

## Handout 6: Pivotal values, and their asymptotics, for C.I. &amp; H.T.

Lecturer &amp; author: Georgios P. Karagiannis

georgios.karagiannis@durham.ac.uk

References: [1, (Ch. 4)], [2, (Ch. 4)]

**Notation 1.** Let  $X, X_1, X_2, \dots, X_n$  be a sequence of IID random variables (unseen observations) generated from a distribution  $f_\theta$  labeled by a  $d$ -dimensional parameter  $\theta \in \Theta \subset \mathbb{R}^d$ , and admitting PDF  $f(\cdot|\theta)$ .

**Assumption 2.** Assume that the conditions from the Cramer Theorem 19 (Handout 4) are satisfied.

**Notation 3.** Consider that  $\hat{\theta}_n$  is the MLE of  $\theta$ .

- Although the methods below use the MLE  $\hat{\theta}_n$ , in fact, any asymptotic equivalent estimator  $\clubsuit_n$  of the MLE can be used; e.g., the one-step-estimators with moment estimator initial guess.

**Recall that:** asymptotic equivalent  $\clubsuit_n - \hat{\theta}_n \xrightarrow{P} 0 \implies$  asymptotic efficient  $\sqrt{n}(\clubsuit_n - \theta_0) \xrightarrow{D} N(0, \mathcal{I}(\theta_0)^{-1})$

- All the tools below are asymptotic equivalent, as you can imagine. However, for smaller samples it seems that the Likelihood ratio is the most powerful.

**Note 4.** We present 3 pivotal values for HT and CI: the Likelihood ratio, Wald, and Score statistics. Their rational is depicted in Figure 1.

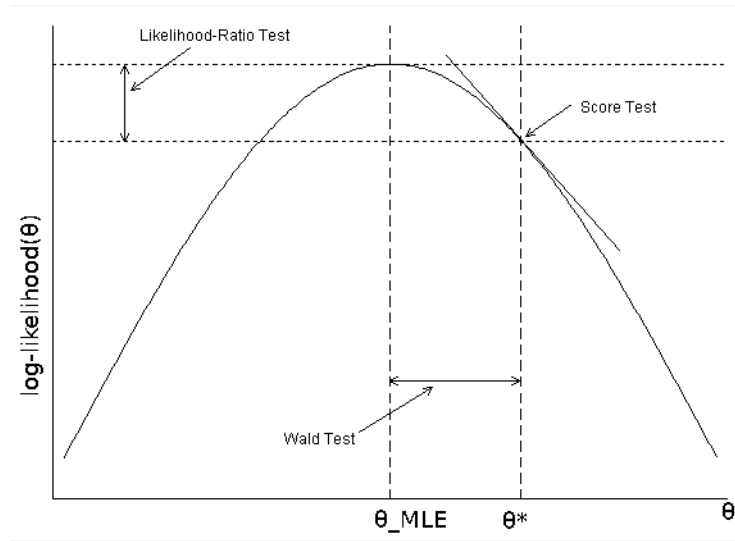


Figure 1: Assume that  $\hat{\theta}$  is the MLE and  $\theta_0$  is the real value of the parameter  $\theta \in \mathbb{R}$ . Comparison, in 1D, of:

Likelihood ratio :  $W_{LR}(\theta_0) = -2(\ell_n(\theta_0) - \ell_n(\hat{\theta}_n))$

Wald :  $W_{Wald}(\theta_0) = n(\theta_0 - \hat{\theta}_n)^2 \mathcal{I}(\theta_0) = -(\theta_0 - \hat{\theta}_n)^2 E(\ddot{\ell}_n(\theta_0))$

Score statistic :  $W_{Score}(\theta_0) = n\dot{\ell}_n(\theta_0)/\mathcal{I}(\theta_0) = -\dot{\ell}_n(\theta_0)/E(\dot{\ell}_n(\theta_0))$

# 1 Likelihood ratio statistic as pivotal statistic

## 1.1 Likelihood ratio pivotal statistic

**Definition 5.** The log likelihood ratio statistic at  $\theta \in \mathbb{R}^d$  is

$$W_{LR}(\theta) = -2(\ell_n(\theta) - \ell_n(\hat{\theta}_n)) \quad (1)$$

**Theorem 6.** Let  $X_1, X_2, \dots, X_n$  be independent random samples generated from a distribution  $f_\theta$  labeled by a  $d$ -dimensional parameter  $\theta \in \Theta \subset \mathbb{R}^d$ , and admitting PDF  $f(\cdot|\theta)$ . Assume the conditions from the Cramer Theorem 19 (Handout 4) are satisfied. Let  $\theta_0$  be the real value of  $\theta$ . Then

$$W_{LR}(\theta_0) = -2(\ell_n(\theta_0) - \ell_n(\hat{\theta}_n)) \xrightarrow{D} \chi_d^2 \quad (2)$$

where  $\hat{\theta}_n$  is the MLE of  $\theta$ .

*Proof.* It is given as an Exercise 31 in Exercise sheet, as well as in the Appendix A.  $\square$

## 1.2 Likelihood ratio hypothesis test

**Proposition 7.** Consider Hypothesis test

$$H_0 : \theta = \theta_* \quad \text{vs.} \quad H_1 : \theta \neq \theta_*$$

The rejection area, at sig. level  $a$ , is

$$RA(X_{1:n}) = \{X_{1:n} : W_{LR}(\theta_*) \geq \chi_{d,1-a}^2\} = \{X_{1:n} : -2(\ell_n(\theta_*) - \ell_n(\hat{\theta}_n)) \geq \chi_{d,1-a}^2\} \quad (3)$$

## 1.3 Likelihood ratio confidence region

**Proposition 8.** The  $(1 - a)$  confidence interval for  $\theta$  is

$$CI(\theta) = \{\theta \in \Theta : W_{LR}(\theta) \leq \chi_{d,1-a}^2\} = \{\theta \in \Theta : -2(\ell_n(\theta) - \ell_n(\hat{\theta}_n)) \leq \chi_{d,1-a}^2\} \quad (4)$$

as produced by inverting the  $RA(x_{1:n})$

## 1.4 Comments

1. Regarding the HT, the comparison relies on the distance of the log-likelihood ratio  $\ell_n(\theta_*)$  and  $\ell_n(\hat{\theta}_n)$ . The larger the distance is, the biggest doubt about the  $H_0$  based on my data. See Figure 1.
2. It is more powerful than the other 2 tests, and hence preferable if it can be practically evaluated. The other 2 were derived possibly because ages ago people did not have computers and wanted to use something computationally cheaper.

## 2 Wald statistics as pivotal statistic

### 2.1 Wald pivotal statistic

**Definition 9.** The Wald statistic, is defined as

$$W_W(\theta_0) = n(\hat{\theta}_n - \theta_0)^T \mathcal{I}(\theta_0)(\hat{\theta}_n - \theta_0) \quad (5)$$

**Definition 10.** Other, more tractable variations of the Wald statistic are

$$W'_W(\theta_0) = n(\hat{\theta}_n - \theta_0)^T \mathcal{I}(\hat{\theta}_n)(\hat{\theta}_n - \theta_0) \quad (6)$$

$$W''_W(\theta_0) = (\hat{\theta}_n - \theta_0)^T \mathcal{J}_n(\hat{\theta}_n)(\hat{\theta}_n - \theta_0) \quad (7)$$

**Proposition 11.** Assume the conditions from the Cramer Theorem 19 (Handout 4) are satisfied. Let  $\theta_0$  be the real value of  $\theta$ . Then  $W_W(\theta_0)$ ,  $W'_W(\theta_0)$ , and  $W''_W(\theta_0)$  are asymptotic equivalent, and

$$W_W(\theta_0) = n(\hat{\theta}_n - \theta_0)^T \mathcal{I}(\theta_0)(\hat{\theta}_n - \theta_0) \xrightarrow{D} \chi_d^2 \quad (8)$$

$$W'_W(\theta_0) = n(\hat{\theta}_n - \theta_0)^T \mathcal{I}(\hat{\theta}_n)(\hat{\theta}_n - \theta_0) \xrightarrow{D} \chi_d^2 \quad (9)$$

$$W''_W(\theta_0) = (\hat{\theta}_n - \theta_0)^T \mathcal{J}_n(\hat{\theta}_n)(\hat{\theta}_n - \theta_0) \xrightarrow{D} \chi_d^2 \quad (10)$$

*Proof.* Asymptotic equivalence is proved by the definition: E.g.,

$$W_W(\theta_0) - W'_W(\theta_0) = \sqrt{n}(\hat{\theta}_n - \theta_0)^T \left[ \mathcal{I}(\theta_0) - \mathcal{I}(\hat{\theta}_n) \right] \sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{p} 0$$

because  $\mathcal{I}(\hat{\theta}_n) \xrightarrow{p} \mathcal{I}(\theta_0)$ . Asymptotic distribution is proved as consequences of Cramer Theorem, Lemma 24 (Handout 4), and Slutsky; E.g. for (8) because  $\sqrt{n}\mathcal{I}(\theta_0)^{1/2}(\hat{\theta}_n - \theta_0) \xrightarrow{D} N(0, I)$  then

$$W_W(\theta_0) = \left[ \sqrt{n}\mathcal{I}(\theta_0)^{1/2}(\hat{\theta}_n - \theta_0) \right]^\top \left[ \sqrt{n}\mathcal{I}(\theta_0)^{1/2}(\hat{\theta}_n - \theta_0) \right] \xrightarrow{D} Z^\top Z$$

where  $Z = (z_1, \dots, z_d) \sim N(0, I)$  and hence  $\sum_{i=1}^d z_i^2 \sim \chi_d^2$ . For 9, and 10, I do the same by using

$$\sqrt{n}\mathcal{I}(\hat{\theta}_n)^{1/2}(\hat{\theta}_n - \theta_0) \xrightarrow{D} N(0, I), \quad \text{and} \quad \mathcal{J}_n(\hat{\theta}_n)^{1/2}(\hat{\theta}_n - \theta_0) \xrightarrow{D} N(0, I)$$

from Lemma 24 (Handout 4), and Slutsky rules.  $\square$

*Remark 12.* (Wald) pivotal statistics in (8, 9, and 10) are asymptotically equivalent for large samples (this is obvious by construction). However, the order of preference is  $W_W(\theta_0)$ ,  $W'_W(\theta_0)$ ,  $W''_W(\theta_0)$  when the sample size is not that large; the proof is out of scope.

### 2.2 Wald hypothesis test

**Proposition 13.** Consider Hypothesis test

$$H_0 : \theta = \theta_* \quad \text{vs.} \quad H_1 : \theta \neq \theta_*$$

The rejection area, at sig. level  $\alpha$ , is

$$RA(X_{1:n}) = \{X_{1:n} : W_{Wald}(\theta_0) \geq \chi_{d,1-\alpha}^2\} = \{X_{1:n} : n(\hat{\theta}_n - \theta_*)^T \mathcal{I}(\theta_*)(\hat{\theta}_n - \theta_*) \geq \chi_{d,1-\alpha}^2\} \quad (11)$$

## 2.3 Wald confidence sets

**Proposition 14.** *The  $(1 - a)$  confidence interval for  $\theta$  is*

$$CI(\theta) = \{\theta \in \Theta : W_{Wald}(\theta) \leq \chi_{d,1-a}^2\} = \{\theta \in \Theta : n(\hat{\theta}_n - \theta)^T \mathcal{I}(\theta)(\hat{\theta}_n - \theta) \leq \chi_{d,1-a}^2\} \quad (12)$$

*produced by inverting the RA( $x_{1:n}$ )*

## 2.4 Comment

1. The Wald pivotal statistics are asymptotically equivalent to the LR one.
2. Regarding the HT, the comparison relies on the distance of the  $\theta_*$  and  $\hat{\theta}_n$ , calibrated by the Information matrix (Fisher or Observed). The larger the distance is, the biggest doubt about the  $H_0$  based on my data. See Figure 1.
3. Wald type of HT, CI are less expensive than the likelihood ratio ones because they require the computation of the expensive likelihood less number of times.

## 3 Score statistics as pivotal statistic

### 3.1 Score statistic

**Definition.** The Score statistic is defined as

$$U(\theta) = \left[ \frac{d}{d\theta} \ell_n(\theta) \right]^T = \left[ \sum_{i=1}^d \underbrace{\frac{d}{d\theta} \log f(X_i|\theta)}_{=\Psi(X_i, \theta)} \right]^T \quad (13)$$

**Proposition 15.** *The asymptotic distribution is*

$$\frac{1}{\sqrt{n}} U(\theta) \xrightarrow{D} N(0, \mathcal{I}(\theta)) \quad (14)$$

*which results as in Example/Proposition 16 (Handout 4).*

**Definition 16.** The following score pivotal statistic is produced from the score statistic:

$$W_{\text{Score}}(\theta) = \frac{1}{n} U(\theta)^T \mathcal{I}(\theta)^{-1} U(\theta) \quad (15)$$

**Definition 17.** Other, more tractable variations of the score pivotal statistic are

$$W'_{\text{Score}}(\theta) = \frac{1}{n} U(\theta)^T \mathcal{I}(\hat{\theta}_n)^{-1} U(\theta) \quad (16)$$

$$W''_{\text{Score}}(\theta) = U(\theta)^T \mathcal{J}_n(\hat{\theta}_n)^{-1} U(\theta) \quad (17)$$

**Proposition 18.** Assume the conditions from the Cramer Theorem 19 (Handout 4) are satisfied. Let  $\theta_0$  be the real value of  $\theta$ . Then  $W_{Score}(\theta_0)$ ,  $W'_{Score}(\theta_0)$ , and  $W''_{Score}(\theta_0)$  are asymptotic equivalent, and

$$W_{Score}(\theta_0) = \frac{1}{n}U(\theta_0)^T \mathcal{I}(\theta_0)^{-1}U(\theta_0) \xrightarrow{D} \chi_d^2 \quad (18)$$

$$W'_{Score}(\theta_0) = \frac{1}{n}U(\theta_0)^T \mathcal{I}(\hat{\theta}_n)^{-1}U(\theta_0) \xrightarrow{D} \chi_d^2 \quad (19)$$

$$W''_{Score}(\theta_0) = U(\theta_0)^T \mathcal{J}_n(\hat{\theta}_n)^{-1}U(\theta_0) \xrightarrow{D} \chi_d^2 \quad (20)$$

*Proof.* The asymptotic equivalent is proven as in Proposition 11. Regarding the asymptotic distributions: by using (14) with Slutsky rules, it can be derived

$$\frac{1}{\sqrt{n}}\mathcal{I}(\theta)^{-1/2}U(\theta) \xrightarrow{D} N(0, I), \quad ; \quad \frac{1}{\sqrt{n}}\mathcal{I}(\hat{\theta}_n)^{-1/2}U(\theta) \xrightarrow{D} N(0, I), \quad \text{and} \quad \mathcal{J}_n(\hat{\theta}_n)^{-1/2}U(\theta) \xrightarrow{D} N(0, I)$$

and hence 18, 19, and (20) can be derived by applied Slutsky rules again.  $\square$

### 3.2 Score hypothesis test

**Proposition 19.** Consider Hypothesis test

$$H_0 : \theta = \theta_* \quad \text{vs.} \quad H_1 : \theta \neq \theta_*$$

Hence the rejection area, at sig. level  $\alpha$ , is

$$RA(X_{1:n}) = \{X_{1:n} : W_{Score}(\theta_0) \geq \chi_{d,1-\alpha}^2\} = \{X_{1:n} : \frac{1}{n}U(\theta_*)^T \mathcal{I}(\theta_*)^{-1}U(\theta_*) \geq \chi_{d,1-\alpha}^2\} \quad (21)$$

### 3.3 Score confidence intervals

**Proposition 20.** The  $(1 - \alpha)$  confidence interval for  $\theta$  is

$$CI(\theta) = \{\theta \in \Theta : W_{Score}(\theta) \leq \chi_{d,1-\alpha}^2\} = \{\theta \in \Theta : \frac{1}{n}U(\theta)^T \mathcal{I}(\theta)^{-1}U(\theta) \leq \chi_{d,1-\alpha}^2\} \quad (22)$$

etc... produced by inverting the  $RA(X_{1:n})$ .

### 3.4 Comments

1. The Score pivotal statistics are asymptotically equivalent to the LR one.
2. Regarding the HT, the comparison relies on the slope of the log-likelihood at  $\theta_*$  (aka the  $U(\theta_*)$ ) calibrated by the curvature (Hessian matrix) at  $\theta_*$ . The larger/steeper the slope, the bigger the distance from the peak (MLE  $\hat{\theta}_n$ ), hence the biggest doubt about the  $H_0$  based on my data. See Figure 1.
3. Score type of HT, CI are less expensive than the likelihood ratio ones because they require the computation of the expensive likelihood less number of times.

4. Score HT and CI are computational convenient, in situations when the practitioner wants to calculate the HT or CI for parameter  $\phi \in \mathbb{R}^k$ , which is a function of parameter  $\theta \in \mathbb{R}^d$  whose Score type HT or CI have already been calculated.

Let  $\phi := \phi(\theta)$ . The score statistic

$$U^*(\phi) = \left[ \frac{d}{d\phi} \ell_n(\theta) \right]^T = \left[ \sum_{i=1}^d \frac{d}{d\phi} f(X_i|\phi) \right]^T$$

is such that

$$U^*(\phi) = \left[ \frac{d}{d\phi} \ell_n(\phi) \right]^T = \left[ \frac{d}{d\theta} \ell_n(\theta) \frac{d\theta}{d\phi} \right]^T = \left[ \frac{d\theta}{d\phi} \right]^T U(\theta)$$

It has expected value

$$E(U^*(\phi)) = E\left(\left[ \frac{d\theta}{d\phi} \right]^T U(\theta)\right) = 0$$

it has variance

$$\text{var}(U^*(\phi)) = \text{var}\left(\left[ \frac{d\theta}{d\phi} \right]^T U(\theta)\right) = \left[ \frac{d\theta}{d\phi} \right]^T \mathcal{I}(\theta) \left[ \frac{d\theta}{d\phi} \right]$$

and hence, by Delta method it has asymptotic distribution

$$\frac{1}{\sqrt{n}} U^*(\phi) \xrightarrow{D} N\left(0, \underbrace{\left[ \frac{d\theta}{d\phi} \right]^T \mathcal{I}(\theta) \left[ \frac{d\theta}{d\phi} \right]}_{=\mathcal{I}^*(\phi)}\right)$$

Hence, one can derive HT and CI as in (21) and (22) by substituting properly <sup>1</sup>. Notice that if the score HT and CI for  $\theta$  are available then the score HT and CI for the transformation  $\phi = \phi(\theta)$  can be computed by avoiding to recompute the expensive likelihood function.

**Example 21.** Let random sample  $x_1, \dots, x_n \stackrel{IID}{\sim} \text{Poi}(\theta)$ ,  $\theta > 0$  with PDF

$$f(x|\theta) = \frac{\theta^x e^{-\theta}}{x!} 1(x > 0)$$

For the hypothesis test

$$H_0 : \theta = \theta_* \quad \text{vs.} \quad H_1 : \theta \neq \theta_*$$

Calculate

1. the log-likelihood ratio RA at  $\alpha$  sig. level
2. the Wald's type RA at  $\alpha$  sig. level (the 3 of them)
3. the Score's type RA at  $\alpha$  sig. level (the 3 of them)

---

<sup>1</sup>please write down the derivation

**Solution.** Ok before that, let's calculate all the quantities required.

$$\begin{aligned}
\log f(x|\theta) &\propto x \log(\theta) - \theta & ;;& & 0 = \dot{\ell}(\theta)|_{\theta=\hat{\theta}} \implies \hat{\theta} = \bar{x} \\
\frac{d}{d\theta} \log f(x|\theta) &= \frac{x}{\theta} - 1 & ;;& & \ddot{\ell}(\theta) = -n\bar{x} \frac{1}{\theta^2} \\
\frac{d^2}{d\theta^2} \log f(x|\theta) &= -\frac{x}{\theta^2} & ;;& & \mathcal{J}_n(\theta) = -\ddot{\ell}(\theta) = n\bar{x} \frac{1}{\theta^2} \\
\mathcal{I}(\theta) &= -E\left(\frac{d^2}{d\theta^2} \log f(x|\theta)\right) = \frac{1}{\theta} & ;;& & \mathcal{J}_n(\hat{\theta}) = \frac{n}{\bar{x}} \\
\ell(\theta) &= n\bar{x} \log(\theta) - n\theta & ;;& & U(\theta) = \dot{\ell}(\theta)^T = n\bar{x} \frac{1}{\theta} - n \\
\dot{\ell}(\theta) &= n\bar{x} \frac{1}{\theta} - n & ;;& & 
\end{aligned}$$

1. It is

$$\begin{aligned}
\text{CI}(\theta) &= \{\theta \in (0, \infty) : -2(\ell_n(\theta) - \ell_n(\hat{\theta}_n)) \leq \chi_{1,1-a}^2\} \\
&= \{\theta \in (0, \infty) : -2((n\bar{x} \log(\frac{\theta}{\bar{x}}) - n(\theta - \bar{x})) \leq \chi_{1,1-a}^2\}
\end{aligned}$$

well, here we do not have the condition  $n\theta - n\bar{x}$  like in the contingency tables ...

2. It is

$$\begin{aligned}
\text{RA}(x_{1:n}) &= \{x_{1:n} : n(\hat{\theta}_n - \theta_*)^T \mathcal{I}(\theta_*)(\hat{\theta}_n - \theta_*) \geq \chi_{1,1-a}^2\} \\
&= \{x_{1:n} : n(\bar{x} - \theta_*)^2 \frac{1}{\theta_*} \geq \chi_{1,1-a}^2\}
\end{aligned}$$

and

$$\begin{aligned}
\text{RA}(x_{1:n}) &= \{x_{1:n} : n(\hat{\theta}_n - \theta_*)^T \mathcal{I}(\hat{\theta}_n)(\hat{\theta}_n - \theta_*) \geq \chi_{1,1-a}^2\} \\
&= \{x_{1:n} : n(\bar{x} - \theta_*)^2 \frac{1}{\bar{x}} \geq \chi_{1,1-a}^2\}
\end{aligned}$$

and

$$\begin{aligned}
\text{RA}(x_{1:n}) &= \{x_{1:n} : (\hat{\theta}_n - \theta_*)^T \mathcal{J}_n(\hat{\theta}_n)(\hat{\theta}_n - \theta_*) \geq \chi_{1,1-a}^2\} \\
&= \{x_{1:n} : (\bar{x} - \theta_*)^2 \frac{n}{\bar{x}} \geq \chi_{1,1-a}^2\}
\end{aligned}$$

where the latter two CI coincide, however this is just a coincidence...

3. It is

$$\begin{aligned}
\text{RA}(x_{1:n}) &= \{x_{1:n} : \frac{1}{n} U(\theta_*)^T \mathcal{I}(\theta_*)^{-1} U(\theta_*) \geq \chi_{1,1-a}^2\} \\
&= \{x_{1:n} : \frac{1}{n} (n\bar{x} \frac{1}{\theta_*} - n)^2 (\frac{1}{\theta_*})^{-1} \geq \chi_{1,1-a}^2\} \\
&= \{x_{1:n} : n(\bar{x} \frac{1}{\theta_*} - 1)^2 \theta_* \geq \chi_{1,1-a}^2\}
\end{aligned}$$

and

$$\begin{aligned} \text{RA}(x_{1:n}) &= \{x_{1:n} : \frac{1}{n}U(\theta_*)^T \mathcal{I}(\hat{\theta}_n)^{-1}U(\theta_*) \geq \chi_{1,1-a}^2\} \\ &= \{x_{1:n} : n(\bar{x}\frac{1}{\theta_*} - 1)^2\hat{\theta}_n \geq \chi_{1,1-a}^2\} \end{aligned}$$

and

$$\begin{aligned} \text{RA}(x_{1:n}) &= \{x_{1:n} : \frac{1}{n}U(\theta_*)^T \mathcal{J}_n(\hat{\theta}_n)^{-1}U(\theta_*) \geq \chi_{1,1-a}^2\} \\ &= \{x_{1:n} : n(\bar{x}\frac{1}{\theta_*} - 1)^2\hat{\theta}_n \geq \chi_{1,1-a}^2\} \end{aligned}$$

where the latter two CI coincide, however it is just a coincidence...

Confidence intervals can be computed by inverting the RA, based on the theory learnt in Concepts in Stats 2 (Term 1).

**Example 22.** Consider an  $M$  way contingency table  $(n_{i,j})$  generated by a Poisson sampling scheme. Consider it is modeled by a log-Linear model with link function

$$\log(\mu) = X^T \beta \quad (23)$$

in the vectorized form, where vector  $\beta \in \mathbb{R}^d$  contains the unknown coefficients. Consider that identifiability constraints have been considered in (23). Show that

1. the asymptotic distribution of the MLEs  $\hat{\beta}_n$  is such that

$$(X \text{diag}(\mu_n) X^T)^{1/2}(\hat{\beta}_n - \beta_0) \xrightarrow{D} N(0, I) \quad (24)$$

where  $\hat{\mu}_n = \exp(X^T \hat{\beta}_n)$ , and  $\beta_0$  is the true value of  $\beta$ .

2. An  $(1 - a)100\%$  asymptotic confidence interval for  $\beta_0$  is

$$\{\beta_0 : (\hat{\beta}_n - \beta_0)^T (X \text{diag}(\hat{\mu}_n) X^T) (\hat{\beta}_n - \beta_0) \leq \chi_{d,1-a}^2\}$$

**Solution.**

1. Since  $\hat{\beta}_n$  is an MLE then it is asymptotically Normal as

$$\mathcal{J}_n(\hat{\theta}_n)^{1/2}(\hat{\beta}_n - \beta_0) \xrightarrow{D} N(0, I) \quad (25)$$

It is

$$\mathcal{J}_n = -\frac{d^2}{d\beta^2} \ell_n(\beta)|_{\beta=\hat{\beta}_n} = -\frac{d^2}{d\beta^2} \sum_{i=1}^n \log(\text{Poi}(n|\mu(\beta)))|_{\beta=\hat{\beta}_n}$$

$$\begin{aligned} \frac{d}{d\beta_j} \ell_n(\mu(\beta)) &= -\sum_i n_i X_{i,j} + \sum_i \exp(\sum_j X_{i,j} \beta_j) X_{i,j} \\ \frac{d^2}{d\beta_j d\beta_k} \ell_n(\mu(\beta)) &= \sum_i \exp(\sum_j X_{i,j} \beta_j) X_{i,j} X_{i,k} \end{aligned}$$



and

$$\mathcal{J}_n = X^T \text{diag}(\hat{\mu}_n) X$$

Therefore from (25)

$$\underbrace{(X^T \text{diag}(\hat{\mu}_n) X)^{1/2} (\hat{\beta}_n - \beta_0)}_{=Z_n} \xrightarrow{D} N(0, I)$$

2. I use the statistic

$$T_n = (\hat{\beta}_n - \beta_0)^T (X^T \text{diag}(\hat{\mu}_n) X) (\hat{\beta}_n - \beta_0) = Z_n^T Z_n = \sum_{i=1}^d Z_{n,i}^2$$

where  $T_n \sim \chi_d^2$  as summation of  $d$  standard normal variables. Then

$$1 - a = P_{\chi_d^2}(T_n < q)$$

where  $q = \chi_{d,1-a}^2$ .

Exercise sheet

Exercise #31, 32, 33, 34, 36

## References

- [1] Robert J Serfling. *Approximation theorems of mathematical statistics*, volume 162. John Wiley & Sons, 2009.
- [2] T. A. Severini. *Likelihood methods in statistics*. Oxford University Press, 2000.

# Appendix

## A Appendix

Sketch of the derivation of the likelihood ratio statistic asymptotic distribution in Theorem 6. For more details see the Exercise sheet.

*Proof.* Let's expand it,

$$\begin{aligned}\ell_n(\theta_0) &= \ell_n(\hat{\theta}_n) + \dot{\ell}_n(\hat{\theta}_n)(\theta_0 - \hat{\theta}_n) + (\theta_0 - \hat{\theta}_n)^T \int_0^1 \int_0^1 u \ddot{\ell}_n(\hat{\theta}_n + uv(\theta_0 - \hat{\theta}_n)) du dv (\theta_0 - \hat{\theta}_n) \\ &= \ell_n(\hat{\theta}_n) + \dot{\ell}_n(\hat{\theta}_n)(\theta_0 - \hat{\theta}_n) + (\theta_0 - \hat{\theta}_n)^T n \int_0^1 \int_0^1 u \frac{1}{n} \ddot{\ell}_n(\hat{\theta}_n + uv(\theta_0 - \hat{\theta}_n)) du dv (\theta_0 - \hat{\theta}_n)\end{aligned}$$

Rearranging the terms

$$-2(\ell_n(\theta_0) - \ell_n(\hat{\theta}_n)) = \underbrace{-\dot{\ell}_n(\hat{\theta}_n)(\theta_0 - \hat{\theta}_n)}_{=0} - n(\theta_0 - \hat{\theta}_n)^T \underbrace{\int_0^1 \int_0^1 u \frac{1}{n} \ddot{\ell}_n(\hat{\theta}_n + uv(\theta_0 - \hat{\theta}_n)) du dv (\theta_0 - \hat{\theta}_n)}_{\xrightarrow{a.s.} -\frac{1}{2}\mathcal{I}(\theta_0)}$$

It is  $\dot{\ell}_n(\hat{\theta}_n) = 0$  because  $\hat{\theta}_n$  is an MLE.

It is

$$\int_0^1 \int_0^1 u \frac{1}{n} \ddot{\ell}_n(\hat{\theta}_n + uv(\theta_0 - \hat{\theta}_n)) du dv \xrightarrow{a.s.} -\frac{1}{2}\mathcal{I}(\theta_0)$$

by using USLLN properly as in the proof of Cramer theorem for MLE asymptotic distribution.  $\square$

So to sum up

$$-2(\ell_n(\theta_0) - \ell_n(\hat{\theta}_n)) \xrightarrow{a.s.} n(\theta_0 - \hat{\theta}_n)^T \mathcal{I}(\theta_0)(\theta_0 - \hat{\theta}_n)$$

which implies that

$$\left[ -2(\ell_n(\theta_0) - \ell_n(\hat{\theta}_n)) \right] - \left[ n(\theta_0 - \hat{\theta}_n)^T \mathcal{I}(\theta_0)(\theta_0 - \hat{\theta}_n) \right] \xrightarrow{p} 0 \quad (26)$$

From Cramer' Theorem I know that

$$\begin{aligned}\sqrt{n}(\hat{\theta}_n - \theta_0) &\xrightarrow{D} N(0, \mathcal{I}(\theta_0)^{-1}) \\ \implies n(\hat{\theta}_n - \theta_0)^T \mathcal{I}(\theta_0)(\hat{\theta}_n - \theta_0) &\xrightarrow{D} \chi_d^2\end{aligned}$$

But  $-2(\ell_n(\theta_0) - \ell_n(\hat{\theta}_n))$  is asymptotic equivalent to  $n(\hat{\theta}_n - \theta_0)^T \mathcal{I}(\theta_0)(\hat{\theta}_n - \theta_0)$  from 26. So by the Slutsky's theorem

$$-2(\ell_n(\theta_0) - \ell_n(\hat{\theta}_n)) \xrightarrow{D} \chi_d^2$$

## B Repeation of results from Handout 4.

**Fact 23.** *Cholesky decomposition: Every symmetric, positive definite matrix  $\Sigma$  can be decomposed into a product of a unique lower triangular matrix  $L$  and its transpose, i.e.  $\Sigma = LL^T$ .*

**Lemma 24.** *(Which will be used as a Proposition later on) Show that given that the assumptions [C.1-C.5] of Theorem 19 are satisfied, and that  $\mathcal{I}(\theta)$  and  $\mathcal{J}_n(\theta)$  are continuous on  $\theta$ , then*

$$\sqrt{n}\mathcal{I}(\theta_0)^{1/2}(\hat{\theta}_n - \theta_0) \xrightarrow{D} N(0, I) \quad (27)$$

$$\sqrt{n}\mathcal{I}(\hat{\theta}_n)^{1/2}(\hat{\theta}_n - \theta_0) \xrightarrow{D} N(0, I) \quad (28)$$

$$\mathcal{J}_n(\hat{\theta}_n)^{1/2}(\hat{\theta}_n - \theta_0) \xrightarrow{D} N(0, I) \quad (29)$$

where  $A^{1/2}$  denotes the lower triangular matrix of the Cholesky decomposition of  $A$ ; i.e.,  $A = A^{1/2}(A^{1/2})^T$ .

**Solution.**

- Eq 27 results from Cramer Theorem, and the properties of covariance matrix.
- Eq. 28 results by using Cramer Theorem and Slutsky theorems. Precisely, because  $\hat{\theta}_n \xrightarrow{a.s.} \theta_0$ , Slutsky implies  $\mathcal{I}(\hat{\theta}_n) \xrightarrow{a.s.} \mathcal{I}(\theta_0)$  which implies  $\mathcal{I}(\hat{\theta}_n)^{1/2}\mathcal{I}(\theta_0)^{-1/2} \xrightarrow{a.s.} I$ . Therefore, by Slutsky

$$\underbrace{\mathcal{I}(\hat{\theta}_n)^{1/2}\mathcal{I}(\theta_0)^{-1/2}\sqrt{n}\mathcal{I}(\theta_0)^{1/2}(\hat{\theta}_n - \theta_0)}_{=\sqrt{n}\mathcal{I}(\hat{\theta}_n)^{1/2}(\hat{\theta}_n - \theta_0)} \xrightarrow{D} \underbrace{I \times N(0, I)}_{=N(0, I)}$$

- Eq. 29 results by using the USLLN and Slutsky theorems. So I just need to show that

$$\frac{1}{n}\mathcal{J}_n(\hat{\theta}_n) \xrightarrow{a.s.} \mathcal{I}(\theta_0)$$

Set  $U(x, \theta) = -\frac{d^2}{d\theta^2} \log(f(x|\theta))$ , and  $\mathcal{I}(\theta) = E(U(x, \theta))$ . Then

$$\begin{aligned} \left| \frac{1}{n} \sum_{i=1}^n \underbrace{\left( -\frac{d^2}{d\theta^2} \log(f(x_i|\hat{\theta}_n)) \right)}_{U(x_i, \hat{\theta}_n)} - \mathcal{I}(\theta_0) \right| &\leq \left| \frac{1}{n} \sum_{i=1}^n U(x_i, \hat{\theta}_n) - \mathcal{I}(\hat{\theta}_n) \right| + |\mathcal{I}(\hat{\theta}_n) - \mathcal{I}(\theta_0)| \quad (30) \\ &\leq \sup_{|\hat{\theta}_n - \theta_0| \leq \delta} \left| \frac{1}{n} \sum_{i=1}^n U(x_i, \theta) - \mathcal{I}(\theta) \right| + |\mathcal{I}(\hat{\theta}_n) - \mathcal{I}(\theta_0)| \quad (31) \end{aligned}$$

The first term converges to zero because the assumptions of the USLLN are satisfied. The second term converges to zero because  $\hat{\theta}_n \xrightarrow{a.s.} \theta_0$  and hence  $\mathcal{I}(\hat{\theta}_n) \xrightarrow{a.s.} \mathcal{I}(\theta_0)$  by using Slutsky theorem.

So by Slutsky  $(\frac{1}{n}\mathcal{J}_n(\hat{\theta}_n))^{1/2}\mathcal{I}(\theta_0)^{-1/2} \xrightarrow{a.s.} I$ , and by Slutsky again

$$\underbrace{\left( \frac{1}{n}\mathcal{J}_n(\hat{\theta}_n) \right)^{1/2}\mathcal{I}(\theta_0)^{-1/2}\mathcal{I}(\theta_0)^{1/2}(\hat{\theta}_n - \theta_0)}_{=\mathcal{J}_n(\hat{\theta}_n)^{1/2}(\hat{\theta}_n - \theta_0)} \xrightarrow{D} \underbrace{I \times N(0, I)}_{=N(0, I)}$$