if  $\psi$  is monotone, and  $p\log(p)/n \to 0$ , and conditions on X, then

there is a solution of (1) satisfying  $||\hat{eta} - eta||^2 = O(p/n)$ 

- "rows of X behave like a sample from a distribution in  $\mathbb{R}^{p}$ "
- if  $p^{3/2} \log n/n \to 0$ , then

$$\max |\mathbf{X}_{i}^{\mathrm{\scriptscriptstyle T}}(\hat{eta}-eta)| \stackrel{p}{
ightarrow} \mathsf{O}$$

• and

$$a_n^{\mathrm{T}}(\hat{\beta} - \beta) \stackrel{d}{\rightarrow} N(\mathbf{0}, \sigma^2)$$

 $\sigma^2 = a_n^{\mathrm{T}}(X^{\mathrm{T}}X)^{-1}a_n \mathsf{E}\psi^2(Z)/\{\mathsf{E}\psi'(Z)\}^2$