

Lecture 5: Statistical Inference

Zhentaο Shi

November 4, 2018

Notation: \mathbf{X} denotes a random variable or random vector. \mathbf{x} is its realization.

1 Hypothesis Testing

- A *hypothesis* is a statement about the parameter space Θ .
- The *null hypothesis* Θ_0 is a subset of Θ of interest, ideally suggested by scientific theory.
- The *alternative hypothesis* $\Theta_1 = \Theta \setminus \Theta_0$ is the complement of Θ_0 .
- *Hypothesis testing* is a decision whether to accept the null hypothesis or to reject it according to the observed evidence.
- If Θ_0 is a singleton, we call it a *simple hypothesis*; otherwise we call it a *composite hypothesis*.
- A *test function* is a mapping

$$\phi : \mathcal{X}^n \mapsto \{0, 1\},$$

where \mathcal{X} is the sample space. We accept the null hypothesis if $\phi(\mathbf{x}) = 0$, or reject it if $\phi(\mathbf{x}) = 1$.

- The *acceptance region* is defined as $A_\phi = \{\mathbf{x} \in \mathcal{X}^n : \phi(\mathbf{x}) = 0\}$, and the *rejection region* is $R_\phi = \{\mathbf{x} \in \mathcal{X}^n : \phi(\mathbf{x}) = 1\}$.
- The *power function* of the test ϕ is

$$\beta_\phi(\theta) = P_\theta(\phi(\mathbf{X}) = 1) = E_\theta(\phi(\mathbf{X})).$$

The power function measures, at a given point θ , the probability that the test function rejects the null.

- The *power* of ϕ at θ for some $\theta \in \Theta_1$ is defined as the value of $\beta_\phi(\theta)$. The *size* of the test ϕ is defined as $\alpha = \sup_{\theta \in \Theta_0} \beta_\phi(\theta)$. Notice that the definition of power depends on a θ in the alternative, whereas that of size is independent of θ as it takes the supremum over the set of null Θ_0 .
- The *level* of the test ϕ is a value $\alpha \in (0, 1)$ such that $\alpha \geq \sup_{\theta \in \Theta_0} \beta_\phi(\theta)$, which is often used when it is difficult to attain the exact supremum. A test of size α is also of level α ; while a test of level α must have size smaller or equal to α .

decision	reject H_0	reject H_0
-----	-----	-----
H_0 true	correct	Type I error
H_0 false	Type II error	correct

- size = $P(\text{reject } H_0 \text{ when } H_0 \text{ true})$
- power = $P(\text{reject } H_0 \text{ when } H_0 \text{ false})$
- The probability of committing Type I error is $\beta_\phi(\theta)$ for some $\theta \in \Theta_0$.
- The probability of committing Type II error is $1 - \beta_\phi(\theta)$ for $\theta \in \Theta_1$.

The philosophy on the hypothesis testing has been debated for centuries. At present the prevailing framework in statistics textbooks is the frequentist perspective. A frequentist views the parameter as a fixed constant, and they keep a conservative attitude about the Type I error. Only if overwhelming evidence is demonstrated should a researcher reject the null. Under the philosophy of protecting the null hypothesis, a desirable test should have a small level. Conventionally we take $\alpha = 0.01, 0.05$ or 0.1 . There can be many tests of the correct size.

Example A trivial test function, $\phi(\mathbf{X}) = 1\{0 \leq U \leq \alpha\}$, where U is a random variable from a uniform distribution on $[0, 1]$, has correct size α but no power. Another trivial test function $\phi(\mathbf{X}) = 1$ has the biggest power but useless size.

Usually, we design a test by proposing a test statistic $T_n : \mathcal{X}^n \mapsto \mathbb{R}^+$ and a critical value $c_{1-\alpha}$. Given T_n and $c_{1-\alpha}$, we write the test function as

$$\phi(\mathbf{X}) = 1\{T_n(\mathbf{X}) > c_{1-\alpha}\}.$$

To ensure such a $\phi(\mathbf{x})$ has correct size, we figure out the distribution of T_n under the null hypothesis (called the *null distribution*), and choose a critical value $c_{1-\alpha}$ according to the null distribution and the desirable size or level α .

The concept of *level* is useful if we do not have sufficient information to derive the exact size of a test.

Example If $(X_{1i}, X_{2i})_{i=1}^n$ are randomly drawn from some unknown joint distribution, but we know the marginal distribution is $X_{ji} \sim N(\theta_j, 1)$, for $j = 1, 2$. In order to test the joint hypothesis $\theta_1 = \theta_2 = 0$, we can construct a test function

$$\phi(\mathbf{X}_1, \mathbf{X}_2) = 1\{\{\sqrt{n}|\bar{X}_1| \geq c_{1-\alpha/4}\} \cup \{\sqrt{n}|\bar{X}_2| \geq c_{1-\alpha/4}\}\},$$

where $c_{1-\alpha/4}$ is the $(1 - \alpha/4)$ -th quantile of the standard normal distribution. The level of this test is

$$\begin{aligned} P_{\theta_1=\theta_2=0}(\phi(\mathbf{X}_1, \mathbf{X}_2)) &\leq P_{\theta_1=0}(\sqrt{n}|\bar{X}_1| \geq c_{1-\alpha/4}) + P_{\theta_2=0}(\sqrt{n}|\bar{X}_2| \geq c_{1-\alpha/4}) \\ &= \alpha/2 + \alpha/2 = \alpha. \end{aligned}$$

where the inequality follows by the *Bonferroni inequality* $P(A \cup B) \leq P(A) + P(B)$. Therefore, the level of $\phi(\mathbf{X}_1, \mathbf{X}_2)$ is α , but the exact size is unknown without the knowledge of the joint distribution. (Even if we know the correlation of X_{1i} and X_{2i} , putting two marginally normal distributions together does not make a jointly normal vector in general.)

Denote the class of test functions of level α as $\Psi_\alpha = \{\phi : \sup_{\theta \in \Theta_0} \beta_\phi(\theta) \leq \alpha\}$. A *uniformly most powerful test* $\phi^* \in \Psi_\alpha$ is a test function such that, for every $\phi \in \Psi_\alpha$,

$$\beta_{\phi^*}(\theta) \geq \beta_\phi(\theta)$$

uniformly over $\theta \in \Theta_1$.

Example Suppose a random sample of size 6 is generated from

$$(X_1, \dots, X_6) \sim \text{i.i.d.} N(\theta, 1),$$

where θ is unknown. We want to infer the population mean of the normal distribution. The null hypothesis is $H_0: \theta \leq 0$ and the alternative is $H_1: \theta > 0$. All tests in

$$\Psi = \left\{ 1 \left\{ \bar{X} \geq c/\sqrt{6} \right\} : c \geq 1.64 \right\}$$

has the correct level. Since $\bar{X} = N(\theta, 1/6)$, the power function for those in Ψ is

$$\beta_\phi(\theta) = P\left(\bar{X} \geq \frac{c}{\sqrt{6}}\right) = P\left(\sqrt{6}(\bar{X} - \theta) \geq c - \sqrt{6}\theta\right) = 1 - \Phi\left(c - \sqrt{6}\theta\right)$$

where Φ is the cdf of standard normal. It is clear that $\beta_\phi(\theta)$ is monotonically decreasing in c . Thus the test function

$$\phi(\mathbf{X}) = 1 \left\{ \bar{X} \geq 1.64/\sqrt{6} \right\}$$

is the most powerful test in Ψ , as $c = 1.64$ is the lower bound that Ψ allows.

Another commonly used indicator in hypothesis testing is p -value:

$$\sup_{\theta \in \Theta_0} P_\theta(T_n(\mathbf{x}) \leq T_n(\mathbf{X})).$$

In the above expression, $T_n(\mathbf{x})$ is the realized value of the test statistic T_n , while $T_n(\mathbf{X})$ is the random variable generated by \mathbf{X} under the null $\theta \in \Theta_0$. The interpretation of the p -value is tricky. p -value is the probability that we observe $T_n(\mathbf{X})$ being greater than the realized $T_n(\mathbf{x})$ if the null hypothesis is true. p -value is *not* the probability that the null hypothesis is true. Under the frequentist perspective, the null hypothesis is either true or false, with certainty. The randomness of a test comes only from sampling, not from the hypothesis.

It measures whether the data is consistent with the null hypothesis, or whether the evidence from the data is compatible with the null hypothesis. p -value is closely related to the corresponding test. When p -value is smaller than the specified test size α , the test rejects the null hypothesis.

2 Confidence Interval

An *interval estimate* is a function $C : \mathcal{X}^n \mapsto \{\Theta' : \Theta' \subseteq \Theta\}$ that maps a point in the sample space to a subset of the parameter space. The *coverage probability* of an *interval estimator* $C(\mathbf{X})$ is defined as $P_\theta(\theta \in C(\mathbf{X}))$. The coverage probability is the frequency that the interval estimator captures the true parameter that generates the sample (From the frequentist perspective, the parameter is fixed while the region is random). It is *not* the probability that θ is inside the given region (From the Bayesian perspective, the parameter is random while the region is fixed conditional on \mathbf{X} .)

Suppose a random sample of size 6 is generated from

$$(X_1, \dots, X_6) \sim \text{i.i.d.} N(\theta, 1).$$

Find the coverage probability of the random interval

$$\left[\bar{X} - 1.96/\sqrt{6}, \bar{X} + 1.96/\sqrt{6} \right].$$

Hypothesis testing and confidence interval are closely related. Sometimes it is difficult to directly construct the confidence interval, but easier to test a hypothesis. One way to construct confidence interval is by *inverting a corresponding test*. Suppose ϕ is a test of size α . If $C(\mathbf{X})$ is constructed as

$$C(\mathbf{x}) = \{\theta \in \Theta : \phi_\theta(\mathbf{x}) = 0\},$$

then its coverage probability

$$P_\theta(\theta \in C(\mathbf{X})) = 1 - P_\theta(\phi_\theta(\mathbf{X}) = 1) = 1 - \alpha.$$

3 Bayesian Credible Set

The Bayesian framework offers a coherent and natural language for statistical decision. However, the major criticism against Bayesian statistics is the arbitrariness of the choice of the prior.

In the Bayesian framework, both the data \mathbf{X}_n and the parameter θ are random variables. Before she observes the data, she holds a *prior distribution* π about θ . After observing the data, she updates the prior distribution to a *posterior distribution* $p(\theta|\mathbf{X}_n)$. The *Bayes Theorem* connects the prior and the posterior as

$$p(\theta|\mathbf{X}_n) \propto f(\mathbf{X}_n|\theta)\pi(\theta)$$

where $f(\mathbf{X}_n|\theta)$ is the likelihood function.

Here is a classical example to illustrate the Bayesian approach of statistical inference. Suppose we have an iid sample (X_1, \dots, X_n) drawn from a normal distribution with unknown θ and known σ . For a researcher with a prior distribution $\theta \sim N(\theta_0, \sigma_0^2)$, her posterior distribution is, by some routine calculation, the posterior is also a normal distribution

$$p(\theta|\mathbf{x}) \sim N(\tilde{\theta}, \tilde{\sigma}^2),$$

where $\tilde{\theta} = \frac{\sigma^2}{n\sigma_0^2 + \sigma^2}\theta_0 + \frac{n\sigma_0^2}{n\sigma_0^2 + \sigma^2}\bar{x}$ and $\tilde{\sigma}^2 = \frac{\sigma_0^2\sigma^2}{n\sigma_0^2 + \sigma^2}$. Thus the Bayesian credible set is

$$(\tilde{\theta} - z_{1-\alpha/2} \cdot \tilde{\sigma}, \tilde{\theta} + z_{1-\alpha/2} \cdot \tilde{\sigma}).$$

This posterior distribution depends on θ_0 and σ_0^2 from the posterior. When the sample size is sufficiently large the posterior can be approximated by $N(\bar{x}, \sigma^2/n)$, where the prior information is overwhelmed by the information accumulated from the data.

In contrast, a frequentist estimates $\hat{\theta} = \bar{x} \sim N(\theta, \sigma^2/n)$. Her confidence interval is

$$(\bar{x} - z_{1-\alpha/2} \cdot \sigma / \sqrt{n}, \bar{x} + z_{1-\alpha/2} \cdot \sigma / \sqrt{n}).$$

4 Application in OLS

4.1 Wald Test

Suppose the OLS estimator $\hat{\beta}$ is asymptotic normal, i.e.

$$\sqrt{n}(\hat{\beta} - \beta) \xrightarrow{d} N(0, \Omega)$$

where Ω is a $K \times K$ positive definite covariance matrix and R is a $q \times K$ constant matrix, then $R\sqrt{n}(\hat{\beta} - \beta) \xrightarrow{d} N(0, R\Omega R')$. Moreover, if $\text{rank}(R) = q$, then

$$n(\hat{\beta} - \beta)' R' (R\Omega R')^{-1} R (\hat{\beta} - \beta) \xrightarrow{d} \chi_q^2.$$

Now we intend to test the null hypothesis $R\beta = r$. Under the null, the Wald statistic

$$W_n = n(R\hat{\beta} - r)' (R\hat{\Omega}R')^{-1} (R\hat{\beta} - r) \xrightarrow{d} \chi_q^2$$

where $\hat{\Omega}$ is a consistent estimator of Ω .

Example (Single test) In a linear regression

$$y = x_i' \beta + e_i = \sum_{k=1}^5 \beta_k x_{ik} + e_i.$$

$$E[e_i x_i] = \mathbf{0}_5,$$

where y is wage and

$$x = (\text{edu}, \text{age}, \text{experience}, \text{experience}^2, 1)'$$

To test whether *education* affects *wage*, we specify the null hypothesis $\beta_1 = 0$. Let $R = (1, 0, 0, 0, 0)$.

$$\sqrt{n}\hat{\beta}_1 = \sqrt{n}(\hat{\beta}_1 - \beta_1) = \sqrt{n}R(\hat{\beta} - \beta) \xrightarrow{d} N(0, R\Omega R') \sim N(0, \Omega_{11}),$$

where Ω_{11} is the $(1, 1)$ (scalar) element of Ω . Therefore,

$$\sqrt{n} \frac{\hat{\beta}_1}{\hat{\Omega}_{11}^{1/2}} = \sqrt{\frac{\Omega_{11}}{\hat{\Omega}_{11}}} \sqrt{n} \frac{\hat{\beta}_1}{\Omega_{11}^{1/2}}$$

If $\hat{\Omega} \xrightarrow{p} \Omega$, then $(\Omega_{11}/\hat{\Omega}_{11})^{1/2} \xrightarrow{p} 1$ by the continuous mapping theorem. As $\sqrt{n}\hat{\beta}_1/\Omega_{11}^{1/2} \xrightarrow{d} N(0, 1)$, we conclude $\sqrt{n}\hat{\beta}_1/\hat{\Omega}_{11}^{1/2} \xrightarrow{d} N(0, 1)$.

The above example is a test about a single coefficient, and the test statistic is essentially a t -statistic. The following example gives a test about a joint hypothesis.

Example (Joint test) We want to simultaneously test $\beta_1 = 1$ and $\beta_3 + \beta_4 = 2$ in the above example. The null hypothesis can be expressed in the general form $R\beta = r$, where the restriction matrix R is

$$R = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 0 \end{pmatrix}$$

and $r = (1, 2)'$. Once we figure out R , it is routine to construct the test.

These two examples are linear restrictions. In order to test a nonlinear regression, we need the so-called *delta method*.

Delta method If $\sqrt{n}(\hat{\theta} - \theta_0) \xrightarrow{d} N(0, \Omega_{K \times K})$, and $f : \mathbb{R}^K \mapsto \mathbb{R}^q$ is a continuously differentiable function for some $q \leq K$, then

$$\sqrt{n}(f(\hat{\theta}) - f(\theta_0)) \xrightarrow{d} N\left(0, \frac{\partial f}{\partial \theta}(\theta_0) \Omega \frac{\partial f}{\partial \theta}(\theta_0)'\right).$$

This result can be easily shown by a mean-value expansion

$$f(\hat{\theta}) - f(\theta_0) = \frac{\partial f(\tilde{\theta})}{\partial \theta} (\hat{\theta} - \theta_0)$$

where $\tilde{\theta}$ lies on the line segment connecting $\hat{\theta}$ and θ_0 . Multiply both sides by \sqrt{n} and notice $\tilde{\theta} \xrightarrow{p} \theta_0$, by Slutsky theorem we have $\sqrt{n}(f(\hat{\theta}) - f(\theta_0)) \xrightarrow{d} \frac{\partial f}{\partial \theta}(\theta_0) N(0, \Omega)$.

In the example of linear regression, the optimal experience level can be found by setting the first order condition with respect to experience to set, $\beta_3 + 2\beta_4 \text{experience}^* = 0$. We test the hypothesis that the optimal experience level is 20 years; in other words,

$$\text{experience}^* = -\frac{\beta_3}{2\beta_4} = 20.$$

This is a nonlinear hypothesis. If $q \leq K$ where q is the number of restrictions, we have

$$n \left(f(\hat{\theta}) - f(\theta_0) \right)' \left(\frac{\partial f}{\partial \theta}(\theta_0) \Omega \frac{\partial f}{\partial \theta}(\theta_0)' \right)^{-1} \left(f(\hat{\theta}) - f(\theta_0) \right) \xrightarrow{d} \chi_q^2,$$

where in this example, $\theta = \beta$, $f(\beta) = -\beta_3 / (2\beta_4)$. The gradient

$$\frac{\partial f}{\partial \beta}(\beta) = \left(0, 0, -\frac{1}{2\beta_4}, \frac{\beta_3}{2\beta_4^2} \right)$$

Since $\hat{\beta} \xrightarrow{p} \beta_0$, by the continuous mapping theorem, if $\beta_{0,4} \neq 0$, we have $\frac{\partial f}{\partial \beta}(\hat{\beta}) \xrightarrow{p} \frac{\partial f}{\partial \beta}(\beta_0)$. Therefore, the (nonlinear) Wald test is

$$W_n = n \left(f(\hat{\beta}) - 20 \right)' \left(\frac{\partial f}{\partial \beta}(\hat{\beta}) \hat{\Omega} \frac{\partial f}{\partial \beta}(\hat{\beta})' \right)^{-1} \left(f(\hat{\beta}) - 20 \right) \xrightarrow{d} \chi_1^2.$$

This is a valid test with correct asymptotic size.

However, we can equivalently state the null hypothesis as $\beta_3 + 40\beta_4 = 0$ and we can construct a Wald statistic accordingly. In general, a linear hypothesis is preferred to a nonlinear one, due to the approximation error in the delta method under the null and more importantly the invalidity of the Taylor expansion under the alternative. It also highlights the problem of Wald test being *variant* for re-parametrization.

4.2 Lagrangian Multiplier Test

Restricted least square

$$\min_{\beta} (y - X\beta)' (y - X\beta) \text{ s.t. } R\beta = r.$$

Turn it into an unrestricted problem

$$L(\beta, \lambda) = \frac{1}{2n} (y - X\beta)' (y - X\beta) + \lambda' (R\beta - r).$$

The first-order condition

$$\begin{aligned} \frac{\partial}{\partial \beta} L &= -\frac{1}{n} X' (y - X\tilde{\beta}) + \tilde{\lambda} R = -\frac{1}{n} X' e + \frac{1}{n} X' X (\tilde{\beta} - \beta^*) + R' \tilde{\lambda} = 0. \\ \frac{\partial}{\partial \lambda} L &= R\tilde{\beta} - r = R(\tilde{\beta} - \beta^*) = 0 \end{aligned}$$

Combine these two equations into a linear system,

$$\begin{pmatrix} \hat{Q} & R' \\ R & 0 \end{pmatrix} \begin{pmatrix} \tilde{\beta} - \beta^* \\ \tilde{\lambda} \end{pmatrix} = \begin{pmatrix} \frac{1}{n} X'e \\ 0 \end{pmatrix},$$

where $\hat{Q} = X'X/n$.

Thus we can explicitly express the estimator as

$$\begin{aligned} \begin{pmatrix} \tilde{\beta} - \beta^* \\ \tilde{\lambda} \end{pmatrix} &= \begin{pmatrix} \hat{Q} & R' \\ R & 0 \end{pmatrix}^{-1} \begin{pmatrix} \frac{1}{n} X'e \\ 0 \end{pmatrix} \\ &= \begin{pmatrix} \hat{Q}^{-1} - \hat{Q}^{-1}R' \left(R\hat{Q}^{-1}R' \right)^{-1} R\hat{Q}^{-1} & \hat{Q}^{-1}R' \left(R\hat{Q}^{-1}R' \right)^{-1} \\ \left(R\hat{Q}^{-1}R' \right)^{-1} R\hat{Q}^{-1} & 0 \end{pmatrix} \begin{pmatrix} \frac{1}{n} X'e \\ 0 \end{pmatrix}. \end{aligned}$$

We conclude that

$$\sqrt{n}\tilde{\lambda} = \left(R\hat{Q}^{-1}R' \right)^{-1} R\hat{Q}^{-1} \frac{1}{\sqrt{n}} X'e \xrightarrow{d} N \left(0, \left(RQ^{-1}R' \right)^{-1} RQ^{-1}\Omega Q^{-1}R' \left(RQ^{-1}R' \right)^{-1} \right).$$

Let $W = \left(RQ^{-1}R' \right)^{-1} RQ^{-1}\Omega Q^{-1}R' \left(RQ^{-1}R' \right)^{-1}$, we have

$$n\tilde{\lambda}'W^{-1}\tilde{\lambda} \xrightarrow{d} \chi_q^2.$$

If homoskedastic, then $W = \sigma^2 \left(RQ^{-1}R' \right)^{-1} RQ^{-1}QQ^{-1}R' \left(RQ^{-1}R' \right)^{-1} = \sigma^2 \left(RQ^{-1}R' \right)^{-1}$. Replace W with the estimated \hat{W} ,

$$\begin{aligned} \frac{n\tilde{\lambda}'R\hat{Q}^{-1}R'\tilde{\lambda}}{\hat{\sigma}^2} &= \frac{1}{n\hat{\sigma}^2} (y - X\tilde{\beta})' X\hat{Q}^{-1}R' \left(R\hat{Q}^{-1}R' \right)^{-1} R\hat{Q}^{-1}X' (y - X\tilde{\beta}) \\ &= \frac{1}{n\hat{\sigma}^2} (y - X\tilde{\beta})' P_{X\hat{Q}^{-1}R'} (y - X\tilde{\beta}). \end{aligned}$$

4.3 Likelihood-Ratio test

For likelihood ratio test, the starting point can be a criterion function $L(\beta) = (y - X\beta)'(y - X\beta)$. It does not have to be the likelihood function.

$$\begin{aligned} L(\tilde{\beta}) - L(\hat{\beta}) &= \frac{\partial L}{\partial \beta}(\hat{\beta}) + \frac{1}{2} (\tilde{\beta} - \hat{\beta})' \frac{\partial^2 L}{\partial \beta \partial \beta}(\hat{\beta}) (\tilde{\beta} - \hat{\beta}) \\ &= 0 + \frac{1}{2} (\tilde{\beta} - \hat{\beta})' \hat{Q} (\tilde{\beta} - \hat{\beta}). \end{aligned}$$

From the derivation of LM test, we have

$$\begin{aligned} \sqrt{n}(\tilde{\beta} - \beta^*) &= \left(\hat{Q}^{-1} - \hat{Q}^{-1}R' \left(R\hat{Q}^{-1}R' \right)^{-1} R\hat{Q}^{-1} \right) \frac{1}{\sqrt{n}} X'e \\ &= \frac{1}{\sqrt{n}} (X'X) X'e - \hat{Q}^{-1}R' \left(R\hat{Q}^{-1}R' \right)^{-1} R\hat{Q}^{-1} \frac{1}{\sqrt{n}} X'e \\ &= \sqrt{n}(\hat{\beta} - \beta^*) - \hat{Q}^{-1}R' \left(R\hat{Q}^{-1}R' \right)^{-1} R\hat{Q}^{-1} \frac{1}{\sqrt{n}} X'e \end{aligned}$$

Therefore

$$\sqrt{n} \left(\tilde{\beta} - \hat{\beta} \right) = -\hat{Q}^{-1} R' \left(R \hat{Q}^{-1} R' \right)^{-1} R \hat{Q}^{-1} \frac{1}{\sqrt{n}} X' e$$

and

$$\begin{aligned} & n \left(\tilde{\beta} - \hat{\beta} \right)' \hat{Q} \left(\tilde{\beta} - \hat{\beta} \right) \\ &= \frac{1}{\sqrt{n}} e' X \hat{Q}^{-1} R' \left(R \hat{Q}^{-1} R' \right)^{-1} R \hat{Q}^{-1} \hat{Q} \hat{Q}^{-1} R' \left(R \hat{Q}^{-1} R' \right)^{-1} R \hat{Q}^{-1} \frac{1}{\sqrt{n}} X' e \\ &= \frac{1}{\sqrt{n}} e' X \hat{Q}^{-1} R' \left(R \hat{Q}^{-1} R' \right)^{-1} R \hat{Q}^{-1} \frac{1}{\sqrt{n}} X' e \end{aligned}$$

In general, it is a quadratic form of normal distributions. If homoskedastic, then

$$\left(R \hat{Q}^{-1} R' \right)^{-1/2} R \hat{Q}^{-1} \frac{1}{\sqrt{n}} X' e$$

has variance

$$\sigma^2 \left(R Q^{-1} R' \right)^{-1/2} R Q^{-1} Q Q^{-1} R' \left(R Q^{-1} R' \right)^{-1/2} = \sigma^2 I_q.$$

We can view the optimization of the log-likelihood as a two-step optimization with the inner step $\sigma = \sigma(\beta)$. By the envelop theorem, when we take derivative with respect to β , we can ignore the indirect effect of $\partial \sigma(\beta) / \partial \beta$.