

4. Heteroscedastic Regression

For the next three problems, consider the following setting.

We measure characteristics of n individuals, sampled randomly from a population. Let (X_i, y_i) be the observed data of the ith individual, where $y_i \in \mathbb{R}$ is the dependent variable and $X_i \in \mathbb{R}^p$ is the vector of p deterministic explanatory variables. Our goal is to estimate the coefficients of $\beta = (\beta_1, \dots, \beta_p)^T$ in the linear regression:

$$y_i = X_i^T eta + \epsilon_i, \qquad i = 1, \dots, n$$

We will consider the case where the model is potentially **heteroscedastic** (i.e. the error terms ϵ_i are **not** i.i.d.).

More specifically, assume that the vector $\boldsymbol{\epsilon} = (\epsilon_1, \dots, \epsilon_n)^T$ is a n-dimensional Gaussian with mean 0 and known **nonsingular** covariance matrix Σ . Denote by \mathbb{X} the matrix in $\mathbb{R}^{n \times p}$ whose rows are $\mathbf{X}_1^T, \dots, \mathbf{X}_n^T$ and by \mathbf{Y} the vector with coordinates y_1, \dots, y_n .

Instead of the usual Least Squares Estimator, instead consider the estimator \hat{eta} that minimizes, over all $eta \in \mathbb{R}^p$,

$$(\mathbf{Y} - \mathbb{X}\beta)^T \Sigma^{-1} (\mathbf{Y} - \mathbb{X}\beta)$$
.

(a) A Generalized Estimator

1/1 point (graded)

Let I_n be the $n \times n$ identity matrix. If $\Sigma = \sigma^2 I_n$ (i.e. homoscedastic ε) for some $\sigma^2 > 0$, then which one of the following statement about $\hat{\beta}$ must be true? Make no assumptions about the rank of X.

- $\hat{\beta}$ has positive entries.
- $m{\hat{eta}}$ is the least squares estimator. \checkmark
- $\hat{oldsymbol{eta}}$ is the unique minimizer of the specified loss.
- igodot The components of $\hat{oldsymbol{eta}}$ are independent.

Solution:

The second choice is the correct answer. Using the fact that $\|v\|^2 = v^T v$, the specified loss simplifies to

$$rac{1}{\sigma^2}(\mathbf{Y} - \mathbb{X}eta)^T I_n \ (\mathbf{Y} - \mathbb{X}eta) = rac{1}{\sigma^2}\|\mathbf{Y} - \mathbb{X}eta)\|^2$$

which is minimized if and only if $\|\mathbf{Y} - \mathbb{X}\boldsymbol{\beta}\|^2$ is minimized, ignoring the term $1/\sigma^2$. Any $\boldsymbol{\beta}$ that attains such a minimum is the least squares estimator.

Now, given the knowledge that $\hat{\beta}$ is the LSE, we can understand the other choices to be false, since they are not properties of the LSE. For example, we would actually want " $\hat{\beta}$ has positive entries" to be false! Imagine the scenario where $\beta_1 = -100$. If our LSE $\hat{\beta}$ somehow had the property that it was always positive, then it would not be a very good estimator to begin with.

1 Answers are displayed within the problem

(b) The Maximum Likelihood Estimator

1.0/1 point (graded)

In this exercise, we will prove that $\hat{\beta}$ is equal to the Maximum Likelihood Estimator, even for general Σ . Recall the form of the n-dimensional Gaussian density from Lecture 10.

Let Σ be an arbitrary $n \times n$ matrix. The maximum likelihood estimator β_{MLE} is the value of β maximizes the log-likelihood function $\ell(\beta) = \ln L(X, Y; \beta)$.

Write down the function ℓ , in terms of X, Y, β , Σ , and n.

(Type **X** for \mathbb{X} , **Y** for \mathbb{Y} , **Sigma** for Σ . Type **trans(X)** for the transpose \mathbb{X}^T , **det(X)** for the determinant $\det \mathbb{X}$, and **X^(-1)** for the inverse \mathbb{X}^{-1} , of a matrix \mathbb{X} .)

$$\ell\left(eta
ight) = \ln L\left(\mathbb{X}, \mathbb{Y}; eta
ight) =$$

-1/2*ln((2*pi)^n*det(Sigma))-1/2*trans(Y-X*beta)*Sigma^-1*(Y-X*beta)

Answer: -(1/2)*trans(Y-X*beta)*(Sigma)^(-1)*(Y-X*beta)-(1/2)*ln((2*pi)^n*det(Sigma))

STANDARD NOTATION

Solution:

Recall that the n-dimensional Gaussian has density

$$f\left(\mathbf{x};\mu,\Sigma
ight) = rac{1}{\sqrt{\left(2\pi
ight)^{n}\mathrm{det}\Sigma}}e^{-rac{1}{2}\left(\mathbf{x}-\mu
ight)^{T}\Sigma^{-1}\left(\mathbf{x}-\mu
ight)}.$$

In our specific setting, the relationship $y_i = X_i^T \beta + \epsilon$ prescribes **Y** to be plugged into **x** in the above density, and $\mu = \mathbb{X}\beta$. The answer is the logarithm of f:

$$-\frac{1}{2}\mathrm{ln}\left(\left(2\pi\right)^{n}\mathrm{det}\Sigma\right)-\frac{1}{2}(\mathbf{Y}-\mathbb{X}\beta)^{T}\Sigma^{-1}\left(\mathbf{Y}-\mathbb{X}\beta\right)$$

Therefore, the MLE maximizes $-\frac{1}{2}(\mathbf{Y}-\mathbb{X}\boldsymbol{\beta})^T\boldsymbol{\Sigma}^{-1}(\mathbf{Y}-\mathbb{X}\boldsymbol{\beta})$ over all possible choices of $\boldsymbol{\beta}$.

To finish the proof, note that we can flip the sign – this turns the maximization into a minimization, so that the MLE minimizes $(\mathbf{Y} - \mathbb{X}\beta)^T \Sigma^{-1} (\mathbf{Y} - \mathbb{X}\beta)$. This matches the provided definition of $\hat{\boldsymbol{\beta}}$.

Submit

You have used 2 of 4 attempts

1 Answers are displayed within the problem

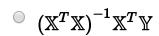
(c)

2/2 points (graded)

Assume that the rank of \mathbb{X} is p. (As a consequence, for any nonsingular $n \times n$ matrix Q, the $p \times p$ matrix product X^TQX is also nonsingular.)

Which of the following is a correct formula for $\hat{\beta}$ in terms of \mathbb{X} , \mathbb{Y} and Σ ?

$$^{\bullet} \ \left(\mathbb{X}^{T} \Sigma^{-1} \mathbb{X}\right)^{-1} \mathbb{X}^{T} \Sigma^{-1} \mathbb{Y} \checkmark$$



$$(\mathbb{X}^T \Sigma^{-1} \mathbb{X})^{-1} \mathbb{Y}^T \Sigma^{-1} \mathbb{X}$$

Using the result from the above, which of the following correctly characterizes the distribution of $\hat{oldsymbol{eta}}$?

- $\circ~\mathcal{N}\left(0,\Sigma
 ight)$
- $^{\circ}~~\mathcal{N}\left(0,\left(\mathbb{X}^{T}\Sigma^{-1}\mathbb{X}
 ight)^{-1}
 ight)$
- $^{\circ}~~\mathcal{N}\left(0,\left(\mathbb{X}^{T}\mathbb{X}
 ight)^{-1}
 ight)$
- $^{\odot} \mathcal{N}(\beta, (\mathbb{X}^T \Sigma^{-1} \mathbb{X})^{-1}) \checkmark$
- $^{\circ}$ $\mathcal{N}\left(eta,\left(\mathbb{X}^{T}\mathbb{X}
 ight)^{-1}
 ight)$

Solution:

For the first question: Start by expanding the product.

$$\begin{aligned} \left(\mathbb{Y} - \mathbb{X}\beta\right)^T \Sigma^{-1} \left(\mathbb{Y} - \mathbb{X}\beta\right) &= \mathbb{Y}^T \Sigma^{-1} \mathbb{Y} - \mathbb{Y}^T \Sigma^{-1} \mathbb{X}\beta - \beta^T \mathbb{X}^T \Sigma^{-1} \mathbb{Y} + \beta^T \mathbb{X}^T \Sigma^{-1} \mathbb{X}\beta \\ &= \mathbb{Y}^T \Sigma^{-1} \mathbb{Y} - \frac{2\beta^T \mathbb{X}^T \Sigma^{-1} \mathbb{Y}}{2\beta^T \mathbb{X}^T \Sigma^{-1} \mathbb{Y}} + \beta^T \mathbb{X}^T \Sigma^{-1} \mathbb{X}\beta. \end{aligned}$$

The second equality uses the fact that $(\mathbb{Y}^T \Sigma^{-1} \mathbb{X} \beta)^T = \mathbb{Y}^T \Sigma^{-1} \mathbb{X} \beta$, since the left hand side is a scalar, and Σ is symmetric. Next, take the gradient with respect to β , which evaluates to $-2\mathbb{X}^T \Sigma^{-1} \mathbb{Y} + 2\mathbb{X}^T \Sigma^{-1} \mathbb{X} \beta$. Set this equal to zero and solve for β to obtain the estimator $\hat{\beta}$:

$$\hat{eta} = \left(\mathbb{X}^T \Sigma^{-1} \mathbb{X} \right)^{-1} \mathbb{X}^T \Sigma^{-1} \mathbb{Y}.$$

Thus, the correct answer is the **first choice**.

In the initial step, one needs to be a bit careful. If instead we wrote $2\mathbb{Y}^T\Sigma^{-1}\mathbb{X}\beta$ for the middle term in the second equality, we might be tempted to choose the third choice instead. By computing the gradient more carefully, the answer does consistently turn out to be the first choice. An easy way to reason about this, without explicitly writing out the calculus, is to remember that the dimensions must match up. For instance, the matrix $(\mathbb{X}^T\Sigma^{-1}\mathbb{X})^{-1}$ is $p \times p$, yet \mathbb{Y}^T is a $1 \times p$ matrix - therefore, the third choice is not a valid expression

For the second question, start with $\hat{\beta} = (\mathbb{X}^T \Sigma^{-1} \mathbb{X})^{-1} \mathbb{X}^T \Sigma^{-1} \mathbb{Y}$ and substitute $\mathbb{Y} = \mathbb{X}\beta + \epsilon$.

$$\begin{split} \hat{\beta} &= \left(\mathbb{X}^T \Sigma^{-1} \mathbb{X}\right)^{-1} \mathbb{X}^T \Sigma^{-1} \left(\mathbb{X} \beta + \epsilon\right) \\ &= \left(\mathbb{X}^T \Sigma^{-1} \mathbb{X}\right)^{-1} \mathbb{X}^T \Sigma^{-1} \mathbb{X} \beta + \left(\mathbb{X}^T \Sigma^{-1} \mathbb{X}\right)^{-1} \mathbb{X}^T \Sigma^{-1} \epsilon \\ &= \beta + \left(\left(\mathbb{X}^T \Sigma^{-1} \mathbb{X}\right)^{-1} \mathbb{X}^T \Sigma^{-1}\right) \epsilon. \end{split}$$

Next, we use the rule of transforming a gaussian $\mathbf{x}\sim\mathcal{N}\left(0,\Sigma\right)$ by a matrix M, which states that $M\mathbf{x}\sim\mathcal{N}\left(0,M\Sigma M^{T}\right)$. In particular,

$$(\hat{eta} - eta) \sim \mathcal{N}\left(0, \left[\left(\mathbb{X}^T \Sigma^{-1} \mathbb{X}
ight)^{-1} \mathbb{X}^T \Sigma^{-1}
ight] \Sigma \left[\left(\mathbb{X}^T \Sigma^{-1} \mathbb{X}
ight)^{-1} \mathbb{X}^T \Sigma^{-1}
ight]^T
ight)$$

Submit	You have used 1 of 3 attempts
) Answei	s are displayed within the problem
A nswei	s are displayed within the problem

© All Rights Reserved