

Lec 3: Parameter Estimation & Inference

Yanyan Han
Sept 17, 2024



Given i.i.d. $y_1, \dots, y_n \sim p_\theta(y) = \exp(\langle \theta, T(y) \rangle - A(\theta)) h(y)$.

This lecture:

Parameter estimation: estimate θ or functions of θ

Inference: test $H_0: \theta = \theta_0$ against $H_1: \theta \neq \theta_0$.

Maximum likelihood estimator (MLE)

$$\begin{aligned}\hat{\theta}_n &= \arg \max_{\theta} \prod_{i=1}^n p_\theta(y_i) \\ &= \arg \max_{\theta} \sum_{i=1}^n \log p_\theta(y_i) \\ &= \arg \max_{\theta} \underbrace{\left\langle \theta, \sum_{i=1}^n T(y_i) \right\rangle - n A(\theta)}_{\text{Concave in } \theta}\end{aligned}$$

$$\text{F.O.C.: } 0 = \sum_{i=1}^n T(y_i) - n \nabla A(\hat{\theta}_n), \text{ or}$$

$$\nabla A(\hat{\theta}_n) = \frac{1}{n} \sum_{i=1}^n T(y_i)$$

- As $\mu_\theta := \mathbb{E}_\theta[T(y)] = \nabla A(\theta)$, the MLE $\hat{\theta}_n$ is chosen so that the "true mean" matches the "sample mean".
(Estimating function view: $\mathbb{E}[\frac{1}{n} \sum_{i=1}^n T(y_i) - \nabla A(\theta)] = 0$)
- The MLE either admits a closed-form expression, or is the solution to a convex optimization problem.

Example: Poisson family.

Recall that $y \sim \text{Poi}(\lambda)$, $\theta = \log \lambda$, $T(y) = y$, $A(\theta) = e^\theta$.

Therefore,

$$\text{MLE for } \theta: e^{\hat{\theta}_n} = \frac{1}{n} \sum_{i=1}^n y_i \Rightarrow \hat{\theta}_n = \log\left(\frac{1}{n} \sum_{i=1}^n y_i\right)$$

$$\text{MLE for } \lambda: \hat{\lambda}_n = e^{\hat{\theta}_n} = \frac{1}{n} \sum_{i=1}^n y_i.$$

Variance of the MLE

1. (Exact) variance for $\hat{\theta}_n = \nabla A(\hat{\theta}_n)$:

$$\begin{aligned}\nabla A(\hat{\theta}_n) &= \frac{1}{n} \sum_{i=1}^n T(Y_i) \\ \Rightarrow \text{Cov}_{\theta}(\nabla A(\hat{\theta}_n)) &= \text{Cov}_{\theta}\left(\frac{1}{n} \sum_{i=1}^n T(Y_i)\right) \\ \Rightarrow \text{Cov}_{\theta}(\nabla A(\hat{\theta}_n)) &= \frac{1}{n} \nabla^2 A(\theta)\end{aligned}$$

In reality we don't know θ , so we typically use

$$\boxed{\text{Cov}_{\theta}(\nabla A(\hat{\theta}_n)) \approx \frac{1}{n} \nabla^2 A(\hat{\theta}_n)}$$

2. Approximate variance: delta method

Question: Suppose $\hat{\theta}_n \approx \theta$ and $f(\cdot)$ is differentiable at θ .

How is $\text{Var}(f(\hat{\theta}_n))$ related to $\text{Var}(\hat{\theta}_n)$?

Idea of delta method: suppose $|\hat{\theta}_n - \theta| = o_p(r_n)$ with $r_n \rightarrow 0$.

Then

$$\begin{aligned}f(\hat{\theta}_n) &= f(\theta) + f'(\theta)(\hat{\theta}_n - \theta) + o_p(r_n) \\ \Rightarrow \text{Var}(f(\hat{\theta}_n)) &= \text{Var}[f(\theta) + f'(\theta)(\hat{\theta}_n - \theta)] + o_p(r_n^2) \\ &= f'(\theta)^2 \cdot \text{Var}(\hat{\theta}_n) + o_p(r_n^2)\end{aligned}$$

So we have:

$$\boxed{\text{1-D delta method : } \text{Var}_{\theta}(f(\hat{\theta}_n)) \approx f'(\theta)^2 \text{Var}_{\theta}(\hat{\theta}_n) \text{ if } \text{Var}_{\theta}(\hat{\theta}_n) \text{ is small}}$$

Similarly, for $f: \mathbb{R}^{d_1} \rightarrow \mathbb{R}^{d_2}$ and $\nabla f(\theta) \in \mathbb{R}^{d_1 \times d_2}$ defined as

$$(\nabla f(\theta))_{ij} = \frac{\partial}{\partial \theta_i} f_j, \quad 1 \leq i \leq d_1, 1 \leq j \leq d_2.$$

then

General delta method: $\text{Cov}_{\theta}(f(\hat{\theta}_n)) \approx \nabla f(\theta)^T \text{Cov}_{\theta}(\hat{\theta}_n) \nabla f(\theta)$
 if $\|\text{Cov}_{\theta}(\hat{\theta}_n)\|$ is small

3. Approximate variance for $\hat{\theta}_n$: by delta method,

$$\begin{aligned} \frac{1}{n} \nabla^2 A(\theta) &= \text{Cov}_{\theta}(\nabla A(\hat{\theta}_n)) \approx \nabla^2 A(\theta) \text{Cov}_{\theta}(\hat{\theta}_n) \nabla^2 A(\theta) \\ \Rightarrow \quad \text{Cov}_{\theta}(\hat{\theta}_n) &\approx \frac{1}{n} (\nabla^2 A(\theta))^{-1} \approx \frac{1}{n} (\nabla^2 A(\hat{\theta}_n))^{-1} \end{aligned}$$

4. Practical way for variance estimation: **bootstrap**

(Central) idea of bootstrap: in order to estimate $\Theta(P)$, one may use $\Theta(P) \approx \Theta(\hat{P})$, with \hat{P} typically being the empirical distribution.

In our case, $\Theta(P) = \text{Variance of MLE based on } y_1, \dots, y_n \sim P$

- if we knew P , we could resample m times from P (say $m = 1,000$):
 - 1) draw $y_1^{(i)}, y_2^{(i)}, \dots, y_n^{(i)} \sim P$,
 - 2) compute the MLE $\hat{\theta}_n^{(i)}$ from $(y_1^{(i)}, \dots, y_n^{(i)})$;
 - 3) compute the sample Variance of $(\hat{\theta}_n^{(1)}, \dots, \hat{\theta}_n^{(m)})$.
- however, we don't know P . Instead, we know $\hat{P} = \text{unif}(\{y_1, \dots, y_n\})$, the empirical distribution of n samples.

- computation of $\theta(\hat{\theta})$:

- 1) draw $y_1^{(i)}, y_2^{(i)}, \dots, y_n^{(i)} \sim \hat{P}$ (i.e. sample from $\{y_1, \dots, y_n\}$ with replacement);
- 2) compute the MLE $\hat{\theta}_n^{(i)}$ from $(y_1^{(i)}, \dots, y_n^{(i)})$;
- 3) compute the sample variance of $(\hat{\theta}_n^{(1)}, \dots, \hat{\theta}_n^{(m)})$.

Some comments on bootstrap:

- bootstrap can be thought of as a general "plug-in" method;
- for example, if $Cov_{\theta}(\nabla A(\hat{\theta}_n)) = \frac{1}{n} \nabla^2 A(\theta)$ for some tractable $\nabla^2 A(\cdot)$, then a simple plug-in method is to use $\frac{1}{n} \nabla^2 A(\theta) \approx \frac{1}{n} \nabla^2 A(\hat{\theta}_n)$;
- however, if the computation of $\nabla^2 A(\cdot)$ is intractable, we can do:
 - nonparametric bootstrap: sample $y_1^{(i)}, \dots, y_n^{(i)} \sim \text{unif}\{y_1, \dots, y_n\}$;
 - parametric bootstrap: sample $y_1^{(i)}, \dots, y_n^{(i)} \sim P_{\theta_n}(y)$.

Example: Fisher's 2x2 table

R. A. Fisher considered the conditional distribution of X_1 given the row & column sums, i.e. (N, r_1, c_1) :

		success	failure	
		X_1	X_2	
treatment	control	π_1	π_2	r_1
		π_3	π_4	r_2
		c_1	c_2	N

$$\begin{aligned}
 p(x_1 | N, r_1, c_1) &\propto \frac{N!}{x_1! (r_1 - x_1)! (c_1 - x_1)! (N - r_1 - c_1 + x_1)!} \pi_1^{x_1} \pi_2^{r_1 - x_1} \pi_3^{c_1 - x_1} \pi_4^{N - r_1 - c_1 + x_1} \\
 &\propto \frac{1}{x_1! (r_1 - x_1)! (c_1 - x_1)! (N - r_1 - c_1 + x_1)!} \underbrace{\left(\frac{\pi_1 \pi_4}{\pi_2 \pi_3} \right)^{x_1}}_{e^{\theta x_1}}
 \end{aligned}$$

log odds: $\theta = \log \left(\frac{\pi_1 \pi_4}{\pi_2 \pi_3} \right)$ ($\theta = 0$: no treatment effect)

log-partition function: $A(\theta) = \log \sum_{x_1} \frac{e^{\theta x_1}}{x_1! (r_1 - x_1)! (c_1 - x_1)! (N - r_1 - c_1 + x_1)!}$

The wldata is on the right.

Numerically one may evaluate:

- $\hat{\theta} = 0.600$
- $A''(\hat{\theta}) = 2.56$

$$\Rightarrow \text{Var}(\hat{\theta}) \approx \frac{1}{A''(\hat{\theta})} = 0.391.$$

		success	failure	
		9	12	21
treatment	π₁		π₂	
	7		17	π₄
control		16	29	45
	π₃			

Question: how would you estimate $\text{Var}(\hat{\theta})$ via bootstrap?

Inference of θ . $H_0: \theta = \theta_0$ vs. $H_1: \theta \neq \theta_0$.

1. 1-D inference ($\theta \in \mathbb{R}$)

- Pearson residual: $\frac{1}{n} \sum_{i=1}^n T(y_i) \xrightarrow{n \rightarrow \infty} N(A'(\theta), \frac{A''(\theta)}{n})$

$$R_p = \frac{\frac{1}{n} \sum_{i=1}^n T(y_i) - A'(\theta_0)}{\sqrt{A''(\theta_0)/n}} \xrightarrow{n \rightarrow \infty} N(0, 1)$$

- Deviance:

$$\begin{aligned} D(\theta_1, \theta_2) &= 2 \mathbb{E}_{\theta_1} \left[\log \frac{P_{\theta_1}(y)}{P_{\theta_2}(y)} \right] \\ &= 2(A(\theta_2) - A(\theta_1) - (\theta_2 - \theta_1) A'(\theta_1)) \geq 0 \end{aligned}$$

Pf of second identity:

$$\begin{aligned} \mathbb{E}_{\theta_1} \left[\log \frac{P_{\theta_1}(y)}{P_{\theta_2}(y)} \right] &= \mathbb{E}_{\theta_1} \left[(\theta_1 - \theta_2) T(y) - A(\theta_1) + A(\theta_2) \right] \\ &= A(\theta_2) - A(\theta_1) - (\theta_2 - \theta_1) A'(\theta_1) \end{aligned}$$

- deviance residual:

$$R_D = \sqrt{n D(\hat{\theta}_n; \theta_0)} \operatorname{sign}\left(\frac{1}{n} \sum_{i=1}^n T(y_i) - A'(\theta_0)\right) \xrightarrow{n \rightarrow \infty} N(0, 1)$$

Intuition: $D(\hat{\theta}_n; \theta_0) = 2(A(\theta_0) - A(\hat{\theta}_n) - (\theta_0 - \hat{\theta}_n)A'(\hat{\theta}_n))$

$$\approx A''(\theta_0)(\hat{\theta}_n - \theta_0)^2$$

$$\approx \frac{1}{n A''(\theta_0)} Z^2 \text{ with } Z \sim N(0, 1)$$

- comparison of Pearson / deviance residuals: see HW.

2. Multivariate inference ($\theta \in \mathbb{R}^d$)

- Wald test: $\sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{n \rightarrow \infty} N(0, \nabla^2 A(\theta_0)^{-1})$ under H_0 .

$$T_{n, \text{Wald}} = n(\hat{\theta}_n - \theta_0)^T \nabla^2 A(\theta_0) (\hat{\theta}_n - \theta_0) \xrightarrow{n \rightarrow \infty} \chi_d^2$$

- Rao's test (score test):

$$\sqrt{n}(\nabla A(\hat{\theta}_n) - \nabla A(\theta_0)) \xrightarrow{n \rightarrow \infty} N(0, \nabla^2 A(\theta_0)) \text{ under } H_0$$

$$\begin{aligned} T_{n, \text{Score}} &= n(\nabla A(\hat{\theta}_n) - \nabla A(\theta_0))^T \nabla^2 A(\theta_0)^{-1} (\nabla A(\hat{\theta}_n) - \nabla A(\theta_0)) \\ &= n\left(\frac{1}{n} \sum_{i=1}^n T(y_i) - \nabla A(\theta_0)\right)^T \nabla^2 A(\theta_0)^{-1} \left(\frac{1}{n} \sum_{i=1}^n T(y_i) - \nabla A(\theta_0)\right) \\ &\xrightarrow{n \rightarrow \infty} \chi_d^2 \end{aligned}$$

- Hoeffding's formula: deviance

$$D(\theta_1; \theta_2) = 2(A(\theta_2) - A(\theta_1) - \langle \theta_2 - \theta_1, \nabla A(\theta_1) \rangle)$$

If $\hat{\theta}_n$ is the MLE based on (y_1, \dots, y_n) , then for every θ ,

$$n D(\hat{\theta}_n; \theta) = 2 \log \frac{P_{\hat{\theta}_n}(y_1, \dots, y_n)}{P_\theta(y_1, \dots, y_n)} \quad (\text{Pf: in class})$$

- likelihood ratio test:

$$T_{n,LRT} = 2 \log \frac{P_{\theta_0}(y_1, \dots, y_n)}{P_{\theta_0}(y_1, \dots, y_n)} = n D(\hat{\theta}_n; \theta_0) \xrightarrow{n \rightarrow \infty} \chi_d^2 \text{ under } H_0$$

(known as Wilks' Theorem)

Intuition: $n D(\hat{\theta}_n; \theta_0) = 2n(A(\theta_0) - A(\hat{\theta}_n) - \langle \theta_0 - \hat{\theta}_n, \nabla A(\hat{\theta}_n) \rangle)$

$$\approx n(\theta_0 - \hat{\theta}_n)^T \nabla^2 A(\theta_0)(\theta_0 - \hat{\theta}_n)^T$$

$$= T_{n,Wald} \xrightarrow{n \rightarrow \infty} \chi_d^2.$$

3. Generalization to $H_0: \theta \in \Theta_0$ with $\dim(\Theta_0) = s < d$

Replace θ_0 by $\hat{\theta}_{0,n} = \underset{\theta \in \Theta_0}{\operatorname{argmax}} \frac{1}{n} \sum_{i=1}^n \log P_\theta(y_i)$, then

$$T_{n,Wald}, T_{n,Score}, T_{n,LRT} \xrightarrow{n \rightarrow \infty} \chi_{d-s}^2.$$