

LECTURE 12 NOTES

1. The likelihood ratio test. We begin by studying the asymptotic distribution of the log-LR statistic

$$\log \lambda_n(\mathbf{x}) = \sup_{\theta \in \Theta} \sum_{i \in [n]} \ell_{\mathbf{x}_i}(\theta) - \sup_{\theta \in \Theta_0} \sum_{i \in [n]} \ell_{\mathbf{x}_i}(\theta).$$

Once we know the asymptotic distribution of $2 \log \lambda_n(\mathbf{x})$, setting the critical value t equal to its $1 - \alpha$ quantile ensures

$$\mathbf{P}_0(2 \log \lambda_n(\mathbf{x}) > t) \rightarrow 1 - F_0(t) \leq \alpha,$$

where $F_0(t)$ is the CDF of the asymptotic distribution of $2 \log \lambda_n(\mathbf{x})$. Equivalently,

$$p(\mathbf{x}) = 1 - F_0(2 \log \lambda_n(\mathbf{x}))$$

is an asymptotically valid p-value: its asymptotic distribution is $\text{unif}(0, 1)$.

THEOREM 1.1. *Consider testing $H_0 : \theta \in \Theta_0$, where Θ_0 is a q -dimensional subspace of \mathbf{R}^p . Let $\mathbf{x}_i \stackrel{\text{i.i.d.}}{\sim} f_{\theta^*}(x)$ for some $\theta^* \in \text{int}(\Theta_0)$. Assume*

1. ℓ_x is twice-continuously differentiable for any $x \in \mathcal{X}$,
2. $\frac{1}{\sqrt{n}} \sum_{i \in [n]} \nabla \ell_{\mathbf{x}_i}(\theta^*) \xrightarrow{d} \mathcal{N}(0, I(\theta^*))$, where $I(\theta^*) = -\mathbf{E}_{\theta^*}[\nabla^2 \ell_{\mathbf{x}_1}(\theta^*)]$.

If the (unrestricted) MLE and restricted MLE are consistent,

$$2 \log \lambda_n(\mathbf{x}) \xrightarrow{d} \chi_{p-q}^2.$$

PROOF. To keep notation manageable, let

1. $\ell_n(\theta) = \sum_{i \in [n]} \ell_{\mathbf{x}_i}(\theta)$ be the log-likelihood of $\{\mathbf{x}_i\}$,
2. $I_n(\theta) = -\nabla^2 \ell_n(\theta)$.

The proof consists of three steps:

1. show that the restricted MLE $\{\tilde{\theta}_n\}$ is asymptotically normal:

$$(1.1) \quad \sqrt{n}(\tilde{\theta}_n - \theta^*) = \left(\frac{1}{n} P \nabla^2 I_n(\theta^*) P \right)^\dagger \frac{1}{\sqrt{n}} \nabla \ell_n(\theta^*) + o_P(1),$$

where P is the orthogonal projector onto $\text{span}(\Theta_0)$.

2. show that $2 \log \lambda(\mathbf{x})$ is essentially a quadratic function of the score.
3. show that the quadratic function of the score is asymptotically χ^2 .

We first show (1.1). Let $I_n(\theta^*) = -\nabla^2 \ell_n(\theta^*)$ be the observed Fisher information. By the optimality of $\tilde{\theta}_n$,

$$\begin{aligned} 0 &= P \nabla \ell_n(\tilde{\theta}_n) \\ &= P(\nabla \ell_n(\theta^*) - I_n(\theta^*)(\tilde{\theta}_n - \theta^*) + r_n), \end{aligned}$$

which, since $\tilde{\theta}_n - \theta^* \in \text{span}(\Theta_0)$,

$$= P(\nabla \ell_n(\theta^*) - I_n(\theta^*)P(\tilde{\theta}_n - \theta^*) + r_n).$$

Rearranging,

$$\sqrt{n}(\tilde{\theta}_n - \theta^*) = (\frac{1}{n} P I_n(\theta^*) P)^\dagger \frac{1}{\sqrt{n}} (\nabla \ell_n(\theta^*) + r_n),$$

It is possible to show the remainder term is asymptotically negligible by appealing to the consistency of $\tilde{\theta}_n$. Thus the restricted MLE is asymptotically normal:

$$\sqrt{n}(\tilde{\theta}_n - \theta^*) = (\frac{1}{n} P I_n(\theta^*) P)^\dagger \frac{1}{\sqrt{n}} \nabla \ell_n(\theta^*) + o_P(1).$$

Combining the asymptotic normality of the (unrestricted) MLE and restricted MLE, we show that $2 \log \lambda(\mathbf{x})$ is approximately quadratic. By a second-order Taylor expansion of ℓ_n at θ^* ,

$$\begin{aligned} (1.2) \quad \ell_n(\hat{\theta}_n) - \ell_n(\theta^*) &= \nabla \ell_n(\theta^*)^T (\hat{\theta}_n - \theta^*) - \frac{1}{2} (\hat{\theta}_n - \theta^*)^T I_n(\theta^*) (\hat{\theta}_n - \theta^*) + r_n, \end{aligned}$$

for any $\theta \in \Theta$. Once again, it is possible to show that the remainder term is asymptotically negligible by appealing to the consistency of $\hat{\theta}_n$. The first term in (1.2) is

$$\begin{aligned} \nabla \ell_n(\theta^*)^T (\hat{\theta}_n - \theta^*) &= \frac{1}{\sqrt{n}} \nabla \ell_n(\theta^*)^T \sqrt{n}(\hat{\theta}_n - \theta^*), \end{aligned}$$

which, by the asymptotic normality of the MLE,

$$\begin{aligned} &= \frac{1}{\sqrt{n}} \nabla \ell_n(\theta^*)^T (\frac{1}{n} I_n(\theta^*))^{-1} (\frac{1}{\sqrt{n}} \nabla \ell_n(\theta^*) + o_P(1)) \\ &= \frac{1}{\sqrt{n}} \nabla \ell_n(\theta^*)^T (\frac{1}{n} I_n(\theta^*))^{-1} \frac{1}{\sqrt{n}} \nabla \ell_n(\theta^*) + o_P(1). \end{aligned}$$

Similarly, the second term is

$$\begin{aligned} &= -\frac{1}{2} (\hat{\theta}_n - \theta^*)^T I_n(\theta^*) (\hat{\theta}_n - \theta^*) \\ &= -\frac{1}{2} \sqrt{n} (\hat{\theta}_n - \theta^*)^T (\frac{1}{n} I_n(\theta^*)) \sqrt{n} (\hat{\theta}_n - \theta^*) \\ &= -\frac{1}{2} (\frac{1}{\sqrt{n}} \nabla \ell_n(\theta^*) + o_P(1))^T (\frac{1}{n} I_n(\theta^*))^{-1} (\frac{1}{\sqrt{n}} \nabla \ell_n(\theta^*) + o_P(1)) \\ &= -\frac{1}{2} \frac{1}{\sqrt{n}} \nabla \ell_n(\theta^*)^T (\frac{1}{n} I_n(\theta^*))^{-1} \frac{1}{\sqrt{n}} \nabla \ell_n(\theta^*) + o_P(1). \end{aligned}$$

We plug the expressions into (1.2) to obtain

$$(1.3) \quad \ell_n(\hat{\theta}_n) - \ell_n(\theta^*) = \frac{1}{2} \frac{1}{\sqrt{n}} \nabla \ell_n(\theta^*)^T \left(\frac{1}{n} I_n(\theta^*) \right)^{-1} \frac{1}{\sqrt{n}} \nabla \ell_n(\theta^*) + o_P(1).$$

There is a similar expression for $\ell_n(\tilde{\theta}_n) - \ell_n(\theta^*)$, and we skip to the conclusion:

$$(1.4) \quad \begin{aligned} \ell_n(\tilde{\theta}_n) - \ell_n(\theta^*) \\ = \frac{1}{2} \frac{1}{\sqrt{n}} \nabla \ell_n(\theta^*)^T \left(\frac{1}{n} P I_n(\theta^*) P \right)^\dagger \frac{1}{\sqrt{n}} \nabla \ell_n(\theta^*) + o_P(1). \end{aligned}$$

Finally, we subtract (1.4) from (1.3) to deduce $2 \log \lambda_n(\mathbf{x})$ is approximately quadratic:

$$\begin{aligned} 2 \log \lambda_n(\mathbf{x}) &= 2(\ell_n(\hat{\theta}_n) - \ell_n(\tilde{\theta}_n)) \\ &= \frac{1}{\sqrt{n}} \nabla \ell_n(\theta^*)^T \left(\left(\frac{1}{n} I_n(\theta^*) \right)^{-1} - \left(\frac{1}{n} P I_n(\theta^*) P \right)^\dagger \right) \frac{1}{\sqrt{n}} \nabla \ell_n(\theta^*) + o_P(1). \end{aligned}$$

Under the conditions of the Theorem, $\frac{1}{\sqrt{n}} \nabla \ell_n(\theta^*) \xrightarrow{d} \mathcal{N}(0, I(\theta^*))$. By Slutsky's theorem,

$$\begin{aligned} 2 \log \lambda_n(\mathbf{x}) &\xrightarrow{d} \left(I(\theta^*)^{\frac{1}{2}} \mathbf{z} \right)^T \left(I(\theta^*)^{-1} - (P I(\theta^*) P)^\dagger \right) \left(I(\theta^*)^{\frac{1}{2}} \mathbf{z} \right) \\ &= \mathbf{z}^T \underbrace{\left(I_p - I(\theta^*)^{\frac{1}{2}} (P I(\theta^*) P)^\dagger I(\theta^*)^{\frac{1}{2}} \right)}_J \mathbf{z}, \end{aligned}$$

where $\mathbf{z} \sim \mathcal{N}(0, I_p)$.

To complete the proof, we check that $I_p - J$ is a projector onto a q -dimensional subspace of \mathbf{R}^p . Indeed, J is symmetric, and, since $\mathcal{R}(P I(\theta^*) P)$ is $\text{span}(\Theta_0)$,

1. J is equivalently $I(\theta^*)^{\frac{1}{2}} P (P I(\theta^*) P)^\dagger P I(\theta^*)^{\frac{1}{2}}$, which begets $J^2 = J$.
- 2.

$$\begin{aligned} \text{tr}(J) &= \text{tr}\left(I(\theta^*)^{\frac{1}{2}} P (P I(\theta^*) P)^\dagger P I(\theta^*)^{\frac{1}{2}}\right) \\ &= \text{tr}\left((P I(\theta^*) P) (P I(\theta^*) P)^\dagger\right) \\ &= \dim(\text{span}(\Theta_0)), \end{aligned}$$

which implies $\text{tr}(I_p - J) = p - q$.

□

EXAMPLE 1.2. Let $\mathbf{x}_i \stackrel{\text{i.i.d.}}{\sim} \text{bin}(n, p)$. Consider testing $H_0 : p = \frac{1}{2}$ versus $H_1 : p \neq \frac{1}{2}$. We showed that the likelihood ratio is

$$\lambda_n(\mathbf{t}) = \left(\frac{1}{2}\right)^{-n} \prod_{i \in [n]} \binom{t}{n}^t \left(1 - \left(\frac{t}{n}\right)\right)^{n-t},$$

By Theorem 1.1, $2 \log \lambda_n(\mathbf{t}) \xrightarrow{d} \chi_{1-0=1}^2$ under H_0 , and the test rejects H_0 if $2 \log \lambda_n(\mathbf{x})$ exceeds $\chi_{1,\alpha}^2$.

2. The Wald and Rao tests. To evaluate the LR, the investigator must evaluate both the (unrestricted) MLE and the restricted MLE. Sometimes, one is easier to obtain than the other. The Wald and score test statistics only depend on one of the MLE's and are appropriate in such cases.

2.1. The Wald test.

THEOREM 2.1. *Under the conditions of Theorem 1.1, the Wald statistic is*

$$(2.1) \quad \hat{\theta}_n^T ((I_p - P)I_n(\hat{\theta}_n)(I_p - P))\hat{\theta}_n,$$

where $\hat{\theta}_n$ is the MLE on Θ and P is the projector onto $\text{span}(\Theta_0)$, is asymptotically χ_{p-q}^2 .

PROOF. The Wald statistic is equivalently

$$\begin{aligned} & (\hat{\theta}_n - \theta_0)^T ((I_p - P)I_n(\hat{\theta}_n)(I_p - P))(\hat{\theta}_n - \theta_0) \\ &= \sqrt{n}(\hat{\theta}_n - \theta_0)^T \left(\frac{1}{n}(I_p - P)I_n(\hat{\theta}_n)(I_p - P) \right) \sqrt{n}(\hat{\theta}_n - \theta_0), \end{aligned}$$

which, by the asymptotic normality of $\hat{\theta}_n$, the LLN, and Slutsky's theorem,

$$\xrightarrow{d} \mathbf{z}^T I(\theta^*)^{-\frac{1}{2}} ((I_p - P)I(\theta^*)(I_p - P))I(\theta^*)^{-\frac{1}{2}} \mathbf{z},$$

where $\mathbf{z} \sim \mathcal{N}(0, I_p)$. It is possible to adapt the third part of the proof of Theorem 1.1 to show that

$$I(\theta^*)^{-\frac{1}{2}} ((I_p - P)I(\theta^*)(I_p - P))I(\theta^*)^{-\frac{1}{2}}$$

is a projector onto a $p - q$ -dimensional subspace of \mathbf{R}^p , which implies the Wald statistic has the stated asymptotic distribution. \square

If $\text{span}(\Theta_0)$ is the kernel of some $R \in \mathbf{R}^{r \times p}$ (i.e. $H_0 : R\theta = 0$) that has full row rank,

$$I_p - P = R^T (R^T)^\dagger = R^T (RR^T)^{-1} R,$$

and the Wald test statistic is

$$\begin{aligned} & (R\hat{\theta}_n)^T ((R^T)^\dagger I_n(\hat{\theta}_n)R^\dagger)^{-1} (R\hat{\theta}_n) \\ &= (R\hat{\theta}_n)^T (R^T I_n(\hat{\theta}_n)^{-1} R)^{-1} (R\hat{\theta}_n) \\ &= \|R\hat{\theta}_n\|_{(R^T I_n(\hat{\theta}_n)^{-1} R)^{-1}}^2. \end{aligned}$$

We see that the Wald test statistic is a measure of the size of the infeasibility of $\hat{\theta}_n$. Under H_1 , the infeasibility is likely large. Thus we reject H_0 if the Wald test statistic is large.

2.2. *The Rao test.* The Rao test is the counterpart to the Wald test: it only depends on the restricted MLE.

THEOREM 2.2. *Under the conditions of Theorem 1.1, the Rao statistic is*

$$\nabla \ell_n(\tilde{\theta}_n)^T I_n(\tilde{\theta}_n)^{-1} \nabla \ell_n(\tilde{\theta}_n),$$

where $\tilde{\theta}_n$ is the MLE on Θ_0 , is asymptotically χ^2_{p-q} .

PROOF. By the optimality of $\tilde{\theta}_n$, $\nabla \ell_n(\tilde{\theta}_n) \in \text{span}(\Theta_0)^\perp$. Thus the Rao statistic is equivalently

$$\frac{1}{\sqrt{n}} \nabla \ell_n(\tilde{\theta}_n)^T (I_p - P) \left(\frac{1}{n} I_n(\tilde{\theta}_n) \right)^{-1} (I_p - P) \frac{1}{\sqrt{n}} \nabla \ell_n(\tilde{\theta}_n),$$

which, by the CLT, the LLN, and Slutsky's theorem,

$$\xrightarrow{d} \mathbf{z}^T I(\theta^*)^{\frac{1}{2}} (I_p - P) I(\theta^*)^{-1} (I_p - P) I(\theta^*)^{\frac{1}{2}} \mathbf{z}.$$

where $\mathbf{z} \sim \mathcal{N}(0, I_p)$. It is possible to adapt the third part of the proof of Theorem 1.1 to show that

$$I(\theta^*)^{\frac{1}{2}} (I_p - P) I(\theta^*)^{-1} (I_p - P) I(\theta^*)^{\frac{1}{2}}$$

is a projector onto a $p - q$ -dimensional subspace of \mathbf{R}^p , which implies the Rao statistic has the stated asymptotic distribution. \square

Under the alternative, the restricted MLE $\tilde{\theta}_n$ is far from the true parameter. Thus $\nabla \ell_n(\tilde{\theta}_n)$ is likely large, and we should reject the null when the Rao statistic is large.

The geometric intuition behind Theorems 1.1, 2.1, 2.2 is worth mentioning. It is most obvious in the form of the score test statistic. By the optimality of $\tilde{\theta}_n$, $\nabla \ell_n(\tilde{\theta}_n)$ falls in the orthocomplement of $\text{span}(\Theta_0)$. Thus it is supported on a $p - q$ -dimensional subspace of \mathbf{R}^p . By the CLT, $\nabla \ell_n(\tilde{\theta}_n)$ is asymptotically normal and quadratic forms of normals are χ^2 , which is the form of the score test statistic.

We remark that the conclusions of Theorems 1.1, 2.1, 2.2 are valid under more general conditions. A practical generalization allows Θ_0 to be a submanifold of \mathbf{R}^p , which includes sets of the form

$$\Theta_0 = \{\theta \in \Theta : h(\theta) = 0\},$$

where $h : \Theta \rightarrow \mathbf{R}^r$ is continuously-differentiable and $\nabla h(\theta)$ is surjective for any $\theta \in \text{int}(\Theta)$. The LR and Rao statistics for testing $H_0 : h(\theta) = 0$ are unchanged; the Wald statistic is

$$h(\hat{\theta}_n)^T (\nabla h(\hat{\theta}_n)^T I_n(\hat{\theta}_n)^{-1} \nabla h(\hat{\theta}_n))^{-1} h(\hat{\theta}_n).$$

Under conditions that ensure the MLE is asymptotically normal, the three statistics are asymptotically χ_r^2 under H_0 .

EXAMPLE 2.3. Let $\mathbf{x} \stackrel{\text{i.i.d.}}{\sim} \text{Cat}(p)$, where p_1, p_2, p_3 are positive and sum to one. We wish to test

$$H_0 : 2\sqrt{p_1}\sqrt{p_2} - p_3 = 0.$$

If p_1, p_2, p_3 are the frequencies of homozygous-dominant, homozygous-recessive, and heterozygous genotypes in a population, H_0 asserts the frequencies are in Hardy-Weinberg equilibrium.

Since $\Theta_0 := \{p \in [0, 1]^3 : 2\sqrt{p_1}\sqrt{p_2} - p_3 = 0\}$ is defined by a non-linear constraint, it is hard to optimize on Θ_0 . However, the (unrestricted) MLE of p_1, p_2, p_3 are easy to obtain: $\hat{p}_j = \frac{1}{n} \sum_{i \in [n]} \mathbf{1}\{\mathbf{x}_i = j\}$. We are led to consider the Wald test.

Since the parameter space

$$\{(p_1, p_2, p_3) \in [0, 1]^3 : p_1 + p_2 + p_3 = 1\}$$

is an affine subset of \mathbf{R}^3 , we cannot directly appeal to Theorem 2.1 to deduce the asymptotic distribution of the Wald test statistic. However, if we reparametrize the model in terms of p_1 and p_2 , the (reduced) parameter space is a full-dimensional set in \mathbf{R}^2 , and Theorem 2.1 is applicable.

After re-parametrization, the hypothesis is equivalently

$$H_0 : 2\sqrt{p_1}\sqrt{p_2} - (1 - p_1 - p_2) = 0$$

Let $h(p) = 2\sqrt{p_1}\sqrt{p_2} - (1 - p_1 - p_2)$. It is possible to show the Wald test statistic is

$$\mathbf{w}_n = h(\hat{p}_n)^T (\nabla h(\hat{p}_n)^T I_n(\hat{p}_n)^{-1} \nabla h(\hat{p}_n))^{-1} h(\hat{p}_n),$$

where $\nabla h(p)^T = \left[\frac{\sqrt{p_2}}{\sqrt{p_1}} + 1 \quad \frac{\sqrt{p_1}}{\sqrt{p_2}} + 1 \right]$ and $I_n(p) = n(\text{diag}(p) - pp^T)^{-1}$.

Since Θ_0 is a 1-dimensional curve in \mathbf{R}^2 , $\mathbf{w}_n \xrightarrow{d} \chi_{1=2-1}^2$ by Theorem 2.1. Thus we reject H_0 if \mathbf{w}_n exceeds $\chi_{1,\alpha}^2$.

3. GMM analogues. To wrap up, we describe the GMM analogues of the LR, Wald, and Rao tests for testing $H_0 : h(\theta) = 0$.

1. The analogue of the LR statistic is the *distance metric statistic*

$$2(Q_n(\tilde{\theta}_n) - Q_n(\hat{\theta}_n)),$$

where $\tilde{\theta}_n$ and $\hat{\theta}_n$ are the restricted and unrestricted GMM estimators.

2. The analogue of the Wald test statistic is

$$h(\hat{\theta}_n)^T (\nabla h(\hat{\theta}_n)^T (n\hat{J}_n)^{-1} \nabla h(\hat{\theta}_n))^{-1} h(\hat{\theta}_n),$$

where \hat{J}_n is a consistent estimator of the asymptotic variance of

$$\sqrt{n} \nabla Q_n(\theta).$$

3. The analogue of the Rao test statistic is

$$n \nabla Q_n(\tilde{\theta}_n)^T \hat{J}_n^{-1} \nabla Q_n(\tilde{\theta}_n),$$

where \hat{J}_n is some consistent estimator of the asymptotic variance of $\sqrt{n} \nabla Q_n(\theta_0)$. In econometrics, where GMM estimators are more popular than the MLE, the Rao test is usually referred to as the *score test* or *Lagrange multiplier test*.

Under conditions that ensure the restricted and unrestricted GMM estimators are asymptotically normal, the three test statistics are asymptotically χ_r^2 .